

# *Software Manual*

## A mathematical model for universal semantics

Weinan E<sup>1,2\*</sup>, Yajun Zhou<sup>2\*</sup>

<sup>1</sup>Department of Mathematics & Program in Applied and Computational Mathematics,  
Princeton University, Princeton, NJ 08544, USA

<sup>2</sup>Beijing Institute of Big Data Research, Beijing 100871, P. R. China

\*Corresponding authors. E-mail: weinan@math.princeton.edu (W.E), yajun.zhou.1982@pku.edu.cn (Y.Z.)

March 2020

## Contents

<b>I Scope and Limitations</b>	<b>1</b>
1 What our software does	
1.1 Markov analysis of word patterns and topic extraction . . . . .	1
1.2 Semantic fingerprints from invariant Markov spectra . . . . .	2
1.3 Automated word translation based on semantic fingerprints . . . . .	5
1.4 Automated question answering based on Markov semantics . . . . .	11
2 What our software does not	
2.1 Non-universal statistics on short time scales . . . . .	11
2.2 Word disambiguation and stylistic adaptability . . . . .	12
2.3 Causal inference and deductive reasoning . . . . .	14
<b>II Protocols for text cleansing</b>	<b>15</b>
3 String manipulations and text normalizations	
3.1 Notational conventions . . . . .	16
3.2 Normalizations of Wikipedia pages and WikiQA dataset . . . . .	18
3.3 Normalizations of ebooks for text mining . . . . .	22
<b>III Protocols for word clustering</b>	<b>24</b>
4 Approximate word clustering in English	
4.1 Capitalization and stop words . . . . .	24
4.2 Modified Porter stemming algorithm for English . . . . .	25
4.2.1 Effective spelling and essential root . . . . .	26
4.2.2 Admissible mutation and approximate clustering . . . . .	30
5 Approximate word clustering in selected Germanic languages	
5.1 Modified Porter stemming algorithm for Danish . . . . .	60
5.1.1 Effective spelling and essential root . . . . .	61
5.1.2 Admissible mutation and approximate clustering . . . . .	62
5.1.3 Heuristic detection of compounds . . . . .	66

5.2	Modified Porter stemming algorithm for German . . . . .	70
5.2.1	Effective spelling and essential root . . . . .	71
5.2.2	Admissible mutation and approximate clustering . . . . .	73
5.2.3	Heuristic detection of compounds . . . . .	77
5.3	Modified Porter stemming algorithm for Dutch . . . . .	81
5.3.1	Effective spelling and essential root . . . . .	81
5.3.2	Admissible mutation and approximate clustering . . . . .	83
5.3.3	Heuristic detection of compounds . . . . .	86
6	Approximate word clustering in selected Romance languages	
6.1	Modified Porter stemming algorithm for Spanish . . . . .	89
6.1.1	Effective spelling and essential root . . . . .	90
6.1.2	Admissible mutation and approximate clustering . . . . .	92
6.2	Modified Porter stemming algorithm for French . . . . .	103
6.2.1	Effective spelling and essential root . . . . .	104
6.2.2	Admissible mutation and approximate clustering . . . . .	106
6.3	Modified Porter stemming algorithm for Latin . . . . .	120
6.3.1	Effective spelling and essential root . . . . .	122
6.3.2	Admissible mutation and approximate clustering . . . . .	125
7	Approximate word clustering in selected Slavic languages	
7.1	Modified Porter stemming algorithm for Polish . . . . .	148
7.1.1	Effective spelling and essential root . . . . .	149
7.1.2	Admissible mutation and approximate clustering . . . . .	152
7.2	Modified Porter stemming algorithm for Russian . . . . .	164
7.2.1	Effective spelling and essential root . . . . .	166
7.2.2	Admissible mutation and approximate clustering . . . . .	169
8	Approximate word clustering in selected Uralic languages	
8.1	Modified Porter stemming algorithm for Finnish . . . . .	187
8.1.1	Effective spelling and essential root . . . . .	190
8.1.2	Admissible mutation and approximate clustering . . . . .	194
8.1.3	Heuristic detection of compounds . . . . .	210
8.2	Modified Porter stemming algorithm for Hungarian . . . . .	216
8.2.1	Effective spelling and essential root . . . . .	218
8.2.2	Admissible mutation and approximate clustering . . . . .	221
8.2.3	Heuristic detection of compounds . . . . .	227
9	Approximate word clustering in various agglutinative languages	
9.1	Modified Porter stemming algorithm for Basque . . . . .	232
9.1.1	Effective spelling and essential root . . . . .	234
9.1.2	Approximate clustering . . . . .	236
9.2	Modified Porter stemming algorithm for Korean . . . . .	241
9.2.1	Stop words and transliterations . . . . .	242
9.2.2	Effective spelling and essential root . . . . .	245
9.2.3	Approximate clustering . . . . .	250
9.3	Modified Porter stemming algorithm for Turkish . . . . .	255
9.3.1	Effective spelling and essential root . . . . .	256
9.3.2	Approximate clustering . . . . .	257

This Software Manual accompanies our article entitled “A mathematical model for universal semantics” [1], supplying linguistic and algorithmic details. Mathematical proofs behind our model are elaborated in [1, Appendices A and B], so they will not be repeated here.

In Part I, we describe the scope and limitations of our current computational methods for automated word translation and text comprehension. Numerical tests on corpora are also briefly summarized in tables and figures.

In Part II, we define some terminologies and notations in linguistics, in the context of our current software implementation. We also list all the procedures required to cleanse our text corpora.

In Part III, we present detailed algorithms for treating 14 languages out of 5 representative language families, including 12 European languages (Danish, German, Dutch, Spanish, French, Latin, Polish and Russian from the Indo-European language family; Finnish and Hungarian from the Uralic language family; Basque from the Vasconic language family) and 2 Asian languages (Korean from the Koreanic language family; Turkish from the Turkic language family). We illustrate these algorithms with automated word translations from English masterpieces to other languages, and with question answering on the WikiQA dataset.

## Part I

# Scope and Limitations

## 1 What our software does

### 1.1 Markov analysis of word patterns and topic extraction

As in [1, §1], we are interested in word patterns, each of which is a collection of morphologically related content words. We refer our readers to Part III for the working definition of content words (which varies from language to language), and how we recognize and group content words that have related morphologies.

In our numerical studies of text documents, we count waiting times between two word patterns  $W_i$  and  $W_j$  by the *effective fragment length*  $L_{ij}$  (defined in [1, Fig. 1]). As we have mentioned in the caption to [1, Fig. 1], a text fragment qualifies as an *long-range transition* from  $W_i$  and  $W_j$ , if

- it is flanked by a word in  $W_i$  to the left and a word in  $W_j$  to the right;
- it does not mention any word in the set  $W_i$ ;
- it is longer than the longest word in the set  $W_i \cup W_j$ .

We always ignore close encounters of words (that is, those failing the last of the three criteria above), while collecting data for effective fragment lengths over long-range transitions.

With  $n_{ij}$  samples of effective fragment lengths  $L_{ij}$  in long-range transitions from  $W_i$  to  $W_j$ ,<sup>1</sup> we model a text by a Markov matrix  $\mathbf{P} = (p_{ij})_{1 \leq i,j \leq N}$ , whose entries are given by [1, (5)]

$$p_{ij} := \frac{n_{ij} e^{-\langle \log L_{ij} \rangle}}{\sum_{k=1}^N n_{ik} e^{-\langle \log L_{ik} \rangle}}. \quad (1.1)$$

Numerical evidence suggests that this empirical Markov matrix  $\mathbf{P} = (p_{ij})_{1 \leq i,j \leq N}$  is a fair approximation to an *ergodic* matrix  $\mathbf{P}^* = (p_{ij}^*)_{1 \leq i,j \leq N}$ , whose equilibrium state  $\boldsymbol{\pi}^* = (\pi_1^*, \dots, \pi_N^*)$  matches word count statistics [1, Fig. 4a], and honors the *detailed balance* condition  $\pi_i^* p_{ij}^* = \pi_j^* p_{ji}^*$  [1, Fig. 4b].

For an ergodic Markov chain satisfying the detailed balance condition, one can show that the probability distribution of recurrence times  $\tau = L_{ii}$  goes like [1, Appendix B.3]

$$\mathbb{P}(\tau > t) \sim \sum_m c_m e^{-k_m t}, \quad \left( \text{where } c_m, k_m > 0, \text{ and } \sum_m c_m = 1 \right), \quad (1.2)$$

a functional form that frequently crops up in dynamic studies of biological macromolecules [7]. This weighted mixture of multiple exponential decay laws imposes an inequality constraint on the recurrence time  $\tau$  [1, (3)]:

$$\langle \log \tau \rangle - \log \langle \tau \rangle + \gamma_0 = \sum_m c_m \log \frac{1}{k_m} - \log \sum_m \frac{c_m}{k_m} \leq 0, \quad (1.3)$$

---

<sup>1</sup>It is worth noting that we construct an empirical Markov matrix by an *in situ* analysis of a text, without digesting a document (or small parts of it) as a scrambled bag of words, a procedure implemented in conventional algorithms [2, 3, 4, 5, 6]. The notation  $p_{ij}$  in the semantic model of Turney–Pantel [4] is unrelated to ours.

where  $\gamma_0 := \lim_{n \rightarrow \infty} (-\log n + \sum_{m=1}^n \frac{1}{m})$  is the Euler–Mascheroni constant. This inequality explains the systematic data trend in [1, Fig. 2e], where most data points reside on one side of the line of Poissonian banality  $\langle \log \tau \rangle - \log(\tau) + \gamma_0 = 0$  (blue line in [1, Fig. 2e]).

In [1, Appendix A], we have shown that the statistic

$$\delta_i := \log\langle L_{ii} \rangle - \langle \log L_{ii} \rangle - \gamma_0 + \frac{1}{2n_{ii}} \quad (1.4)$$

satisfies

$$|\delta_i| < \frac{2}{\sqrt{n_{ii}}} \sqrt{\frac{\pi^2}{6} - 1 - \frac{1}{2n_{ii}}} \quad (1.5)$$

with probability 95%, when there are  $n_{ii}$  independent samples of exponentially distributed (hence banal)  $L_{ii}$ . In our current work, we automatically extract topics (Figs. S3–S16 in Part III), by identifying all the word patterns that violate the inequality above. Being consistent with the weighted exponential mixture model in (1.2), almost all of our automatically identified topics in Figs. S3–S16 satisfy

$$\delta_i \geq \frac{2}{\sqrt{n_{ii}}} \sqrt{\frac{\pi^2}{6} - 1 - \frac{1}{2n_{ii}}} > 0. \quad (1.6)$$

If a concept (representable by a word pattern) qualifies as a topical (resp. non-topical) pattern in a certain document, according to the algorithm above, then this concept usually remains topical (resp. non-topical) if we examine the long-range statistic  $\delta_i$  in a parallel document written in a different language [1, Fig. 3b,b']. In other words, topicality/non-topicality is an invariant (*i.e.* a language-independent property) under translations.

Generalizing this further, we hypothesize that temporal structures are nearly universal across different languages, on all the scales longer than next-to-nearest neighboring words. In fact, later in §1.2, we will demonstrate that the statistics of effective fragment lengths  $L_{ij}$  almost represent a complete set of semantic invariants, quantifying the semantic similarity between  $W_i$  and  $W_j$  in a language-independent way.

## 1.2 Semantic fingerprints from invariant Markov spectra

What gives meanings to things? If we threw this ontological question to paleolithic cavemen or preschool children, we would probably expect tangible answers that define objects by physical and physiological stimuli. Putatively speaking, cavemen would conceptualize “fire” by associating it to “brightness” (visual impression), “warmth” (feeling at a distance) and “pain” (reaction upon touch); children would recognize an “apple” as something that is “round” (in shape), “red” (in color) and “sweet” (in taste). While this reductionist view of semantics does not generalize well to abstract concepts (“love”, “pride”, ...), fictional characters (“Elizabeth Bennet”, “Harry Potter”, ...) or non-existent places (“Meryton”, “Hogwarts”, ...), we still find it an attractive idea to generate semantic tags by association—perhaps a universal way to clarify the meaning of a word is to link it to a small set of other words (even if they do not represent bodily stimuli).<sup>2</sup>

Such an ontological picture motivates us to define the meaning of words in terms of a language-independent numerical fingerprint, using Markov matrices that characterize association and connectivity. As we will soon see, certain invariance properties of Markov matrices account for the translatability of concepts across a wide variety of languages. Put differently, in our Markov language model, the semantic context of an individual word pattern is (up to numerical errors) independent of the language in which it is expressed—the mental states underlying a text thus remain (approximately) invariant under translation.

Through numerical experiments ([1, Fig. 2e] and Fig. S1a,a',a''), we find that the spectrum of a Markov matrix (collection of eigenvalues, counting multiplicity) is language-independent.

In [1, §2.2.3], we have designed a thought experiment (estimation of brainstorming rates from words in language A to words in language B) involving two monolingual human subjects, which has led us to an equation [1, (6)]

$$\mathbf{P}_A \mathbf{T}_{A \rightarrow B} = \mathbf{T}_{A \rightarrow B} \mathbf{P}_B. \quad (1.7)$$

Here,  $\mathbf{P}_A$  and  $\mathbf{P}_B$  are, respectively, the Markov matrices characterizing the mental language of Alice (knowing only language A) and Bob (knowing only language B), while  $\mathbf{T}_{A \rightarrow B}$  is the common dictionary matrix shared by both people. If we have an invertible Markov matrix  $\mathbf{T}_{A \rightarrow B}$  (when translations are lossless), then  $\mathbf{P}_A$  and  $\mathbf{P}_B$  share the same characteristic polynomial  $\det(\lambda \mathbf{I} - \mathbf{P}_A) = \det(\lambda \mathbf{I} - \mathbf{P}_B)$ , hence the same spectrum. We note that the dictionary matrix  $\mathbf{T}_{A \rightarrow B}$  here does not need to be

<sup>2</sup>It is worth mentioning that Jane Austen has written several novels, besides *Pride and Prejudice*, with themes of love and marriage, despite her very lack of personal experience in these matters—she was able to create vivid characters out of her pure imagination, while tagging them with appropriate contexts.

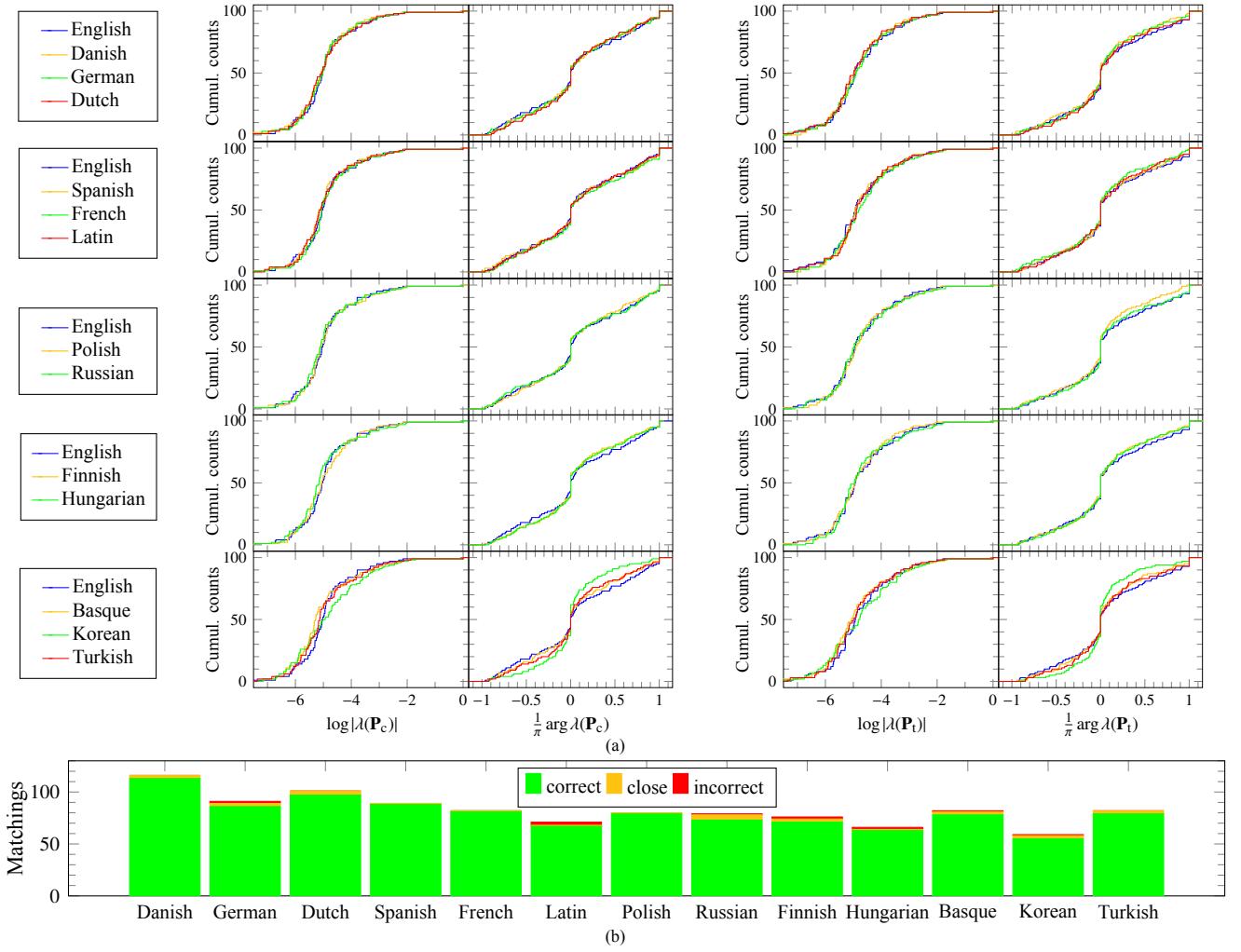


Fig. S1. Language-independence of Markov spectra. (a) Approximate spectral invariance in the Markov matrices  $\mathbf{P}_c$  (resp.  $\mathbf{P}_t$ ) for top 100 content (resp. topic) word clusters from the English original and 13 different translations of *Pride and Prejudice*. (See §9.2 for difficulties in analyzing Korean texts.) (b) Evaluations of automated semantic alignments ([1, Fig. 6] and Figs. S4–S7, S9–S16) between topics in the English version of *Pride and Prejudice* and those in the 13 different translations.

a permutation matrix<sup>3</sup> to qualify for spectral invariance. An invertible  $\mathbf{T}_{A \rightarrow B}$  may assume non-permutation forms (which partially accommodates to polysemy of words in realistic languages), such as  $\begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ 0 & 1 \end{pmatrix}$ . Furthermore, the inverse matrix  $\mathbf{T}_{A \rightarrow B}^{-1}$  does not necessarily coincide with the dictionary matrix  $\mathbf{T}_{B \rightarrow A}$ . This is because  $\mathbf{P}_A(\mathbf{A}\mathbf{T}_{A \rightarrow B}) = (\mathbf{A}\mathbf{T}_{A \rightarrow B})\mathbf{P}_B$  and  $\mathbf{P}_A(\mathbf{T}_{A \rightarrow B}\mathbf{B}) = (\mathbf{T}_{A \rightarrow B}\mathbf{B})\mathbf{P}_B$  hold for any matrices  $\mathbf{A}$  and  $\mathbf{B}$  satisfying  $\mathbf{A}\mathbf{P}_A = \mathbf{P}_A\mathbf{A}$  and  $\mathbf{B}\mathbf{P}_B = \mathbf{P}_B\mathbf{B}$ .

In [1, Fig. 4c], we have demonstrated approximate spectral invariance, by localizing Markov matrices on the top 100 content word patterns, in the English, French, Russian and Finnish versions of *Pride and Prejudice*. Such invariance properties remain unscathed when we consider translations of the same novel into other languages, and when we localize on the top 100 topical patterns instead (Fig. S1a). The numerical agreement is understandable from the language-independence of semantic content. In fact, we have pushed this spectral invariance down to an even more local scale in [1, Fig. 5b]: after translations, there are only small perturbations in the dominant eigenvalues for Markov matrices localized on *semantic cliques* surrounding specific concepts.

Here, to define a semantic clique  $\mathcal{S}_i$  surrounding a specific word pattern  $W_i$ , we need to specify numerical semantic fields through an affinity score  $\alpha_{ij}(\ell)$ . If we have a Markov process on a semantic web that honors detailed balance, then for each fixed word pattern  $W_i$ , we can determine [8] the distribution of hitting times  $L_{ij}, 1 \leq j \leq N$  from that of the return times  $L_{ii}$ . Further assuming that the distribution of  $\langle \log L_{ij} \rangle$  is nearly Gaussian, we can use the following approximation [1, (7)]:

$$\mathbb{P}(\langle \log L_{ij} \rangle > \ell) \approx \alpha_{ij}(\ell) := \sqrt{\frac{n_{ij}}{2\pi\beta_i}} \int_{\ell}^{\infty} e^{-\frac{n_{ij}(x-\ell_i)^2}{2\beta_i}} dx, \quad (1.8)$$

<sup>3</sup>A permutation matrix is a Markov matrix whose entries are either 0 or 1. When the dictionary matrix  $\mathbf{T}_{A \rightarrow B}$  is a permutation matrix, translation becomes a bijection of words between two languages. This picture, of course, is much too restrictive.

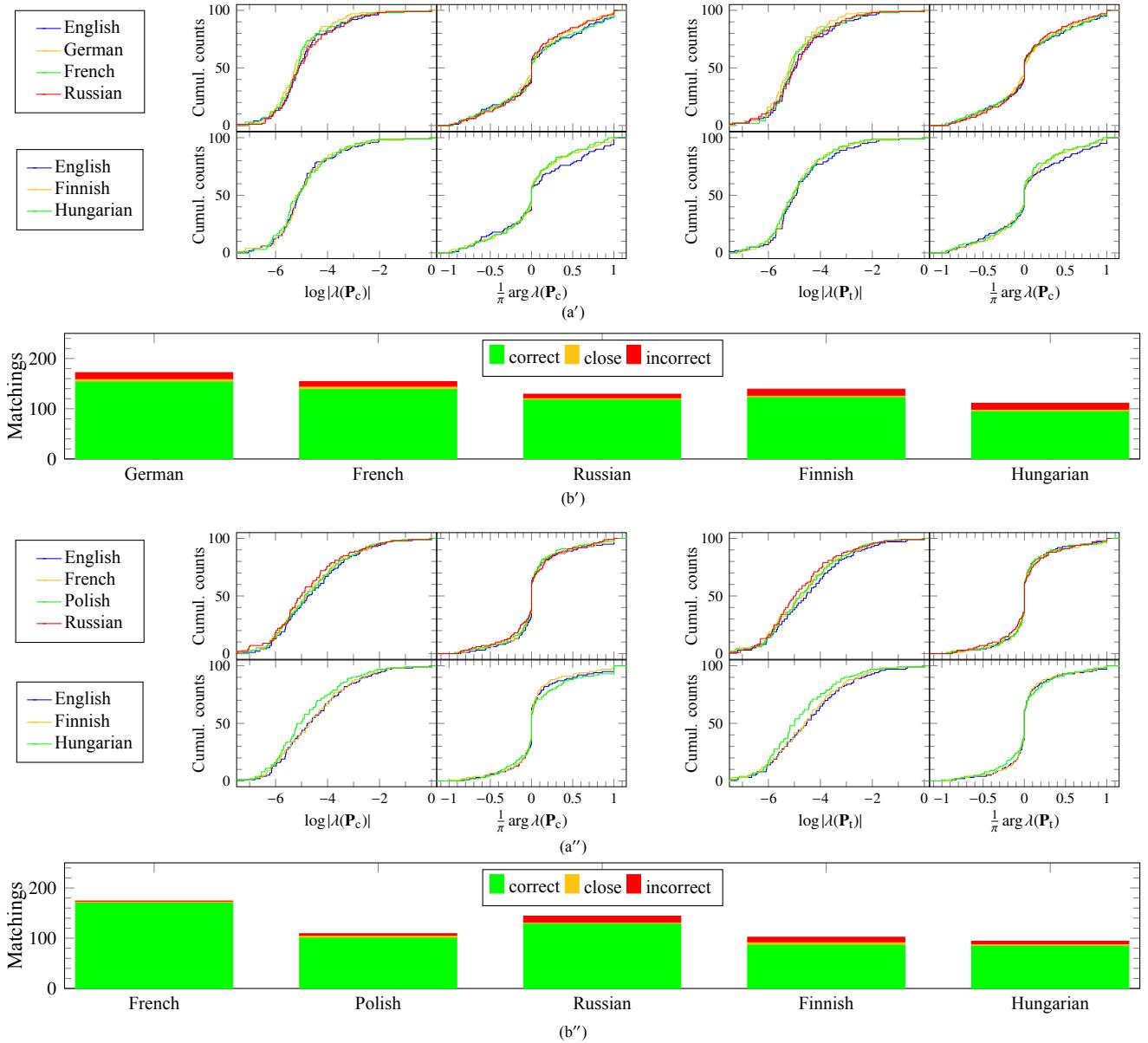


Fig. S1. Language-independence of Markov spectra. (*Continued*) Similar service on multiple versions of (a')–(b') *Jane Eyre* and (a'')–(b'') *Origin of Species*.

where [1, (8) and Theorem 4 in Appendix B.4]

$$\ell_i := \frac{\langle L_{ii} \log L_{ii} \rangle}{\langle L_{ii} \rangle} - 1, \quad \beta_i := \frac{\langle L_{ii}(\ell_i - \log L_{ii})^2 \rangle}{\langle L_{ii} \rangle}. \quad (1.9)$$

Our semantic clique  $\mathcal{S}_i$  [1, Fig. 5a,b] admits  $W_j$  as its member if and only if

$$\min\{\alpha_{ij}(\langle \log L_{ij} \rangle), \alpha_{ji}(\langle \log L_{ji} \rangle)\} > \frac{1}{\sqrt{2\pi}} \int_{-\infty}^1 e^{-\frac{x^2}{2}} dx. \quad (1.10)$$

Here, we give some caveats on the spectral invariance for localized Markov models of semantic cliques, and further clarify the meaning of “dominant eigenvalues”. The criterion (1.10) for membership in semantic cliques and the 95% confidence interval for topicality [see (1.5)] both involve arbitrary cut-offs. Whenever there are close calls, our algorithm may include certain topics in a particular semantic clique  $\mathcal{S}_i^A$  in language A, but exclude them from a parallel clique  $\mathcal{S}_i^B$  in language B. Thus, the size of semantic cliques might not be stable across languages (see dashed lines in [1, Fig. 5b]), due to the presence of marginal topics. To combat this instability, we have proposed to use the entropy production rate<sup>4</sup> to compress dimensions (see

<sup>4</sup>Heuristically speaking, the entropy production rate  $\eta(\mathbf{P}) = -\sum_{i,j} \pi_i p_{ij} \log p_{ij}$  [9, (4.27)] of a Markov matrix  $\mathbf{P}$  represents the weighted average (assigning

solid lines in [1, Fig. 5b]). For an  $N \times N$  Markov matrix  $\mathbf{P}$  with strictly positive entries, the entropy production rate  $\eta(\mathbf{P})$  satisfies a sharp inequality<sup>5</sup>  $\eta(\mathbf{P}) \leq \log N$  [10, Theorem 14.1], so  $\lfloor e^{\eta(\mathbf{P})} \rfloor$  represents the effective dimension of our localized Markov matrix  $\mathbf{P}$  for a semantic clique. This effective dimension is less susceptible to marginal topics than the number of members in a numerically constructed semantic clique (= the size of matrix  $\mathbf{P}$ ), hence we use it to decide how many dominant eigenvalues to keep before vectorizing a topic from the recurrence eigenvalues [1, Fig. 5b] in its surrounding semantic clique. The vector that arranges these recurrence eigenvalues (corresponding to decay modes in (1.2)) in descending order is a *language-independent semantic fingerprint* for a topic.

### 1.3 Automated word translation based on semantic fingerprints

Our machine translation experiments solve bipartite matching problems involving the Ružička similarities [11, 12]

$$s_R(\mathbf{a}, \mathbf{b}) := \frac{\|\mathbf{a} \wedge \mathbf{b}\|_1}{\|\mathbf{a} \vee \mathbf{b}\|_1} \equiv \frac{\sum_i \min(a_i, b_i)}{\sum_i \max(a_i, b_i)} \quad (1.11)$$

between vectorized topics.<sup>6</sup> The Ružička dissimilarity function  $d_R(\mathbf{a}, \mathbf{b}) = 1 - s_R(\mathbf{a}, \mathbf{b})$  measures the distance between two non-zero vectors with non-negative components, and satisfies the three axioms for a distance metric:

1. (Discernibility)  $d_R(\mathbf{a}, \mathbf{b}) = 0 \iff \mathbf{a} = \mathbf{b}$ ;
2. (Symmetry)  $d_R(\mathbf{a}, \mathbf{b}) = d_R(\mathbf{b}, \mathbf{a})$ ;
3. (Subadditivity)  $d_R(\mathbf{a}, \mathbf{c}) \leq d_R(\mathbf{a}, \mathbf{b}) + d_R(\mathbf{b}, \mathbf{c})$ .

To verify subadditivity, one can use Gilbert's bounds [13]

$$1 - \frac{\|\mathbf{v}_1 \wedge \mathbf{v}_2 \wedge \mathbf{v}_3\|_1}{\|\mathbf{v}_1 \vee \mathbf{v}_2 \vee \mathbf{v}_3\|_1} \geq d_R(\mathbf{v}_i, \mathbf{v}_j) = 1 - \frac{\|\mathbf{v}_i \wedge \mathbf{v}_j\|_1}{\|\mathbf{v}_i \vee \mathbf{v}_j\|_1} \geq \frac{\|\mathbf{v}_i \vee \mathbf{v}_j\|_1 - \|\mathbf{v}_i \wedge \mathbf{v}_j\|_1}{\|\mathbf{v}_1 \vee \mathbf{v}_2 \vee \mathbf{v}_3\|_1} \quad (1.12)$$

for  $i, j \in \{1, 2, 3\}$ , together with the fact that  $\|\mathbf{a} \vee \mathbf{b} \vee \mathbf{c}\|_1 - \|\mathbf{a} \wedge \mathbf{b} \wedge \mathbf{c}\|_1 \leq (\|\mathbf{a} \vee \mathbf{b}\|_1 - \|\mathbf{a} \wedge \mathbf{b}\|_1) + (\|\mathbf{b} \vee \mathbf{c}\|_1 - \|\mathbf{b} \wedge \mathbf{c}\|_1)$ . Since the Ružička dissimilarity satisfies the axioms for a distance metric, our bipartite matching model for machine translation is a discrete version of the optimal transport in variational calculus.

During the ballpark screening stage of our machine translation algorithm, we divide our target document into  $K$  chapters and vectorize each topic pattern  $W_i$  by its chapter-wise counts  $\mathbf{b}_i \in \mathbb{R}^K$ . To pass a ballpark screening test, a pair of topics  $W_i^A$  and  $W_j^B$  from parallel versions (in languages A and B) of the same document have to satisfy

$$d_R(\mathbf{b}_i^A, \mathbf{b}_j^B) \leq \varepsilon \sqrt{K} \quad \text{and} \quad d_R(\mathbf{b}_i^A, \mathbf{b}_j^B) \leq \sqrt{\frac{\|\mathbf{b}_i^A \wedge \mathbf{b}_j^B\|_0}{\|\mathbf{b}_i^A \vee \mathbf{b}_j^B\|_1}}. \quad (1.13)$$

Here, the coefficient  $\varepsilon = 0.07$  in the first inequality above is chosen empirically (see control experiments in Figs. S5, S8, S10, S11, S12, S13), so that the right-hand side extrapolates to a reasonable error margin<sup>7</sup> for total word counts (when  $K = 1$ ). The second inequality in (1.13) draws upon a Poisson approximation to word count fluctuations, with  $\|\mathbf{b}_i^A \wedge \mathbf{b}_j^B\|_0$  being the number of corresponding chapters that contain both  $W_i^A$  and  $W_j^B$ . Through numerical experiments, we find that the screening criteria in (1.13) lead to acceptable recall and precision of machine translation later afterwards.

The vector  $\mathbf{b}_i \in \mathbb{R}^K$  captures only the temporal behaviour of  $W_i$  on the chapter scale. To paint a fuller kinetic picture of  $W_i$  above the discourse level, we vectorize it by  $\mathbf{v}_i$ , the list of dominant recurrence eigenvalues in its semantic clique, and use  $s_R(\mathbf{v}_i^A, \mathbf{v}_j^B)$  in our final similarity score for bipartite matching (see Fig. S1b,b',b'', Tables S2–S5 for summaries of performance; see [1, Fig. 6] and Figs. S4–S16 for detailed bipartite matching of semantic fingerprints  $\mathbf{v}_i$ ). In the current work, we do impose an additional censorship  $s_R(\mathbf{v}_i^A, \mathbf{v}_j^B) \geq 0.7$  before throwing similarity scores to the Hungarian algorithm for bipartite matching. This threshold value 0.7 is somewhat arbitrary. Relaxing this arbitrary censorship (*i.e.* reducing the threshold to 0) seems to have minimal consequences for translations between English and its close relative (see a control experiment for English–French in Table S6 and Fig. S8c), but seems to boost the recall rate significantly for translations between English and a non-Indo-European language (see a control experiment for English–Korean in Table S6 and Fig. S15c). Perhaps more work is needed to understand whether/how the threshold here adapts to the specific language pairs.

<sup>5</sup>probability mass  $\pi_i$  to the  $i$ th Markov state) of Boltzmann's partition entropies  $-\sum_j p_{ij} \log p_{ij}$  [10, §8.2]. For a formal discussion of entropy production and information compressibility in Markov chains, see [9, Chaps. 4–5] and [10, Chaps. 3, 8 and 14].

<sup>6</sup>The equality is attained when  $\mathbf{P} = (p_{ij} = \frac{1}{N})_{1 \leq i, j \leq N}$ . This is the case of the least compressible Markov chain, whose maximal effective dimension is equal to  $N$ .

<sup>7</sup>Here as in [1, (4)],  $\mathbf{a} \wedge \mathbf{b} = (\min\{a_1, b_1\}, \dots, \min\{a_N, b_N\})$  and  $\mathbf{a} \vee \mathbf{b} = (\max\{a_1, b_1\}, \dots, \max\{a_N, b_N\})$  for  $\mathbf{a} = (a_1, \dots, a_N)$ ,  $\mathbf{b} = (b_1, \dots, b_N)$ ;  $\|\mathbf{u}\|_1 = \sum_{i=1}^N |u_i|$  for  $\mathbf{u} = (u_1, \dots, u_N)$ .

<sup>7</sup>Possible error sources are imperfections in translation and word clustering. See discussions below.

In Figs. S3–S16, we run our experiments on multilingual corpora (Table S1) involving 9 members (Danish, Dutch, English, French, German, Latin, Polish, Russian and Spanish) in the Indo-European language family, 2 members (Finnish and Hungarian) in the Uralic language family, and 3 more agglutinative languages (Basque, Korean, Turkish) that are neither Indo-European nor Uralic, which exhibit diverse word orders<sup>8</sup> for subject (S), verb (V) and object (O):

SVO: Danish, English, French, Spanish;<sup>9</sup>

SV<sub>1</sub>O(V<sub>2</sub>): Dutch, German;<sup>10</sup>

SOV: Basque, Latin, Korean, Turkish;<sup>11</sup>

Free word order: Finnish, Hungarian, Polish, Russian.<sup>12</sup>

The partial success of these experiments demonstrates that our Markov model (of topics on the long-range length scale) is translatable, regardless of a language's genealogical affiliation and syntactical structure.

**Table S1. Provenances of multilingual texts used in our experiments and discussions**

Language	Author/Adaptor	Title	Publisher	Electronic availability
<i>Indo-European</i>				
English	Jane Austen	<i>Pride and Prejudice</i>	T. Egerton (Whitehall, UK, 1813)	ExampleData in <i>Mathematica</i> (essentially identical to the 1813 version, with modernized spellings) <a href="https://www.imsdb.com/scripts/Pride-and-Prejudice.html">https://www.imsdb.com/scripts/Pride-and-Prejudice.html</a>
Danish	Deborah Moggach	<i>Pride and Prejudice</i>	Internet Movie Scripts Database (2005)	<a href="https://www.imsdb.com/scripts/Pride-and-Prejudice.html">https://www.imsdb.com/scripts/Pride-and-Prejudice.html</a>
German	Vibeke Houstrup	<i>Stolz und Vorurteil</i>	Lindhardt og Ringhof Forlag (Copenhagen, Denmark, 2015)	<a href="https://play.google.com/store/books/details?id=dd0-CgAAQBAJ&amp;rid=book-dd0-CgAAQBAJ">https://play.google.com/store/books/details?id=dd0-CgAAQBAJ&amp;rid=book-dd0-CgAAQBAJ</a> (Google Play)
	Werner Beyer	<i>Stolz und Vorurteil</i>	Fischer Klassik Plus (Frankfurt am Main, Germany, 2011)	<a href="https://www.amazon.de/Stolz-Vorurteil-Roman-Fischer-Klassik-ebook/dp/B006GH31IC/">https://www.amazon.de/Stolz-Vorurteil-Roman-Fischer-Klassik-ebook/dp/B006GH31IC/</a> (Amazon)
	Christian Grawe	<i>Stolz und Vorurteil</i>	Philipp Reclam jun. GmbH & Co. KG (Stuttgart, Germany, 2015)	<a href="https://www.amazon.de/Stolz-Vorurteil-Roman-Reclam-Taschenbuch-ebook/dp/B01BHWFX6/">https://www.amazon.de/Stolz-Vorurteil-Roman-Reclam-Taschenbuch-ebook/dp/B01BHWFX6/</a> (Amazon)
	Ursula Grawe			
	Karin von Schwab	<i>Stolz und Vorurteil</i>	Null Papier (Düsseldorf, Germany, 2018)	<a href="https://www.amazon.de/Stolz-Vorurteil-Vollständige-Ausgabe-Klassiker-ebook/dp/B007506EGS/">https://www.amazon.de/Stolz-Vorurteil-Vollständige-Ausgabe-Klassiker-ebook/dp/B007506EGS/</a> (Amazon)
Dutch	Margret Stevens & Trots en vooroordeel		Athenaeum-Polak & Van Gennep (Amsterdam, The Netherlands, 2014), 2nd edn	<a href="https://www.amazon.de/Trots-vooroedeel-Dutch-Jane-Austen-ebook/dp/BOONYLO1W/">https://www.amazon.de/Trots-vooroedeel-Dutch-Jane-Austen-ebook/dp/BOONYLO1W/</a> (Amazon)
	Annelies Roeleveld			
Spanish	Ana M. Rodriguez	<i>Orgullo y prejuicio</i>	Penguin Clásicos (Barcelona, Spain, 2016)	<a href="https://www.amazon.com/Orgullo-prejuicio-mejores-clasicos-Spanish-ebook/dp/B00XJQ09BG/">https://www.amazon.com/Orgullo-prejuicio-mejores-clasicos-Spanish-ebook/dp/B00XJQ09BG/</a> (Amazon)
	José de Utrries	<i>Y Orgullo y prejuicio</i>	Biblioteca Luna (Madrid, Spain, 2017)	<a href="https://www.amazon.com/Orgullo-Prejuicio-Spanish-Jane-Austen-ebook/dp/B07SGKTPP/">https://www.amazon.com/Orgullo-Prejuicio-Spanish-Jane-Austen-ebook/dp/B07SGKTPP/</a> (Amazon)
	Azara			
	Marciano Guerrera	<i>Primeras Impresiones</i>	CreateSpace Independent Publishing Platform (2018)	<a href="https://www.amazon.com/Primeras-Impresiones-Orgullo-Prejuicio-Illustrated-ebook/dp/B07CMJ8S4V/">https://www.amazon.com/Primeras-Impresiones-Orgullo-Prejuicio-Illustrated-ebook/dp/B07CMJ8S4V/</a> (Amazon)
		<i>u Orgullo y Prejuicio</i>		
French	Valentine Leconte	<i>Les cinq filles de Mrs Bennet</i>	Plon (Paris, France, 1932)	<a href="https://fr.wikisource.org/wiki/Les_Cinq_Filles_de_Mrs_Bennet/Texte_entier">https://fr.wikisource.org/wiki/Les_Cinq_Filles_de_Mrs_Bennet/Texte_entier</a> (Wikisource)
	& Charlotte Pres-Bennet			
		<i>soir</i>		
Latin	Thomas M. Cotton	<i>Superbia Et Odium</i>	Phaselus Publishing (Rhodes, UK, 2015)	<a href="http://www.lulu.com/shop/tom-cotton/superbia-et-odium/ebook/product-22288563.html">http://www.lulu.com/shop/tom-cotton/superbia-et-odium/ebook/product-22288563.html</a> (Lulu)
Polish	Anna Przedelska	<i>Duma i uprzedszenie</i>	Proszynski Media (Warsaw, Poland, 2012)	<a href="http://www.rulit.me/books/duma-i-uprzedszenie-read-270058-1.html">http://www.rulit.me/books/duma-i-uprzedszenie-read-270058-1.html</a> (Electronic Library Rulit)
	Treczakowska			
Russian	И. С. Маршак	<i>Гордость и предубеждение</i>	Художественная литература (Москва, CCCP [Moscow, USSR], 1988)	<a href="http://lib.ru/INOOLD/OSTEN/gord.txt">http://lib.RU/INOOLD/OSTEN/gord.txt</a> (Maxim Mashkov's Library)
Finnish	O. A. Joutsen	<i>Vlypeys ja ennakkoluulo</i>	Werner Söderström Osakeyhtiö (Porvoo, Finland, 1922)	<a href="http://www.gutenberg.org/cache/epub/45186/pg45186.txt">http://www.gutenberg.org/cache/epub/45186/pg45186.txt</a> (Project Gutenberg)
Hungarian	Szenczi Miklós	<i>Büszkeség és báltiélet</i>	Europa Könyvkiadó (Budapest, Hungary, 1958)	<a href="http://mek.oszk.hu/00300/00317/">http://mek.oszk.hu/00300/00317/</a> (Hungarian Electronic Library)
Basque	Ana Isabel Morales	<i>Harrostanua eta aurre-Elkar argitaletxea</i>	(Donostia, Spain, 2013)	<a href="https://play.google.com/store/books/details/Austen_Jane_Harrostanua_eta_aurrejuzguak?id=cMrNDQAAQBAJ">https://play.google.com/store/books/details/Austen_Jane_Harrostanua_eta_aurrejuzguak?id=cMrNDQAAQBAJ</a> (Google Play)
		<i>juzgiak</i>		
Korean	원유경(元裕卿)	오만과 편견(傲慢과 주식회사)	열린책들(株式會社 열린책들) (대한민국 [South Korea], 1986)	<a href="https://play.google.com/store/books/details/%EC%A0%9C%EC%9D%B8_%EC%98%A4%EC%8A%A4%ED%8B%B4_%EC%98%偏見)"></a> (Google Play)
Turkish	İşil Önder	<i>Gurur ve Önyargı</i>	İletişim Yayımları (Istanbul, Turkey, 2019)	<a href="https://play.google.com/store/books/details/Jane_Austen_Gurur_ve_%C3%96nyarg%C4%B1?id=EfCEDwAAQBAJ">https://play.google.com/store/books/details/Jane_Austen_Gurur_ve_%C3%96nyarg%C4%B1?id=EfCEDwAAQBAJ</a> (Google Play)
<i>Indo-European</i>				
English	Charlotte Brontë	<i>Jane Eyre</i>	Service & Paton (London, UK, 1897)	<a href="http://www.gutenberg.org/cache/epub/1260/pg1260.txt">http://www.gutenberg.org/cache/epub/1260/pg1260.txt</a> (Project Gutenberg)
	Moira Buffini	<i>Jane Eyre</i>	Internet Movie Scripts Database (2008)	<a href="https://www.imsdb.com/scripts/Jane-Eyre.html">https://www.imsdb.com/scripts/Jane-Eyre.html</a>
German	Maria von Borch	<i>Jane Eyre</i>	Jazzybee Verlag (Altenmünster, Germany, 2012)	<a href="https://www.amazon.com/gp/product/B009XCD1D4/">https://www.amazon.com/gp/product/B009XCD1D4/</a> (Amazon)
French	Noémie Lébazeilles	<i>Jane Eyre ou Les mé-Librarie Hachette</i>	(Paris, France, 1847)	<a href="http://www.gutenberg.org/cache/epub/16235/pg16235.txt">http://www.gutenberg.org/cache/epub/16235/pg16235.txt</a> (Project Gutenberg)
	Souvestre	<i>moires d'une institutrice</i>		
Russian	Б. Станевич	<i>Джейн Эйр</i>	Издательство «Правда» (Москва, CCCP [Moscow, USSR], 1988)	<a href="http://lib.RU/INOOLD/BRONTE/janeair.txt">http://lib.RU/INOOLD/BRONTE/janeair.txt</a> (Maxim Mashkov's Library)
Finnish	Tyne Haapanen	<i>Kotipoettajattaren ro</i>	Werner Söderström Osakeyhtiö (Porvoo, Finland, 1921)	<a href="http://www.gutenberg.org/cache/epub/47275/pg47275.txt">http://www.gutenberg.org/cache/epub/47275/pg47275.txt</a> (Project Gutenberg)
	Tallgren	<i>maani</i>		
Hungarian	Ruzitska Mária	<i>Jane Eyre</i>	Europa Könyvkiadó (Budapest, Hungary, 1972)	<a href="http://mek.oszk.hu/05600/05680/">http://mek.oszk.hu/05600/05680/</a> (Hungarian Electronic Library)
<i>Indo-European</i>				
English	Charles Darwin	<i>Origin of Species</i>	John Murray (London, UK, 1873), 6th edn.	<a href="http://www.gutenberg.org/cache/epub/2009/pg2009.txt">http://www.gutenberg.org/cache/epub/2009/pg2009.txt</a> (Project Gutenberg)
French	Edmond Barbier	<i>L'origine des espèces</i>	Schleicher Frères éditeurs (Paris, France, 1906)	<a href="http://www.gutenberg.org/cache/epub/14158/pg14158.txt">http://www.gutenberg.org/cache/epub/14158/pg14158.txt</a> (Project Gutenberg)
Polish	Szymon Dickstein	<i>O powstaniu gatun-Przegląd Tygodniowy</i>	(Warsaw, Poland, 1884)	<a href="https://wolnelektury.pl/katalog/lektura/darwin-o-powstaniu-gatunkow/">https://wolnelektury.pl/katalog/lektura/darwin-o-powstaniu-gatunkow/</a> (Polish Free Reading)
	& Józef Nusbaum	<i>ków droga doboru naturalnego</i>		
Russian	К. А. Тимирязев & П. Павлов	<i>Природа и происхождение</i>	Наука (Санкт-Петербург, CCCP [Saint Petersburg, USSR], 1991)	<a href="http://darwin-online.org.uk/content/frameset?itemID=F763b&amp;viewtype=text&amp;pageseq=1">http://darwin-online.org.uk/content/frameset?itemID=F763b&amp;viewtype=text&amp;pageseq=1</a> (Darwin Online)
		<i>и видов</i>		
		<i>А. П. Павлов &amp; П. А. Петровский</i>		
Finnish	Aarno Rafael Kos-Lajien synty		Arvi A. Karisto (Hämeenlinna, Finland, 1913)	<a href="http://www.gutenberg.org/cache/epub/52187/pg52187.txt">http://www.gutenberg.org/cache/epub/52187/pg52187.txt</a> (Project Gutenberg)
	kimkies			
Hungarian	Kampis György	<i>A fajok eredete</i>	Neumann Kht. (Budapest, Hungary, 2004)	<a href="http://mek.oszk.hu/05000/05011/">http://mek.oszk.hu/05000/05011/</a> (Hungarian Electronic Library)

Notes:  
 (1) Translators' names are given in native script and native order (surnames go first for Hungarian and Korean translators). Names of collaborating translators are separated by an & sign, irrespective of native practices. Books' titles and publishers' names are printed in their original languages. Publishers' addresses are translated into English if originally written in the Latin script. Otherwise, English translations (in brackets) accompany the original addresses in non-Latin script.  
 (2) Freely available electronic resources are marked with an asterisk. Descriptions of websites are translated into English, wherever applicable.  
 (3) The German (resp. Spanish) version of *Pride and Prejudice* used in the preparation of Table S2 is translated by Christian Grawe & Ursula Grawe (resp. Ana M. Rodriguez). The same will also apply to Fig. S5 (resp. Fig. S7).

<sup>8</sup>We do not consider questions or dependent clauses, whose word orders might differ from those of the declarative sentences as independent clauses.

<sup>9</sup>Both French and Spanish switch to SOV order when the object is a pronoun. Since we ignore pronouns (which belong to function words, as opposed to content words) in our Markov model, we can regard both French and Spanish as fully SVO languages.

<sup>10</sup>Here, the (optionally present) second verb (V<sub>2</sub>) is in a non-finite form, such as an infinitive or a participle. In addition, Danish, Dutch and German are V2 (verb-second, not to be confused with the aforementioned V<sub>2</sub> notation) languages, placing the conjugated verb in the second position of a sentence, even when its preceding constituent is not the subject.

<sup>11</sup>In principle, one can freely switch word orders in Latin without affecting the meaning of a sentence: “Darcy loves Elizabeth” can be translated by all the six possible permutations of the words *Darceius*, *amat* and *Elizabetham*. In practice, the SOV order (such as *Darceius Elizabetham amat*) is dominant. The situation of Basque is similar to that of Latin. Korean and Turkish are both strictly SOV.

<sup>12</sup>Like Latin, these languages mark overtly subjects and objects through declensions, so word order is not important. It is arguable whether these languages should be classified as (predominantly) SVO, just as Latin is classified as SOV (by default).

**Table S2. Statistical dictionary generated from bipartite matching of semantic fingerprints, using Jane Austen's *Pride and Prejudice*, and its translation into an Indo-European language**

(en) English	(da) Danish	(de) German	(nl) Dutch	(es) Spanish (fr) French	(la) Latin	(pl) Polish	(ru) Russian	(en) English	(da) Danish	(de) German	(nl) Dutch	(es) Spanish (fr) French	(la) Latin	(pl) Polish	(ru) Russian
absence		aufwezigkeit		absence				laugh		lachen	lachen	rire		list	письмо
admire	bewundern	Bewunderung						letter	Brief	brief	carta	lettre		bibliotek	библиотеку
agreeable	bequiglich							library	biblioteket	bibliothek	biblioteca		Lizzy	Lizzy	Лиззи
asked								Lizzy	Lizzy	Lizzy	Lizzy	Lizzy	Lizzy	Lizzy	Лондон
aunt	Tante	tante	ta	ta	tante	matertera	cioika	тетушка							
ball	Ball	bal					balu	бал							
beauty				hermosa											
believe	tror	geloven													
Bennet	Bennet	Bennet	Bennet	Bennet	Bennet	Bennet	Bennet	Bennet	Беннет	Беннет	Беннет	Беннет	Беннет	Беннет	Беннет
Bingley	Bingley	Bingley	Bingley	Bingley	Bingley	Bingley	Bingley	Bingley	Бингли	Бингли	Бингли	Бингли	Бингли	Бингли	Бингли
book	Buch	boek	boek												
Bourgh	de														
breakfast															
Brighton	Brighton	Brighton	Brighton	Brighton	Brighton	Brighton	Brighton	Brighton	Брайтон	Брайтон	Брайтон	Брайтон	Брайтон	Брайтон	Брайтон
brother	Bruder	broer	broer	hermano	frere	fratris									
carriage	Kutsche	rijtuig	coche	vouiture											
case	tilfælde														
Catherine	lady	Lady	lady	Catherine	Catherine	Caterina	Katarzyna	Катарина							
character															
Charlotte	Charlotte	Charlotte	Charlotte	Charlotte	Charlotte	Carlotta	Charlotte	Шарлотта	характер						
children	born	Kind	kinderen		enfants				детей						
choose															
civility															
Collins	Collins	Collins	Collins	Collins	Collins	Collins	Collins	Collins	Коллинз	Коллинз	Коллинз	Коллинз	Коллинз	Коллинз	Коллинз
colonel	oberst	kolonel	coronel	colonel	colonel	colonel	colonel	colonel							
comfort	trost	Trost													
compliment	kompliment			compliment	cumplido										
congratulations															
consolation															
cousin	Cousine			cousin	cousin	cousin	кузина	кузина							
cried	udbrod	rip	exclamó	écria											
dance	dansen	dansen	dansen												
Darcy	Darcy	Darcy	Darcy	Darcy	Darcy	Darcieus	Darcy	Darcy	Дарси	Дарси	Дарси	Дарси	Дарси	Дарси	Дарси
daughter	datter	Tochter	dochter	hija	dia	filia	córki	дочь							
day	dag	Nachricht	dag	dia	die	die	dmia	день							
de	Bourgh	Bourgh	Bourgh	Bourgh	Bourgh	Bourgh			Бург						
dear	kære	Reve	chère	carissima											
deceived															
Derbyshire	Derbyshire	Derbyshire	Derbyshire	Derbyshire	Derbyshire	Derbyshire	Derbyshire	Derbyshire	Дербишир	Дербишир	Дербишир	Дербишир	Дербишир	Дербишир	Дербишир
deserve	fortjent														
dine	middag	Essen	dineren		diner		obiad	обяд							
disappointment	skuffelse	Enttäuschung													
dislike															
distance							distance								
door	doren	Tür		puerta			drzwi								
elegant		elegante													
Elizabeth	Elizabeth	Elizabeth	Lizzy	Elizabeth	Elizabeth	Elizabeth	Elizabeth	Elizabeth	Элизабет	Элизабет	Элизабет	Элизабет	Элизабет	Элизабет	Элизабет
entered															
estate															
evening	aflenen	Abend	avond	velada	soir	vesperem	wieczór	вечер							
eyes				ogen	ojos	yeux	oculos								
father	far	Vater	vader	padre	pere	pater	ojec	отец							
Fitzwilliam	Fitzwilliam	Fitzwilliam	Fitzwilliam	Fitzwilliam	Fitzwilliam	Fitzwilliam	Fitzwilliam	Fitzwilliam							
five	fem	fünf	vijf	cinco			bijć	мти,							
forget															
Forster	Forster	Forster	Forster	Forster	Forster	Forster	Forster	Forster	Форстер	Форстер	Форстер	Форстер	Форстер	Форстер	Форстер
fortnight															
fortune	formue	Vermögen		fortuna	fortune										
four	Gardiner	Gardiner	Gardiner	vier	cuatro	Gardiner	Gardiner	Gardiner	четыре	четыре	четыре	четыре	четыре	четыре	четыре
Gardiner	gentlemen														
girls															
glad															
good	godt	gut			mejor	mejor									
hand															
handsome															
happy															
honour	ære	Ehre	eer	honor			honor								
house	huset	Haus	huis				domum								
Hurst															
husband	egtemand														
ignorant	uvidende														
impossible	umungkin	unmöglich													
intended	hensigt														
invitation	Einladung	uitnodiging	invitación	invitation	invitation		zaproszenie								
Jane	Jane	Jane	Jane	Jane	Ioanna	Jane	bodźzy	Джейн							
journey															
Kitty	Kitty	Kitty	Kitty	Kitty	Kitteja	Kitty	Kitty	Kitty							
know	weiß	weet													
lady							lady	леди							

(1) See Table S1 for provenances of texts.

(2) This table is sorted alphabetically according to the English entries.

(3) For clarity, only the most frequent word within a word pattern is tabulated. See Figs. S4–S11 for details of bipartite matching.

(4) Color encoding follows [1, Fig. 5c]: green = correct match; amber = close match; red = incorrect match.

(5) Column separator demarcates different subgroups (Germanic, Romance & Latin, Slavic) of the Indo-European language family. A Germanic compound noun (such as Danish *morbor*) may appear as the translation of more than one word.

Table S3. Statistical dictionary generated from bipartite matching of semantic fingerprints, using Jane Austen's *Pride and Prejudice*, and its translation into a non-Indo-European language

(en) English	(fi) Finnish	(hu) Hungarian	(eu) Basque	(ko) Korean	(tr) Turkish	(en) English	(fi) Finnish	(hu) Hungarian	(eu) Basque	(ko) Korean	(tr) Turkish
admire				노우란(好物) 외	bayranlık	know			uko		
arrival				oxikamendu		laugh	aura				
attachment				zeko		letter	kirjeen	levelet			
aunt	äiti			누노 외(嫗翁) 외		library	kirjasto	könyvtárszabába			
ball	bal					like					
believe	uskua			Bennet		Lizzy	Lizzy	Lizzy			
Bennet	Bennet	Bennet	Bennet	베넷(Bennet)	Bennet	London	London	London			
Bingley	Bingley	Bingley	Bingley	빌링(Bingley)	Bingley	long	päästä	akarta			
book						Lucas	Lucas	Lucas			
Bourgh	Bourgh	Bourgh	Bourgh		Bourgh	Lydia	Lydia	Lydia			
breakfast						man	mies	ember			
Brighton	Brightonba	Brightonba	Brightonera	브라이턴(Brighton) 외	Brighton'a	marriage	Mary	Mary			
brother	väijensä	bätyja	neba			Mary	Mary	Maryk			
carriage	vaunut		koxtea			master		agazaba			
Catherine	lady	Catherine	lajelmet	캐시 린(Catherine)	araba	means					
character						Meryton	Merytonissa	Merytonban			
Charlotte	Charlotte	Charlotte	izaera	찰리(Charlotte) 외		Miss	ni	anderehou			
children	lapseni	gyermekem	Charlotte			money	rahoja	pénzt			
civility		luvariasan				mother	atti	anyja			
Collins	Collins	ezredes	Collins			Mr	hma				
colonel	eversti	ezredes	koronelak			Mrs	rva	Mrs.			
come		jou	etori	오는		neighborhood					
compliment						nephew	sisarenpoikaansa	unokaöccese			
congratulations	omniteli					Netherfield	Netherfieldin	Netherfieldi			
consolation						niece		smokahunga			
country						officers	ah				
cousin	erkkuna					oh					
cried						pain	puiston	park			
dance						park	Pemberley	Pemberleyben			
Darcy	Darcy					Philips	Philips	Philips			
daughter	tytärensä					poor					
day	päivänä	lädyes	de	시간(時間) 읍 데(De)	gun	poor					
de		Bourgh				pounds					
dear						pretty					
Derbyshire	Derbyshire-ben					pride					
different						promised					
dim						read					
dinner	vacorara					refuse					
disappointment	csalodott					relations					
door						respect					
elegant						room	Rosingsissa	Rosingsi			
Elizabeth	Elizabeth					Rosings	Rosingsissa	Rosingsen			
evening	illalla	este				said					
eyes						sat	istui				
face	kasvot					see					
father	isä	apja	anta	아빠(父) 가		send					
felt		erzett				serious					
fine	Fitzwilliam					servant					
Fitzwilliam	Fitzwilliam					sir					
five	vitsi					suffer					
forget						surprise					
Forster	Forster					table					
four		negy	lau	포스터(Forster)		ten					
Gardiner	Gardiner		Gardner	가디너(Gardiner)		thank					
gentlemen			zaldun			thousand					
girls			neska			town					
good			yna			turned					
hand	kezét					two					
happy	boldog					uncle					
honour	kunnia		oboreca	행복(幸福) 외		vain					
house	ház					visit					
Hurst	írghez					way					
husband	mahdoton		senarra	남친(男友) 외		week					
impossible						watus					
intended						Wickham	Wickham				
interest						vaimosa					
introduction						William	William				
invitation	kutsua	Jane	Janek	초개(前介) 외		akkunasta	ablaikhoz				
Jane	Jane	Jane	Kitty	초대(招待) 외		davet	irt				
Kitty	Kitty	Kitty	Kitty	제인(Jane) 외		write	irt				
				제인(Kitty) 외		years	irt				
						young	nuori				

(1) See Table S1 for provenances of texts.

(2) This table is sorted alphabetically according to the English entries.

(3) For clarity, only the most frequent word within a word pattern is tabulated. See Figs. S12–S16 for details of bipartite matching.

(4) Color encoding follows [1, Fig. 5c]: green = correct match; amber = close match; red = incorrect match.

(5) Column separator demarcates different language families.

(6) Korean stems that derive from Chinese or English are accompanied by their etymologies in parentheses.

**Table S4. Statistical dictionary generated from bipartite matching of semantic fingerprints, using Charlotte Brontë's *Jane Eyre*, and its translation into a European language**

Table S5. Statistical dictionary generated from bipartite matching of semantic fingerprints, using Charles Darwin's *Origin of Species*, and its translation into a European language

Table S6. Statistical dictionaries compiled from *Pride and Prejudice* and its translations into French and Korean, with different thresholds for Ružička similarities

## 1.4 Automated question answering based on Markov semantics

When we are given a document (of moderate length) and a natural language question as input, we rate and rank the sentences within the document by their relevance to the question, containing topical patterns  $\mathcal{Q} = \{W_{q_1}, \dots, W_{q_K}\}$ . We expand the query into  $\mathcal{Q} \cup \mathcal{Q}'$ , a union of semantic cliques:  $\mathcal{Q} \cup \mathcal{Q}' = \bigcup_{k=1}^K \mathcal{S}_{q_k}$ . We build a localized Markov matrix  $\mathbf{P} = (p_{ij})_{1 \leq i,j \leq N}$  on  $\mathcal{Q} \cup \mathcal{Q}'$ . We further use the Brin–Page damping [15] to form an ergodic Markov matrix  $\tilde{\mathbf{P}} = (\tilde{p}_{ij})_{1 \leq i,j \leq N}$ , where  $\tilde{p}_{ij} = 0.85 p_{ij} + \frac{0.15}{N}$ .

If our question  $Q$  contains words from  $W_{Q_1}, \dots, W_{Q_m} \in \mathcal{Q}$  and a candidate answer  $A$  contains words  $W_{A_1}, \dots, W_{A_n} \in \mathcal{Q} \cup \mathcal{Q}'$  (counting multiplicities, but excluding function words and patterns with fewer than 3 occurrences in the reference document), then we assign the following entropy production score [1, (11) and Footnote 4]

$$\mathcal{F}[Q, A] := - \sum_{i=1}^m \sum_{j=1}^n \tilde{\pi}_{Q_i} \tilde{p}_{Q_i A_j} \log \tilde{p}_{Q_i A_j} \quad (1.14)$$

to this question-answer pair.

In [1, Fig. 7a,b], we test our Markov language model on the WikiQA dataset. The distribution of average precisions (AP) and reciprocal ranks (RR) for individual questions in the WikiQA dataset is summarized in Fig. S2. One may wish to check Tables S11–S12 in §4 and see how our algorithm performs on each of the 1242 questions.

## 2 What our software does not

## 2.1 Non-universal statistics on short time scales

As recapitulated in §1.1, we reject text fragments shorter than or equal to the lengths of the flanking words. This is because we find such short-range features highly sensitive to grammatical rules that are specific to certain languages, thus lacking universality.

In [1, Footnote 2], we have already alluded to the highly versatile reduplications (or lack of) in different languages: German *liebe Studentinnen und Studenten* “dear (female and male) students”, Malay *orang-orang* “people” (vs. *orang* “person”), and Hawaiian *wiki wiki* “very quick” (vs. *wiki* “quick”). We choose to ignore such short-range recurrence statistics in our current study, which are typically not transferrable across languages and cultures.

The numerical instability of short-range textual features can be further illustrated by a comparison of Turkish and English morphemes (adapted from [17, p. 61]):

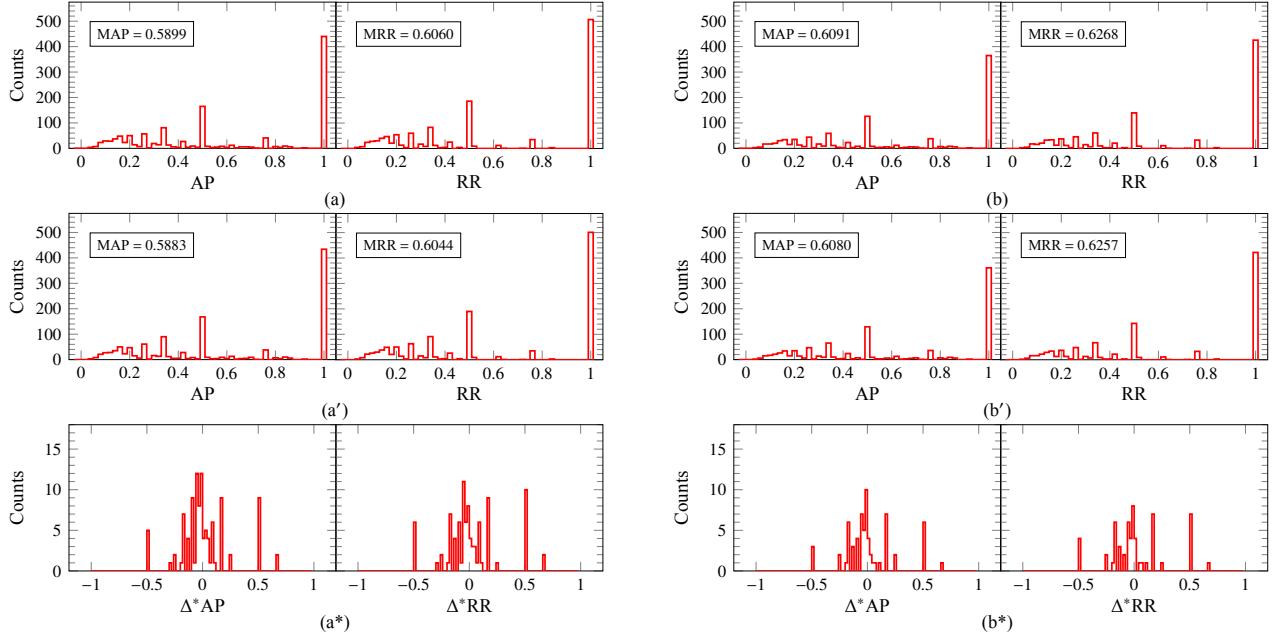


Fig. S2. Score distributions in WikiQA tests. (a) Test results on 1242 answerable questions the WikiQA dataset: histograms for average precisions (AP) and reciprocal ranks (RR) of correct answers, with their respective means displayed as inset. Tied scores are resolved by the McSherry–Najork algorithm [14], which effectively averages AP and RR over all possible permutations of tied entries. (b) Test results on a subset of 990 WikiQA questions (see Algorithm 3.12 in §3 for detailed screening criteria) that do not require logic inference beyond associative reasoning. (The Brin–Page damping factor 0.85 [15] is used in both panels a and b.) (a')–(b') Control experiments using the Bressan–Peserico damping factor 0.5 [16]. (a\*)–(b\*), Non-zero differences between Brin–Page and Bressan–Peserico values (upper panel minus lower panel).

Avrupa hlaştırılamayacaklardan sizin

You (all) are among those who will not be able to be caused to become Europeanized  $\approx$  You (all) are among those who will not be made Europeanizable

Here, in both languages, the relevance of a morpheme roughly decays monotonically with respect to the distance from the central morpheme (Avrupa/Europe). The distance correspondence is generally stabler for larger separations, but is highly unstable when we count distances between nearest (or next-to-nearest) neighboring morphemes.

Table S7. Length/time scales in natural language processing (NLP)

Scale	Friederici's hierarchy [18]	Jakobson's metaphor	Rôle in NLP	Evolutionary susceptibility
short-range	phonological level	“elementary particle” [19]	acoustic encoding	climatic/environmental [20, 21, 22]
↓	lexical level	“atom”	concept tagging	cognitive/cultural/social [23, 24, 25, 26]
long-range	sentence level	“molecule”	syntactic computing	cognitive/cultural/social [23, 24, 25, 26]
	text/discourse level	“bulk material”	semantic processing	—

In our current work, we mainly concern ourselves with theories and applications for the recurrence statistics of word patterns that exhibit universal behavior, with a focus on semantic processing [27] rather than syntactic computation [28]. Practically, this means that we will suppress the fine structures in the three short scales in Friederici's hierarchy (Table S7), which vary by language typology. Similar exclusions of short-range word contacts are also implemented in the  $n$ -gram language model of Brown *et al.* [29].

## 2.2 Word disambiguation and stylistic adaptability

In theory, we expect that each Markov state corresponds to exactly one concept (a family of words that are related to each other by inflection and derivation). In practice, there are homographs (words with the identical spelling but unrelated meanings), whose denotations must be resolved through careful examination of context (either by experienced human readers, or by sophisticated algorithms in supervised learning). In numerical implementations of our model, we do not make any attempt to solve the polysemy puzzle, so we do run the risk of conflating homographs into one Markov state. In languages like Danish and Dutch, this conflation is serious enough to affect the performance of certain computational tasks related to the following

Table S8. Parallel versions for a conversation in Chapter 1 of *Pride and Prejudice*

English	French	Russian	Finnish
“My dear Mr. Bennet,” <u>said his</u> lady to him one day, “have you heard that Netherfield Park is let at last?” Mr. Bennet replied that he <b>had not</b> . <u>I</u> gnorait.	— Savez-vous, mon cher ami, <u>dit</u> Netherfield Park est enfin loué ? Mr. Bennet répondit qu’il	— Дорогой мистер Беннет, — <u>сказала как-то раз</u> миссис Беннет своему мужу, — слышали вы, что Незерфилд-парк наконец больше не будет пустовать?	“Rakas Bennet”, <u>sanoi</u> tämän arvon herran puoliso michelleen eräänä päivänä, “oletko kuullut, että Netherfield Parkin kartano on vihdoinkin saanut vuokraajan?”
		Мистер Беннет ответил, что он этого <b>не слышал</b> .	Herra Bennet vastasi, ettei hän <b>ollut kuullut</b> .

Notes:

(1) See Table S1 for provenances.

(2) The underlined texts are roughly equivalent across four versions; the same can be said for the texts in **boldface**.Table S9. Parallel translations of the opening line in *Pride and Prejudice*

German (Werner Beyer)	(Christian Grawe & Ursula Grawe)	(Karin von Schwab)
In der ganzen Welt gilt es als ausgemachte Wahrheit, daß ein begüterter Junggeselle unbedingt nach <b>einer Frau</b> Ausschau halten muß.	Es ist eine allgemein anerkannte Tatsache, dass ein allein stehender Mann im Besitz eines hübschen Vermögens an- geblich nichts dringender braucht als <b>eine Frau</b> .	Es ist eine Wahrheit, über die sich alle Welt einig ist, daß ein unbewiebter Mann von einem Vermögen unbedingt auf der Suche nach <b>einer Lebensgefährtin</b> sein muß.
Spanish (Ana M. Rodriguez)	(José de Urries y Azara)	(Marciano Guerrero)
Es una verdad reconocida por todo el mundo que un soltero dueño de una gran fortuna siente un día u otro la necesidad de <b>una mujer</b> .	Es verdad universalmente admitida que un soltero poseedor de una buena fortuna tiene que necesitar <b>una mujer</b> .	Es una verdad mundialmente reconocida que un hombre soltero y poseedor de una buena fortuna, debe estar en busca de <b>una esposa</b> .

Notes:

(1) See Table S1 for provenances.

(2) The underlined texts are roughly equivalent across four versions; the same can be said for the texts in **boldface**.(common but ambiguous) words:<sup>13</sup>

Dutch	English gloss	(German cognate)	Danish	English gloss
meer	{ more lake }	(mehr) (Meer)	så	{ saw (past of see) so (to) sow }
vroeg	{ asked early }	(frug) (früh)	ved	{ at knows wood }

We refer our readers to §9.2 for a discussion of homographs in Korean.

Besides homographs, there are three more subtler noise sources that may negatively impact the translational invariance of the Markov spectrum, hence our automated word translation via cross-lingual matching of semantic fingerprints.

The first kind of noise is attributed to free variations of translators (Tables S8 and S9). As shown in Table S8, translators may take the liberty of using free (as opposed to literal) interpretations: the underlined texts in the three translations essentially read “Mrs. Bennet said to her husband one day”, which paraphrases the interpersonal relationship expressed in the original English text; the bold texts in the three translations all explicitly supply the negated verb (“had not known”/“had not heard”), which is omitted in the English original. In Table S9, all the underlined texts correspond to “a truth universally acknowledged” in the English original and all the phrases in boldface mean “a wife”, despite nuances in their wording and phrasing.

The second kind of noise is due to algorithmic imperfection in morphological classifications of words into word patterns. There are irregularities in word morphologies that cannot be fully covered by our automated word clustering algorithms described in Part III. Furthermore, the heterogeneous etymological origins<sup>14</sup> of the English language may bring us additional challenges. For example, our programs do not recognize *heart* and *cordial* as morphologically related in English, even though their German or Hungarian translations display such relationship manifestly:

English	German	Hungarian
heart	Herz	szív
hearty, cordial	herzlich	szíves
heartily, cordially	herzlich	szívesen

The third kind of noise is due to our insufficient treatment of space, time and causality in the lexicosemantics of verbs (see [30, Chap. 4] or [31, Chap. 5]). For the time being, we have swept all prepositions (which are vital to the space component of verbs) into the class of function words, and have ignored them in our text mining tasks. Thus, our numerical translation algorithm cannot handle polysemic verbs like *go*, *make*, *take*, *turn*, whose meanings are best clarified by subsequent prepositions

<sup>13</sup>In modern German, *Meer* means “sea”. In archaic German, *Meer* means “lake”. In modern German, the strong conjugation *frug* for first- or third-person singular preterite is rare, while the weak conjugation *fragte* is common. We have included these German cognates, to serve as visual aids to their Dutch counterparts.

<sup>14</sup>English lexicon possesses a native Germanic stock (exemplified by *hearty*, *heartily* and so forth), but is simultaneously under strong Latinate influence (shown in near-synonyms of the aforementioned words *cordial*, *cordially*), via Norman French. A more challenging scenario appears in Korean word families with disparate (Chinese/English/Korean) origins. See §9.2 for a discussion.

(if present). We have chosen to ignore prepositions (in either isolated form or as separable prefixes<sup>15</sup>) in our numerical text analysis, because they do not exhibit universal behavior, in at least two senses: (1) Short-range relations between verbs and prepositions (like *takes up*) in English may translate into long-range relations in certain Germanic languages (like *nimmt ... auf* in German and *neemt ... op* in Dutch, where the written space occupied by the ellipsis can be arbitrarily long). (2) The rôles of locative prepositions (like *on*, *in* and *from*) are largely subsumed by certain case endings in Uralic languages (such as Finnish and Hungarian). In our work, we not only ignore prepositions, but also strip away case endings in a stemming procedure (Part III). We believe that such simultaneous treatments of prepositions and case endings are consistent: they generate similar levels of coarse-graining in words across different language families. In addition to neglecting prepositions, we have designed our stemming procedure so as to merge verb forms in different tenses and aspects (time component), but it may also conflate causative and non-causative forms (causality component), in a fashion that is not necessarily consistent across all languages. For example, our stemming algorithm treats Finnish *syödä*/Korean 먹다 “(to eat)” (non-causative) and Finnish *syöttää*/Korean 먹이다 “(to feed)” (causative) as the same verb in their respective language, but not so for their English counterparts.

The language-specific usage of function words is another practical hindrance to lossless translation, but this hindrance does not contribute to our Markov spectrum (if we can ignore effects from pronominal coreferences, as discussed below) for content words. As a glance of the highly diverse typologies of function words found in different languages, we say a few more words about the gender agreement of pronouns. All the 8 modern Indo-European languages under our consideration in Part III distinguish “he” from “she” when referring to humans as subjects in a sentence, while the 5 non-Indo-European languages do not.<sup>16</sup> Gender agreement of possessive pronouns is a subtler issue for Indo-European languages. Expanding the example given by Willim and Chomsky [32, p. 100], we have

English:	<i>John told Mary about his father.</i>	<i>John told Mary about her father.</i>
French:	<i>Jean a parlé de son père à Marie.</i>	<i>Jean a parlé de son père à Marie.</i>
German:	<i>Johann hat Maria von seinem Vater erzählt.</i>	<i>Johann hat Maria von ihrem Vater erzählt.</i>
Polish:	<i>Jan opowiadał Marii o swoim ojcu.</i>	<i>Jan opowiadał Marii o jej ojcu.</i>
English:	<i>John told Mary about his mother.</i>	<i>John told Mary about her mother.</i>
French:	<i>Jean a parlé de sa mère à Marie.</i>	<i>Jean a parlé de sa mère à Marie.</i>
German:	<i>Johann hat Maria von seiner Mutter erzählt.</i>	<i>Johann hat Maria von ihrer Mutter erzählt.</i>
Polish:	<i>Jan opowiadał Marii o swojej matce.</i>	<i>Jan opowiadał Marii o jej matce.</i>

Here, the pronouns *his*, *her*, *seinem*, *seiner*, *ihrem*, *ihrer* and *jej* agree with the gender of the possessor; meanwhile the pronouns *son*, *sa*, *seinem*, *seiner*, *ihrem*, *ihrer*, *swoim* and *swojej* agree with the gender of the possessed. Automated resolution of such agreements is beyond the scope of our current research, as this language-specific problem occurs on a different neuro-linguistic scale from that of our concern. However, if pronominal coreferences become ambiguous enough after literal translation, then the translator might choose to supply personal names in place of pronouns, and such rephrasing may partly affect the performance of our Markov algorithms. (See Example 4.15.2 and Fig. S3e.)

### 2.3 Causal inference and deductive reasoning

Admittedly, our semantic cliques (see [1, §2.3.2] and §1.2 above) have some limitations. At relatively low computational cost, our criterion for semantic dependence [see (1.5)] does not distinguish between causal [33, 34] and non-causal relations, nor does it generate semantic webs with hierarchical topologies [35, 36]. For comprehension tasks requiring high-precision causal inference and relationship mining, Bayesian networks [33, 34] and persistent homologies [35, 36] offer more reliable guidance than our approach.

The lack of causal inference and deductive reasoning in our question-answering machine places a cap on its MAP and MRR scores ([1, Fig. 7b] and Fig. S2). No matter how we improve the stemming/clustering algorithms for English (§4), there are certain fractions of WikiQA questions that our method cannot handle properly (Table S12). Typical failures include, but are not limited to the following categories (see Table S10 for evaluations of our performance):

- Our semantic cliques do not always enable us to address causation questions (like *What causes ...?*) correctly, with only about a fifty-fifty chance of hitting the right answer. These cliques indicate association and connection between concepts, which may or may not be of causal nature.
- Our semantic cliques do not automatically sieve answers with numbers (spelt out in English or Arabic numerals), even if the question (such as *How big ...? How long ...? How many ...? How old ...? What year ...?*) elicits one. Perhaps

<sup>15</sup>For example, the prefixes (*ab/af*, *auf/op*, and *zu/toe*) in German *abnehmen/Dutch afnemen* “(to) decrease”, German *aufnehmen/Dutch opnemen* “(to take up”, and German *zunehmen/Dutch toenemen* “(to) increase” must be separated from the verb stems in certain conjugated forms. Note however, that the prefixes in English words “decrease” and “increase” are not separable as prepositions.

<sup>16</sup>In some rare circumstances, the lack of gender-specific pronouns may also hamper topic analysis in our numerical experiments. (See Example 4.15.2 and Fig. S3e.)

Table S10. Distributions of hits and misses in certain types of questions in our WikiQA experiments

Question type	#Hits	#Misses	Question type	#Hits	#Misses
What caused ...?	1	1	How much ...?	5	9
What causes ...?	4	4	How often ...?	0	2
What is caused ...?	1	0	How old ...?	4	5
How big ...?	0	2	What became of ...?	0	1
How deep ...?	0	1	What happened to ...?	0	4
How long ...?	3	5	What year ...?	7	14
How many ...?	19	51	When ...?	46	67

the statistical behavior of recurring numbers in a text document does not neatly fit into our Markov model for word patterns, so the semantic cliques do not easily detect numbers as synonyms to certain question words.

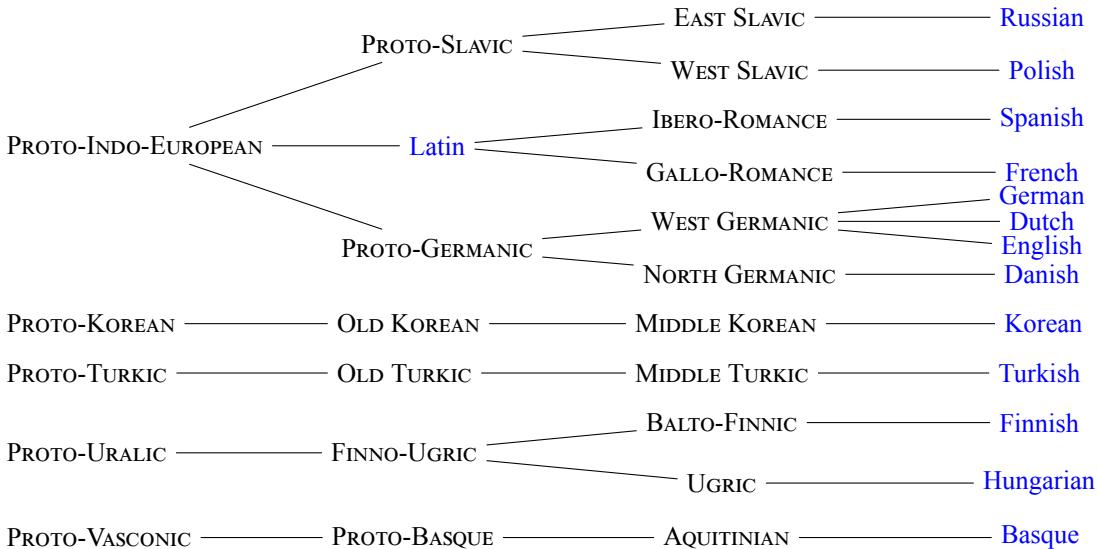
- Our semantic cliques stumble over questions that ask for superlatives, either explicitly or implicitly. Some of these questions (*What became of ...?* *What happened to ...?*) can be very subtle, as they require finding the most recent date out of a list, going beyond the reach of our Markov model.

## Part II

# Protocols for text cleansing

After a brief summary of notations (§3.1), we present in this Part cleansing protocols for certain articles in the English Wikipedia (§3.2) and electronic books written in 14 languages spoken in Europe and Asia (§3.3), in preparation for word clustering and word translation experiments in Part III.

Hereafter in Parts II and III, aside from English (§4), we pick Danish, German and Dutch as representative Germanic languages (§5); Spanish, French and Latin<sup>17</sup> as representative Romance languages (§6); Polish and Russian as representative Slavic languages (§7); Finnish and Hungarian as representative Uralic languages (§8); Basque, Korean and Turkish as representatives from various other language families (§9). A putative evolutionary history of these languages is sketched below:



From now on, the languages other than English are first sorted alphabetically according to their genealogical affiliation (Germanic, Romance, Slavic, Uralic, “Various”), and then sorted within their respective affiliation according to their ISO codes [e.g. Spanish (es) precedes French (fr)]<sup>18</sup>.

All these Asian and European languages under investigation have clearly defined word boundaries (typically marked by spaces and punctuation marks, and additionally by apostrophes and hyphens in the case of French) and rich word morphologies

<sup>17</sup>Strictly speaking, Latin is not a Romance language, but an ancestor to all the modern Romance languages. Since a modern Romance language, Romanian, still shares a lot of conservative features with Latin, it may not hurt to label Latin as a Romance language.

<sup>18</sup>The individual sections from §5 to §9 can be read in any order, depending on need. To help our readers navigate through these protocols, we have reproduced certain explanatory texts from one section to another.

(e.g. adjective comparisons, noun declensions and verb conjugations). Inspired by the sequence alignment procedure [37, 38] in bioinformatics, we devise algorithms that employ string matching to find affinity between words, with a limited amount of input for grammatical rules.

The Germanic, Romance and Slavic languages all belong to the Indo-European family, which distinguishes from the other language families under consideration, both lexically and syntactically. The languages within the same affiliation share a lot of common features in vocabulary and syntax. We note however that English, though officially a Germanic language, has borrowed extensively from Latin and (Norman) French, so it has blended characteristics of both Germanic and Romance subgroups. In the light of this, we have singled out the algorithmic treatment of English texts in §4.

Our word clustering algorithms are based on Porter stemming [39, 40], along with extensive additions that accommodate to irregularities in word morphologies. Since regular and irregular word morphologies are handled by different neurobiological mechanisms in the human brain [41, 42], we are not going to unify their treatment under a single algorithmic framework. Regular morphologies will be identified by substitution rules (as in Porter stemming), while irregular morphologies in each language will be accommodated by a finite (though sometimes long) list of exceptions. It has been observed that grammatical irregularities tend to be associated with highly common words [43]: the frequent uses of these words make their peculiarities harder to be forgotten. It is our hope that, via many explicitly coded “exceptions to general rules”, we can achieve higher precision in text mining tasks related to those highly frequent yet highly irregular words.

## 3 String manipulations and text normalizations

### 3.1 Notational conventions

In the source codes accompanying this work, we have implemented our text mining algorithms in *Mathematica* (including both the cleansing procedures in Part II and the clustering procedures in Part III). However, we hope that readers of this document could easily adapt our algorithms to other programming languages. Towards this end, we will present our methods of text processing using some standard notations in linguistic analysis (applicable to both Parts II and III), as declared below.

**Definition 3.1** (String length). The length of a text string  $\hat{\sigma}$  is a non-negative integer that equals the number of Unicode characters used to write out the string, denoted by  $\ell(\hat{\sigma})$ .

The empty string  $\emptyset$  (not to be confused the Danish letter  $\emptyset/\circ$ ) is the only text string with length zero:  $\ell(\emptyset) = 0$ .

For a positive integer  $n$  less than or equal to  $\ell(\hat{\sigma})$ , the  $n$ th position in string  $\hat{\sigma}$  is denoted by  $\hat{\sigma}^{\{n\}}$ , and the first  $n$  characters in string  $\hat{\sigma}$  is denoted by  $\hat{\sigma}^{\{n\}}$ . For a non-empty string  $\hat{\sigma}$ , its “first character” is  $\hat{\sigma}^{\{1\}}$  and its “last character” is  $\hat{\sigma}^{\{\ell(\hat{\sigma})\}}$ . For the empty string, its “first character” and “last character” are both defined as the empty string. In general, the last character of a string  $\hat{\sigma}$  is denoted by  $\Omega(\hat{\sigma})$ .

The notation  $\hat{\sigma}^{-1}$  stands for the reverse of the string  $\hat{\sigma}$ . In particular, we have  $\emptyset^{-1} = \emptyset$ ; when  $\ell(\hat{\sigma}) > 0$ , we have  $(\hat{\sigma}^{-1})^{\{n\}} = \hat{\sigma}^{\{\ell(\hat{\sigma})+1-n\}}$  for all integers in the range  $1 \leq n \leq \ell(\hat{\sigma})$ .  $\square$

**Definition 3.2** (String spelling and string alternatives). For all the languages written in either the Latin or the Cyrillic script, letters (including apostrophes and hyphens) within a text string in an algorithm is always spelt out in slant typeface, such as *word* (English) and *слово* (Russian). Boxed characters set in the *typewriter* typeface are taken verbatim. A boxed expression in red, in the form of U+xxxx, refers to the Unicode character with hexadecimal code `xxxx`. The expression  $\hat{\sigma}^+$  (resp.  $\hat{\sigma}_-$ ) results from converting  $\hat{\sigma}$  to upper (resp. lower) case.

The notation for “alternatives”  $\hat{\sigma}_1|\dots|\hat{\sigma}_n$  refers to an arbitrary text string in the set  $\{\hat{\sigma}_1, \dots, \hat{\sigma}_n\}$ . Parentheses in string spellings, usually in the form of  $\hat{\sigma}(\hat{\sigma}_1|\dots|\hat{\sigma}_n)$ , are used to enclose alternative portions that are not shared by a group of words.

The rules above also apply to algorithms for the Korean language, except that slant typefaces will not be used, and the native Korean alphabet does not distinguish upper and lower cases.  $\square$

*Example 3.2.1.* The boxed pattern , is interpreted literally, as a comma followed by a white space. The pattern (a|b) is a concatenation of five Unicode characters, while  $(a|b)$  refers to either member of the set  $\{a, b\}$ .

*Example 3.2.2.* For English, the string pattern *word*( $\emptyset|s|s's'$ ) refers to any text string in the set  $\{\textit{word}, \textit{words}, \textit{word's}, \textit{words'}$ . For Russian, the string pattern *слово*( $\emptyset|a|\emptyset$ ) refers to any text string in the set  $\{\textit{слово}, \textit{слова}, \textit{слов}\}$ . The expression  $(\textit{word}(\emptyset|s|s's'))^+$  has the same effect as *WORD*( $\emptyset|S|S'S'$ ), while  $((СЛОВ(O|A|\emptyset))_-$  is equivalent to *слово*( $\emptyset|a|\emptyset$ ).

**Definition 3.3** (Multiple occurrences of string patterns). For a string pattern  $(\hat{\sigma}_1|\dots|\hat{\sigma}_n)_m$  (resp.  $(\hat{\sigma}_1|\dots|\hat{\sigma}_n)_{m_0}$ ) denotes one (resp. zero) or more successive occurrences of the same pattern in a text string, and the notation  $(\hat{\sigma}_1|\dots|\hat{\sigma}_n)_{[n_1, n_2]}$  refers to between  $n_1$  and  $n_2$  (inclusive) consecutive appearances of the pattern  $(\hat{\sigma}_1|\dots|\hat{\sigma}_n)$ .<sup>19</sup> When  $v = n_1 = n_2$ , we also write  $(\hat{\sigma}_1|\dots|\hat{\sigma}_n)_{[v, v]}$  as  $(\hat{\sigma}_1|\dots|\hat{\sigma}_n)_{\times v}$ , to denote exactly  $v$  repeats of the pattern  $(\hat{\sigma}_1|\dots|\hat{\sigma}_n)$ .  $\square$

<sup>19</sup>In *Mathematica* codes,  $(\hat{\sigma})_m$  is represented by  $(\sigma)\dots$  or `Repeated[\sigma]`,  $(\hat{\sigma})_{m_0}$  by  $(\sigma)\dots$  or `RepeatedNull[\sigma]`, and  $(\hat{\sigma})_{[n_1, n_2]}$  by `Repeated[\sigma, {n1, n2}]`.

*Example 3.3.1.* The pattern  $(d|e)_m$  may be matched by strings  $d, e, dd, de, ed, ee, ddd, dde, ded, dee, edd, ede, eed, eee$  and so on. If we consider the pattern  $(d|e)_{m_0}$ , then the empty string  $\emptyset$  also forms a match. The pattern  $(d|e)_{[0,2]}$  is only matched by the strings  $\emptyset, d, e, dd, de, ed, ee$ .

**Definition 3.4** (Wildcard symbols). We write  $\mathbf{X}$  for a generic text string of zero or positive length. At times, we may restrict  $\mathbf{X}$  to a specific set, or may use subscripts (like  $\mathbf{X}_1$  and  $\mathbf{X}_2$ ) to put the generic string in context.

We write  $\mathbf{V}$  and  $\mathbf{C}$  for generic vowels and consonants, respectively. The classification schemes for vowels and consonants may depend on language and context. Like the case of  $\mathbf{X}$ , the exact identity of the wildcard symbol  $\mathbf{V}$  or  $\mathbf{C}$  may vary from entry to entry, line to line.

The symbol  $\hat{x}$  denotes a single letter character (a Unicode character other than digit or punctuation<sup>20</sup>). When we repeat it  $n$  times (as in  $\underbrace{\hat{x} \cdots \hat{x}}_{n \text{ times}}$ ), we mean exactly  $n$  letter characters (not necessarily identical) appearing successively in a text string.

The symbol  $\hat{x}$  stands for a single letter character (which serves the same function as  $\hat{x}$  when used alone). However, whenever  $\hat{x}$  is repeated, it always stands for the same character.  $\square$

*Example 3.4.1.* The string pattern  $\mathbf{X}aint$  ( $\mathbf{X} \in \{rt, mpl, tr\}$ ) is synonymous with the notation  $(rt|mpl|tr)aint$ , denoting any member among  $\{rtaint, mplaint, traunt\}$ . For succinctness, we also abbreviate  $\mathbf{X}aint$  ( $\mathbf{X} \in \{rt, mpl, tr\}$ ) in a superscript notation, as  $\mathbf{X}^e(rt|mpl|tr)aint$ . It is sometimes useful to specify the identity of a wildcard symbol by the set it does not belong to, such as  $\hat{x}e$  ( $\hat{x} \notin \{e, i\}$ ) denotes a string pattern where the letter  $e$  does not follow a letter  $e$  or  $i$ . In such cases, we will introduce a shorthand using subscript:  $\hat{x}_{\notin}(e|i)e$ , or use an overline<sup>21</sup> for the negation of a parenthesized string pattern:  $\hat{x}^{\infty}(e|i)e$ .

*Example 3.4.2.* The multiplicity notation (introduced in Definition 3.3) also applies to negated string patterns. For instance, in the notation  $\mathbf{X}^e(\overline{(eli)}_{m_0})e$ , the wildcard  $\mathbf{X}$  is restricted to zero or more repeated occurrences of any letter other than  $e$  or  $i$ . When  $n$  repetitions of the same letter character are intended, we write  $\hat{x}_{\times n}$ . Thus  $\hat{x}_{\times 2}$  refers to double letters, which contrasts with  $\hat{x}\hat{x}$  (two letter characters, either identical or distinct).

*Example 3.4.3.* In Finnish, there are word endings in the form of  $(ahan|ehen|ihin|ohon|ähän|öön)$ . Such patterns can be represented by  $\hat{x}\hat{h}\hat{e}\hat{n}$  for  $\hat{x}^e(a|e|i|o|ä|ö)$ .

*Example 3.4.4.* The string pattern  $\mathbf{CVCe}$  in English, when  $\mathbf{V}$  and  $\mathbf{C}$  are appropriately defined, usually implies that the  $\mathbf{V}$  in question is pronounced as a long vowel.

*Example 3.4.5.* We write  $\hat{\alpha} = \hat{\beta}$  if the two strings in question match exactly. With string alternatives and wildcard symbols, we can interpret many other string equations and string inequalities. For example,  $\hat{\beta} = \hat{\alpha}\mathbf{X}$  means that string  $\hat{\beta}$  can be constructed from string  $\hat{\alpha}$ , together with zero or more trailing characters (which is usually a linguistically meaningful scenario, where  $\hat{\beta}$  is a morphological derivative of  $\hat{\alpha}$ ); the equation  $\hat{\beta} = \hat{\alpha}(\hat{\sigma}_1 | \cdots | \hat{\sigma}_n)$  means that at least one of the  $n$  statements  $\hat{\beta} = \hat{\alpha}\hat{\sigma}_1, \dots, \hat{\beta} = \hat{\alpha}\hat{\sigma}_n$  holds (where  $\hat{\sigma}_1, \dots, \hat{\sigma}_n$  can be regarded as  $n$  candidate suffixes); the string inequality  $\hat{\beta} \neq \hat{\alpha}(\hat{\sigma}_1 | \cdots | \hat{\sigma}_n)$  means that none of the  $n$  statements  $\hat{\beta} = \hat{\alpha}\hat{\sigma}_1, \dots, \hat{\beta} = \hat{\alpha}\hat{\sigma}_n$  is true.

**Definition 3.5** (Partial and full matches). A string  $\hat{\sigma}$  written in the plain may occur anywhere in a word. The underlined notation  $\underline{\hat{\sigma}}$  refers to a string that matches the entire word exactly.

The notation  $\hat{\sigma} \sim$  (resp.  $\sim \hat{\sigma}$ ) refers to any word that “starts” (resp. “ends”) with the string  $\hat{\sigma}$ . Here, an isolated word may “start” or “end” at external and internal word boundaries. The external word boundaries are the two ends of the text string representing the isolated word, while the internal word boundaries occur at any non-letter character (apostrophe, digit, hyphen, etc.) in the string.<sup>22</sup>  $\square$

**Definition 3.6** (Substitution rules). When we say “do  $\hat{\sigma} \rightarrow \hat{\sigma}'$ ”, we mean to replace any occurrence(s) of the string  $\hat{\sigma}$  in a word by a new string  $\hat{\sigma}'$ . Here in the substitution rule, the string  $\hat{\sigma}$  may also be augmented with wildcard symbols and partial/full match notations. The letters behind the wildcard symbols are left intact in substitutions, if  $\hat{\sigma}'$  contains the respective wildcard symbols.

When we say “do  $\hat{\sigma}_1 \rightarrow \hat{\sigma}'_1, \dots, \hat{\sigma}_n \rightarrow \hat{\sigma}'_n$ ” on a particular word, we mean to perform  $n$  types of substitution tasks in a single sweep — there might be multiple replacements in a single word, but the replacements do not overlap each other.

Instead of saying “do  $\hat{\sigma}_1 \rightarrow \hat{\sigma}'_1, \dots, \hat{\sigma}_n \rightarrow \hat{\sigma}'_n$ ”, we may also tabulate the substitution rules as follows:

$\hat{\sigma}_1$	$\cdots$	$\hat{\sigma}_n$
$\hat{\sigma}'_1$	$\cdots$	$\hat{\sigma}'_n$

<sup>20</sup>Roughly speaking, this corresponds to `LetterCharacter` in *Mathematica* for versions 10.3 and higher. Earlier versions of *Mathematica* did not treat some symbols in non-Latin writing systems as `LetterCharacter`.

<sup>21</sup>In general, we avoid applying the overline directly to letters, such as  $\bar{e}$ , because this may be confused with the macron diacritic (as in the long vowel markings  $\bar{a}, \bar{e}, \bar{i}, \bar{o}, \bar{u}, \bar{y}$  employed in Latin dictionaries). When  $\hat{x}_{\notin}(e)$  is intended, we write  $\hat{x}^{\infty}(e)$ .

<sup>22</sup>Both external and internal word boundaries are represented by `WordBoundary` in *Mathematica*.

where the strings in the shaded entries are replaced by their counterparts that lie immediately below. The tabulation is particularly helpful when  $n$  is large and/or some structures in the substitution rules are better visualized in table form.  $\square$

*Example 3.6.1.* Application of the substitution rule “do  $\sim s \rightarrow \emptyset$ ” will remove the plural markers in the word *daughters* as well as the compound *daughters-in-law*. If we “do  $r \sim R$ ” on the word *rock-and-roll*, we obtain *Rock-and-Roll*. If we “do  $ass \rightarrow \emptyset$ ,  $sas \rightarrow \emptyset$ ” on the word *assassinate*, we obtain *inate*, as a result of *assassinate*  $\xrightarrow{\text{remove } sas}$  *assinate*  $\xrightarrow{\text{remove } ass}$  *inate*.<sup>23</sup>

*Example 3.6.2.* The following tabulated substitution rules

$\sim(go(es ing ne) went)$	$\frac{\text{taught}}{\text{teach}}$	$\frac{\mathbf{x} \in (\emptyset   in   mis   off   out   over   p   under) led}{X lead}$
----------------------------	--------------------------------------	---

regularize some verb forms in English.

**Definition 3.7** (NW and SW). For two text strings  $\hat{\sigma}_1$  and  $\hat{\sigma}_2$ , the output of the function  $\text{NW}(\hat{\sigma}_1, \hat{\sigma}_2)$  [resp.  $\text{SW}(\hat{\sigma}_1, \hat{\sigma}_2)$ ] is a string pattern that results from the sequence alignment from  $\hat{\sigma}_1$  to  $\hat{\sigma}_2$  according to the Needleman–Wunsch (resp. Smith–Waterman) algorithm.  $\square$

*Example 3.7.1.*  $\text{NW}(keep, kept) = \text{SW}(keep, kept) = k[e, \emptyset]ep[\emptyset, t]$ , where the bracketed portions show mismatches of the two strings. Here, in the presentation of the output string pattern, we have used  $[e, \emptyset]$  and  $[\emptyset, t]$  to show the respective contributions from the string  $\hat{\sigma}_1 = keep$  and  $\hat{\sigma}_2 = kept$ . Unlike the situation in notation for “alternatives”  $(\hat{\sigma}'_1 | \dots | \hat{\sigma}'_n)$ , the arguments within the square brackets no longer commute, and must reflect the order in the input arguments of the functions  $\text{NW}(\hat{\sigma}_1, \hat{\sigma}_2)$  and  $\text{SW}(\hat{\sigma}_1, \hat{\sigma}_2)$ .

*Example 3.7.2.* The outputs of the functions  $\text{NW}(\hat{\sigma}_1, \hat{\sigma}_2)$  and  $\text{SW}(\hat{\sigma}_1, \hat{\sigma}_2)$  may not necessarily be identical. For example,  $\text{NW}(infer, inferior) = infer[\emptyset, rio]r$  while  $\text{SW}(infer, inferior) = infer[\emptyset, ior]$ .

## 3.2 Normalizations of Wikipedia pages and WikiQA dataset

Each question in the WikiQA dataset [44] is associated with a unique Wikipedia page. We need to pick a sentence from a small subset (usually just the opening paragraph, as provided by the WikiQA dataset) of this Wikipedia page that best answers the question. In the main text of this research, we use the single Wikipedia page (in full, not just the opening paragraph) as the sole training source of our artificial neural network, before we score and screen individual sentences in the candidate answer pool.

To construct the knowledge base used in this work, we downloaded the static dump of English Wikipedia dated Sept. 1, 2015 (file name: `enwiki-20150901-pages-articles-xml.bz2`) from the URL [https://meta.wikimedia.org/wiki/Data\\_dump\\_torrents#enwiki](https://meta.wikimedia.org/wiki/Data_dump_torrents#enwiki) and saved it as an XML file. We chose this static dump because its time stamp is close to the publication date of the WikiQA dataset [44].

Once we have located the exact Wikipedia page suitable for a particular WikiQA task, we need to process the raw XML code of such a page, to generate (nearly human-readable) plain text strings. This normalization procedure is described below.

**Algorithm 3.8** (Normalization of English Wikipedia Articles). *Let  $\mathbf{D} = (0|1|2|3|4|5|6|7|8|9)$  be any one of the digit characters, and  $\mathbf{W} = (\hat{\chi}|\mathbf{D})$  be either a letter character or a digit character. The XML codes of a page from the filtered knowledge base are normalized by sequential applications of the following substitution rules, where all the boxed portions are spelt out in verbatim mode.*

- (1) Do  $\boxed{\&amp;}\rightarrow\&$ .
- (2) Replace the shaded patterns by their respective counterparts that lie immediately below:

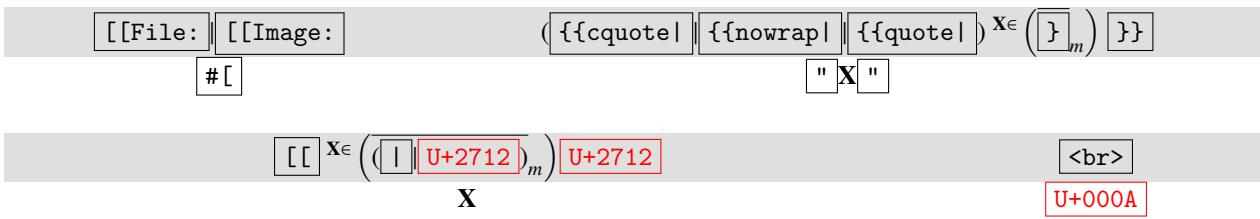
$\&nbs;$	$\&lt;$	$\&gt;$	$\&quot;$	]]
$\emptyset$	$<$	$>$	$" = \textcolor{red}{U+0022}$	$\textcolor{red}{U+2712}$

- (3) Replace

<sup>23</sup>Written in *Mathematica* codes, these statements amount to the following results:

```
StringReplace["daughters", "s"~~WordBoundary :> ""] yields "daughter";
StringReplace["daughters-in-law", "s"~~WordBoundary :> ""] yields "daughter-in-law";
StringReplace["rock-and-roll", WordBoundary~~"r" :> "R"] yields "Rock-and-Roll";
StringReplace["assassinate", {"ass" :> "", "sas" :> ""}] yields "inate".
```

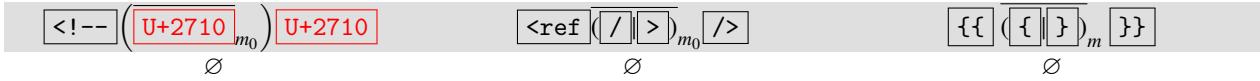
Note in particular that in the last example, the substitution rules  $\{"ass" :> "", "sas" :> "\"} and  $\{"sas" :> "", "ass" :> "\"}$  will generate the same output.$



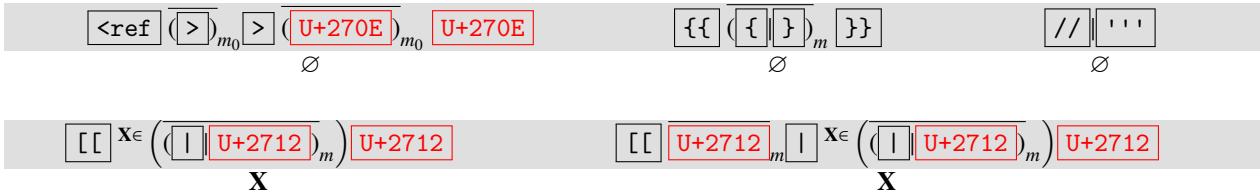
(4) Replace



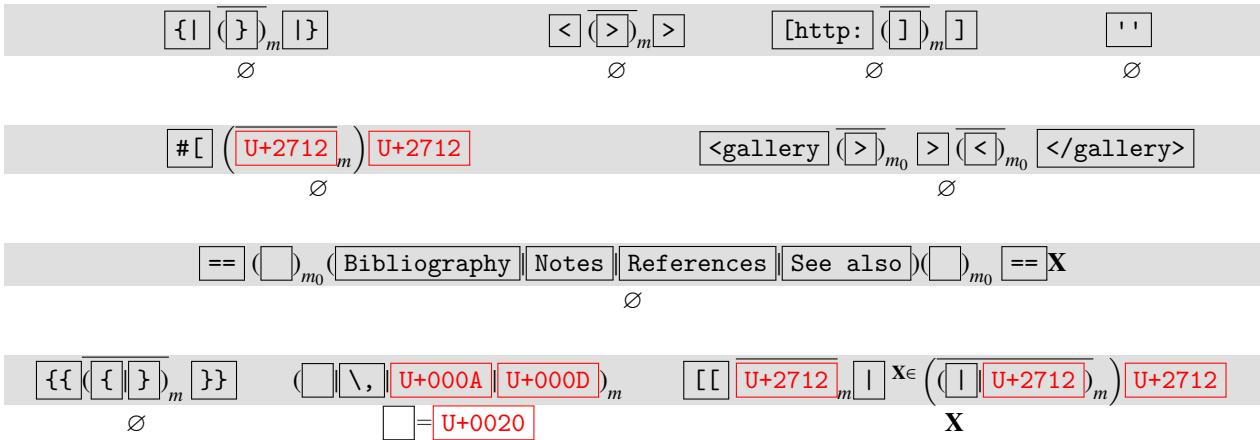
(5) Replace



(6) Replace



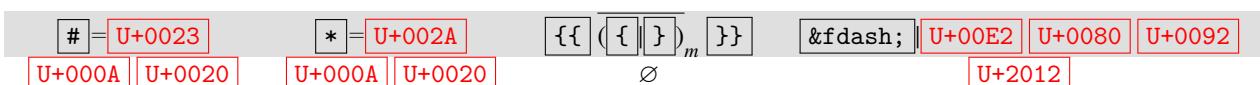
where  $\boxed{/\!/}$  denotes two consecutive appearances of SOLIDUS  $\boxed{U+002F}$  and  $\boxed{''''}$  denotes three consecutive appearances of APOSTROPHE  $\boxed{U+0027}$ .

(7) Replace<sup>24</sup>

and also delete the string pattern  $\hat{T}\text{X}$   $\text{text/x-wiki}$ , where  $\hat{T}$  is the text string matching the title of the Wikipedia page.

(8) Delete the string pattern  $\boxed{{\{} {\{} {\}}_m }}$ .

(9) Replace



<sup>24</sup>Here, the string pattern  $\boxed{==(\>)_{m_0}}$   $(\boxed{\text{Bibliography}}|\boxed{\text{Notes}}|\boxed{\text{References}}|\boxed{\text{See also}})(\>)_{m_0}$   $\boxed{==}$  must be deleted along with anything thereafter, because the wildcard symbol  $\text{X}$  does not reappear in the destination string.

$(\boxed{-} = \text{U+2013}) \& \text{ndash}; (\boxed{\text{U+00E2}} \text{U+0080} \text{U+0093})$	$\& \text{mdash}; (\boxed{\text{U+00E2}} \text{U+0080} \text{U+0094})$
$\boxed{--} = \text{U+002D} \times 2$	$\boxed{—} = \text{U+2014}$

(10) Do  $\overline{\mathbf{W}} \text{U+0027} \rightarrow \boxed{\phantom{X}}$ .

(11) If the text string so far ends with  $\mathbf{W}$ , then append  $\boxed{\phantom{X}}$  to it; otherwise, leave it as is.

(12) Call the result so far as string  $P_0$ . Denote the input WikiQA question by  $Q$  and the Wikipedia page title by  $T$  (both of which must be “corrected” by Algorithm 3.11 below). If  $P_0 = \emptyset$ , then send  $P = \boxed{\phantom{X}}$  as output; otherwise send  $P = Q \boxed{\text{U+000A}} \times 2 T \boxed{\text{U+000A}} \times 2 P_0 \boxed{\text{U+000A}} \times 2 Q \boxed{\phantom{X}}$  as output.

Here, in the algorithm above, we have tried our best to eliminate structured data (infoboxes, references, as well as hyperlinks to other pages and external files) from the XML file, and restore all clickable links to normal text. Some replacements are done in multiple sweeps to accommodate to nested structures in Wikipedia XML files. However, we note that due to the versatility of Wikipedia XML files, our cleansing algorithm may still have some blind spots. In practice, this heuristic algorithm handles most Wikipedia pages well, so long as there are no special characters in hyperlinks and there is no extensive coverage of TeXnical content.

The WikiQA dataset contains some Wikipedia pages that are simply lists or tables. If we process such Wikipedia pages with Algorithm 3.8, then we will end up with (almost) an empty string. To compensate for this, we generate backup versions of Wikipedia pages, using the algorithm below.

**Algorithm 3.9** (Backup version of Wikipedia pages). Define  $B = Q \boxed{\text{U+000A}} \times 2 T \boxed{\text{U+000A}} \times 2 A \boxed{\text{U+000A}} \times 2 Q \boxed{\phantom{X}}$ , where  $A$  is the concatenation of all the candidate answers ( $\approx$  the opening section of the accompanying Wikipedia page, by design of the WikiQA dataset), riffled by  $\boxed{\phantom{X}}$ . If  $\ell(P) < \ell(B)$ , then redefine  $P$  as the backup version  $B$ . If  $T$  matches List of  $X$ , regardless of letter case, then redefine  $P$  as  $B \boxed{\phantom{X}} B \boxed{\phantom{X}} B$ .

After processing a Wikipedia page by Algorithms 3.8 and 3.9, we arrive at a text string  $P$ , which will be suitable for sentence extractions later. To facilitate the measurement of gaps between words, we need to regularize spacing, punctuation and capitalization, as indicated below.

**Algorithm 3.10** (Regularisation of Text String). Let  $D = (0|1|2|3|4|5|6|7|8|9)$  be any one of the digit characters, and  $\mathbf{W} = (\hat{\chi}|D)$  be either a letter character or a digit character. A processed Wikipedia page must be regularized by the following substitution rules before temporal statistics are computed.

(1) Replace

$(\boxed{\phantom{X}} \boxed{\text{U+000A}})_m$	$= (\boxed{=})_m$	$\dots = \boxed{\text{U+002E}} \times 3$	$\boxed{\phantom{X}} \in ((\hat{\chi}   \boxed{-})_m)$
$\boxed{\phantom{X}} = \text{U+0020}$	$\boxed{\text{U+002E}} \boxed{\text{U+000A}}$	$\boxed{\phantom{X}} \times 3 = \boxed{\text{U+0020}} \times 3$	$\mathbf{X}$ if $\ell(\mathbf{X}) = 1$ ; $\mathbf{X}^{(1)}(\mathbf{X}^{[2,\ell(\mathbf{X})]})_-$ if $\ell(\mathbf{X}) > 1$ .
$\boxed{\cdot} D = \boxed{\text{U+002E}} D$	$\boxed{,} D = \boxed{\text{U+002C}} D$	$\boxed{:} D = \boxed{\text{U+003A}} D$	
$\boxed{\cdot} D = \boxed{\text{U+00B7}} D$	$\boxed{D}$		$\boxed{h} D$

where  $\mathbf{X}^{(1)}(\mathbf{X}^{[2,\ell(\mathbf{X})]})_-$  reduces  $\mathbf{X}$  to lower case except the initial character.<sup>25</sup>

(2) Replace

$\boxed{,} \boxed{;} \boxed{:} \boxed{()} \boxed{\text{U+000A}} \boxed{\text{U+2012}}$	$\boxed{\cdot} \boxed{?} \boxed{!} \boxed{''}$	$\boxed{\text{U+2014}} \boxed{(\text{U+2013})} \boxed{\text{U+002D}}$
$\boxed{\phantom{X}} = \text{U+0020}$	$\boxed{\phantom{X}} \times 2 = \text{U+0020} \times 2$	$\boxed{\phantom{X}} \times 3 = \text{U+0020} \times 3$

(3) Prepend a single white space character  $\boxed{\text{U+0020}}$  to the entire text string.

(4) Record the positions of all the  $\boxed{\text{U+0020}}_m$  patterns.

(5) Split the text string at  $\boxed{\text{U+0020}}_m$ , to obtain an ordered list of separate words.

Upon conclusion of the aforementioned algorithm, we have effectively obtained the start and end positions of every word appearing in the normalized Wikipedia page  $P$ .

After scrutinizing the WikiQA dataset, we have detected a few typographical flaws:

<sup>25</sup>This means that the capitalization status of a word is determined by its initial character: eBay → ebay, ELIZABETH → Elizabeth, McDonald → McDonald etc.

- A few Wikipedia page names listed in WikiQA do not match our record.
- There are non-standard uses of capitalizations and diacritics.
- There are non-standard abbreviations that cannot be inferred from the related Wikipedia pages.
- There are other kinds of misspellings (including unnecessary use of foreign language) in the questions that do not match the information available from the related Wikipedia pages.

Since our purpose in the current research is to test a numerical question answering machine, rather than an automated spell checker, we decide to normalize the WikiQA dataset according to the following *ad hoc* procedures.

**Algorithm 3.11** (Normalization of WikiQA page names and questions). *We normalize the reference Wikipedia page names by the following replacements:*

Automatic Document Feeder	High-Sticking	Judas (song)	Julius caesar
Automatic document feeder	High-sticking	Judas (Lady Gaga song)	Julius Caesar
<i>What It Takes (song)</i>			World Trade Center
<i>What It Takes (Aerosmith song)</i>		World Trade Center (2001 <span style="border: 1px solid red; padding: 2px;">U+00E2</span> <span style="border: 1px solid red; padding: 2px;">U+0080</span> <span style="border: 1px solid red; padding: 2px;">U+0093</span> present)	

We use the following heuristics to determine the correct capitalizations and diacritics in a question Q, by checking it against all the candidate answers ( $\approx$  the opening section of the accompanying Wikipedia page, by design of the WikiQA dataset) A:

- (1) Take the union of all the words in A, and call it  $\mathcal{A}$ . Remove diacritics from  $\mathcal{A}$ , and call it  $\mathcal{A}'$ . Use set-theoretic complements to define  $\mathcal{D} = \mathcal{A}' \setminus \mathcal{A}$  and  $\mathcal{D}_0 = \mathcal{A} \setminus \mathcal{A}'$ .<sup>26</sup> Define  $\mathcal{A}_-$  (resp.  $\mathcal{D}_-$ ) as the lowercase form of  $\mathcal{A}$  (resp.  $\mathcal{D}$ ).
- (2) For each word in Q, convert it to lower case, remove diacritics if any, call the result  $\hat{q}^*$ , and then do the following:
  - If  $\hat{q}^*$  matches  $\mathcal{A}_-(\emptyset|s|'s|s')$  and does not match  $\mathcal{A}(\emptyset|s|'s|s')$ , then capitalize  $\hat{q}^*$  and quit. Otherwise, go to next step.
  - If  $\hat{q}^*$  matches  $\mathcal{D}_-(\emptyset|s|'s|s')$ , landing on the n-th member of alphabetized  $\mathcal{D}$  upon its first hit, then determine the correct spelling by the n-th member of alphabetized  $\mathcal{D}_0$ . Otherwise, leave  $\hat{q}^*$  as is and quit.

Furthermore, we correct misspelt questions by the following replacements:

a full job time	cono sur	Điện Chí Minh	<span style="border: 1px solid black; padding: 2px;">fy</span> <span style="border: 1px solid black; padding: 2px;">fiscal year</span>	general chu chicken
a full-time job	Southern Cone	Hồ Chí Minh	<span style="border: 1px solid black; padding: 2px;">fiscal year</span>	General Tso's Chicken
Google in math	jagger bomb	st patty	ti 82	trinity 5 7
googol in math	Jägerbomb	Saint Patrick	TI-82	Trin-i-tee 5:7
tri tip	<span style="border: 1px solid black; padding: 2px;">ua's</span>	<span style="border: 1px solid black; padding: 2px;">ww</span>	<span style="border: 1px solid black; padding: 2px;">wwii</span>	
tri-tip	<span style="border: 1px solid black; padding: 2px;">urinalysis</span>	<span style="border: 1px solid black; padding: 2px;">World War</span>	<span style="border: 1px solid black; padding: 2px;">World War II</span>	

In Table 2 and Fig. S2, we showed that our numerical question answering machine had better performance on a subset of the WikiQA questions. Our criteria for defining such a subset (via its complement) are stated below.

**Algorithm 3.12** (Heuristic classification of WikiQA questions). *We classify a WikiQA question as “quantitative” (requiring rule-based reasoning), if it matches one of the following string patterns (regardless of letter case):*

how( )<sub>m</sub>(big|deep|far|frequent|high|long|many|much|often|old|tall)X,  
whenX,  
X( )#|number|percentage|proportion)X,  
what( )(bec|happen|year)X.

Otherwise, we classify a WikiQA question as “qualitative” (requiring only associative reasoning).

<sup>26</sup>Note that if the reference page uses diacritics consistently and we alphabetize both  $\mathcal{D}$  and  $\mathcal{D}_0$ , then the corresponding positions are occupied by the same words, without and with diacritics. For example, set A = *Estée Lauder* and *Eva Perón* and Q = *péron esTee lauder*, then we have  $\mathcal{D} = \{\text{Estee, Peron}\}$  and  $\mathcal{D}_0 = \{\text{Estée, Perón}\}$ . After going through the remaining procedures in Step (2), we will end up with a corrected version of Q as *Perón Estée Lauder*. Our heuristic algorithm will malfunction, however, if P contains both *Estée* and *Estee* (i.e. the reference page uses diacritics inconsistently), or if alphabetic ordering is seriously affected by diacritics. In practice, the only glitch occurs in WikiQA-Q2158, where *Ho Chí Minh* becomes misspelt as *Điện Chí Minh*, due to unexpected alphabetic ordering of Vietnamese characters in *Mathematica* v11.3 (which brought us  $\mathcal{D} = \{\text{Ai, Bien, Chi, Cong, Dien, Ho, Nguyen, Phu, Quoc, Tat, Thanh, Viet}\}$  and  $\mathcal{D}_0 = \{\text{Ái, Biên, Chí, Công, Hồ, Điện, Nguyễn, Phú, Quốc, Tất, Thành, Việt}\}$ ). This is subsequently corrected by a substitution rule.

### 3.3 Normalizations of ebooks for text mining

For the topic extraction and word translation experiments (Figs. S3–S16), we use the following algorithms for cleansing electronic books written in 12 European languages and 2 Asian languages, sorted in the same order as the remaining sections of this document. Hereafter, in accordance with *Mathematica* terminology, we use the word character  $\mathbf{W} = (\hat{\chi}|D)$  to refer to either a letter or a digit.

**Algorithm 3.13** (Text cleansing for English ebooks). *We process our text (ebooks from Gutenberg Project<sup>27</sup> and elsewhere) according to the following procedures:*

(1) Do  $\mathbf{W} \rightarrow \square_{\times 2}$ .

(2) Replace

$$(\square_m | \textcolor{red}{U+000A}_m | \textcolor{red}{--} | * | \square | ( | ) |) \xrightarrow{\square} \begin{array}{c} \mathbf{x} \in ((\hat{\chi} | \square | \square)_m) \\ \mathbf{X} \text{ if } \ell(\mathbf{X}) = 1; \mathbf{X}^{(1)}(\mathbf{X}^{[2, \ell(\mathbf{X})]})_- \text{ if } \ell(\mathbf{X}) > 1. \end{array}$$

where  $\mathbf{X}^{(1)}(\mathbf{X}^{[2, \ell(\mathbf{X})]})_-$  reduces  $\mathbf{X}$  to lower case except the initial character.

**Algorithm 3.14** (Text cleansing for Danish ebooks). *We normalize our text by doing*

$$(\square_m | \textcolor{red}{U+000A}_m | \textcolor{red}{--} | * | \square | \textcolor{red}{U+002D} | \textcolor{red}{U+2013} | \textcolor{red}{U+2014} | \textcolor{red}{U+201E} | ( | ) | \textcolor{red}{U+00AB} | \textcolor{red}{U+00BB} | \textcolor{red}{U+2026}) \rightarrow \square.$$

**Algorithm 3.15** (Text cleansing for German ebooks). *We normalize our text by doing*  $(\textcolor{red}{U+0027} | \textcolor{red}{U+2019}) \rightarrow \square$ ,  $(\textcolor{red}{U+2039} | \textcolor{red}{U+203A}) \rightarrow \square$ ,  $(\square_m | \textcolor{red}{U+000A}_m | \textcolor{red}{--} | * | \square | \textcolor{red}{U+002D} | \textcolor{red}{U+2013} | \textcolor{red}{U+2014} | \textcolor{red}{U+201E} | ( | ) | \textcolor{red}{U+00AB} | \textcolor{red}{U+00BB} | \textcolor{red}{U+2026}) \rightarrow \square$ ,  $\text{u.s.w.} \rightarrow \square_{\times 6}$ .

**Algorithm 3.16** (Text cleansing for Dutch ebooks). *We normalize our text by doing*  $(\textcolor{red}{U+0027} | \textcolor{red}{U+2019}) \rightarrow \square$ ,  $(\textcolor{red}{U+2039} | \textcolor{red}{U+203A}) \rightarrow \square$ ,  $(\square_m | \textcolor{red}{U+000A}_m | \textcolor{red}{--} | * | \square | \textcolor{red}{U+002D} | \textcolor{red}{U+2013} | \textcolor{red}{U+2014} | \textcolor{red}{U+201E} | ( | ) | \textcolor{red}{U+00AB} | \textcolor{red}{U+00BB} | \textcolor{red}{U+2026}) \rightarrow \square$ .

**Algorithm 3.17** (Text cleansing for Spanish ebooks). *We normalize our text by doing*  $(\textcolor{red}{\dot{z}} | \textcolor{red}{i}) \rightarrow \emptyset$ ,  $(\square_m | \textcolor{red}{U+000A}_m | \textcolor{red}{--} | * | \square | \textcolor{red}{U+002D} | \textcolor{red}{U+2013} | \textcolor{red}{U+2014} | \textcolor{red}{U+201E} | ( | ) | \textcolor{red}{U+00AB} | \textcolor{red}{U+00BB} | \textcolor{red}{U+2026}) \rightarrow \square$ .

**Algorithm 3.18** (Text cleansing for French ebooks). *We normalize our text (ebooks from Wikisource and Gutenberg Project) according to the following procedures:*

(1) Do  $M. \rightarrow Mr$ .

(2) Do  $(\square_m | \textcolor{red}{U+000A}_m | \textcolor{red}{--} | * | \square | \textcolor{red}{U+002D} | \textcolor{red}{U+2013} | \textcolor{red}{U+2014} | \textcolor{red}{U+201E} | ( | ) | \textcolor{red}{U+00AB} | \textcolor{red}{U+00BB} | \textcolor{red}{U+2026}) \rightarrow \square$ ,  $(\textcolor{red}{U+0027} | \textcolor{red}{U+2019}) \rightarrow \square$ .

**Algorithm 3.19** (Text cleansing for Latin ebooks). *We normalize our text by doing*

$$(\square_m | \textcolor{red}{U+000A}_m | \textcolor{red}{--} | * | \square | \textcolor{red}{U+002D} | \textcolor{red}{U+2013} | \textcolor{red}{U+2014} | \textcolor{red}{U+201E} | ( | ) | \textcolor{red}{U+00AB} | \textcolor{red}{U+00BB} | \textcolor{red}{U+2026}) \rightarrow \square.$$

**Algorithm 3.20** (Text cleansing for Polish ebooks). *We normalize our text according to the following procedures:*

(1) Do  $(\emptyset | \textcolor{red}{U+002D} | \textcolor{red}{U+2014})(\textcolor{red}{U+0020} | \textcolor{red}{U+000A})_{m_0} | \mathbf{D}_{m_0} | (\textcolor{red}{U+0020} | \textcolor{red}{U+000A})_{m_0} \rightarrow \emptyset$ ,  $\hat{\chi}_1(\textcolor{red}{U+0027} | \textcolor{red}{U+2019}) \hat{\chi}_2 \rightarrow \hat{\chi}_1 \# \hat{\chi}_2$ .

(2) Do  $(\square_m | \textcolor{red}{U+000A}_m | \textcolor{red}{--} | * | \square | \textcolor{red}{U+002D} | \textcolor{red}{U+2013} | \textcolor{red}{U+2014} | \textcolor{red}{U+201E} | ( | ) | \textcolor{red}{U+00AB} | \textcolor{red}{U+00BB} | \textcolor{red}{U+2026}) \rightarrow \square$ ,  $(\textcolor{red}{U+0027} | \textcolor{red}{U+2019}) \rightarrow \square$ .

**Algorithm 3.21** (Text cleansing for Russian ebooks). *We normalize our text according to the following procedures:*

(1) Do  $(\emptyset | \textcolor{red}{U+002D} | \textcolor{red}{U+2014})(\textcolor{red}{U+0020} | \textcolor{red}{U+000A})_{m_0} | \overline{[\square]}_{m_0} | (\textcolor{red}{U+0020} | \textcolor{red}{U+000A})_{m_0} \rightarrow \emptyset$ ,  $\text{BookTitle}(\square | \square)_{m_0} | \mathbf{D}_{m_0} | \mathbf{D}_{m_0} | (\square | \square)_{m_0} \rightarrow \square$ .

(2) Do  $\hat{\chi}_1 | ? | \hat{\chi}_2 \rightarrow \hat{\chi}_1 y \hat{\chi}_2$ .

<sup>27</sup>It is always understood that the front and back matter (publisher's boilerplate, table of contents, introductory chapter, glossary, index, endnotes, etc.) in ebooks are manually removed. Footnotes that scatter in a few pages are left intact.

(3) Do  $\left( \square_m \text{U+000A}_m \square_{--} \square_* \square_{\_} \text{U+002D} \text{U+2013} \text{U+2014} \text{U+201E} (\square) \text{U+00AB} \text{U+00BB} \text{U+2026} \right) \rightarrow \square$ ,  
 $(\text{U+0027} \text{U+2019}) \rightarrow \square$

**Algorithm 3.22** (Text cleansing for Finnish ebooks). *We normalize our text (ebooks from the Gutenberg Project) according to the following procedures:*

(1) Do  $(\boxed{\quad}_m | \text{U+000A}_m | \boxed{-} | \boxed{*} | \boxed{/} | \text{U+002D}_m | \text{U+2013}_m | \text{U+2014}_m | \text{U+201E}_m | (\quad) | \text{U+00AB}_m | \text{U+00BB}_m | \text{U+2026}_m) \rightarrow \boxed{\quad}$ ,  
 $(\text{U+201C}_m | \text{U+201D}_m) \rightarrow \boxed{\quad}, [\mathbf{D}_{m_0}] \rightarrow \emptyset, \mathbf{W} : \hat{x} \rightarrow \mathbf{W} \# \hat{x}$ .

(2) Do **U+0027** **W** → , **W** **U+0027** → .

**Algorithm 3.23** (Text cleansing for Hungarian ebooks). *We normalize our text (ebooks from <http://mek.oszk.hu>) according to the following procedures:*

(1) Do  $\langle \emptyset | \text{U+002D} | \text{U+2014} \rangle \langle \text{U+0020} | \text{U+000A} \rangle_{m_0} [\mathbf{D}_{m_0}] \langle \text{U+0020} | \text{U+000A} \rangle_{m_0} \rightarrow \emptyset,$   
 $\text{BookTitle}(\boxed{\phantom{a}} - )_{m_0} \mathbf{D}_{m_0} / [\mathbf{D}_{m_0} (\boxed{\phantom{a}} - )_{m_0} \rightarrow \boxed{\phantom{a}}].$

(2) Do  $(\boxed{\quad}_m \text{U+000A}_m \boxed{-} \text{*} \boxed{_} \text{U+002D} \text{U+2013} \text{U+2014} \text{U+201E} (\boxed{}) \text{U+00AB} \text{U+00BB} \text{U+2026}) \rightarrow \boxed{\quad}$ ,  
 $(\text{U+0027} \text{U+2019} \text{U+201C} \text{U+201D}) \rightarrow \boxed{\quad}$ .

**Algorithm 3.24** (Text cleansing for Basque ebooks). *We normalize our text according to the following procedures:*

$$(1) \text{ } Do (\emptyset \boxed{\text{U+002D}} \boxed{\text{U+2014}})(\boxed{\text{U+0020}} \boxed{\text{U+000A}})_{m_0} [\boxed{\mathbf{D}_{m_0}}] (\boxed{\text{U+0020}} \boxed{\text{U+000A}})_{m_0} \rightarrow \emptyset.$$

(2) *Do eliz* → *ëliz*.<sup>28</sup>

(3) Do  $(\boxed{\quad}_m | \text{U+000A}_m | - | * | _ | \text{U+002D}_m | \text{U+2013}_m | \text{U+2014}_m | \text{U+201E}_m | (\quad) | \text{U+00AB}_m | \text{U+00BB}_m | \text{U+2026}_m) \rightarrow \boxed{\quad}$ ,  
 $(\text{U+0027}_m | \text{U+2018}_m | \text{U+2019}_m | \text{U+201C}_m | \text{U+201D}_m) \rightarrow \boxed{\quad}$ .

**Algorithm 3.25** (Text cleansing for Korean ebooks). *We normalize our text according to the following procedures:*

(1) *Do* ( $\emptyset$  | U+002D | U+2014) (U+0020 | U+000A) $_{m_0}$  [ D $_{m_0}$  ] (U+0020 | U+000A) $_{m_0}$   $\rightarrow \emptyset$ .

(2) Do 제  $\mathbf{D}_m(\text{권}|\text{장}) \rightarrow$  □.

(3) Do  $(\boxed{\quad}_m \text{U+000A}_m \text{--} * \boxed{\quad} \text{U+002D} \text{U+2013} \text{U+2014} \text{U+201E} (\ ) \text{U+00AB} \text{U+00BB} \text{U+2026}) \rightarrow \boxed{\quad}$ ,  
 $(\text{U+0027} \text{U+2018} \text{U+2019} \text{U+201C} \text{U+201D} \text{U+2500} \text{U+25CB} \text{U+3008} \text{U+3009} \text{U+300C} \text{U+300D}) \rightarrow \boxed{\quad}$ .

(4) Adjust spacing in the text as follows:<sup>29</sup>

(4.1) **Do** **D**백 → **D**□백, **D**천 → **D**□천, **D**만 → **D**0□천, **D**일 → **D**□, 지난W → 지난□W, 다음W → 다음□W,  
□다시□<sub>속</sub>(씨|양) → □<sub>속</sub>ㄔ(씨|양).

신나십|나자|나지|납니|납시|났|보(으)|본|볼|봄|봄|봐|봤(자))→□ 꾸Ӯ, 고□ Ӯ(말|실)→□ 고Ӯ.  
 $\hat{x}^{\epsilon}$ (과|앤|와) □ 맙 췄 □ → □ 다시 □, 까□  $X^{\epsilon}$ (말|보(으)|본|볼|봄|봄|봐|봤(자)|실)→□ 까X.

$\hat{\chi}^{\hat{1}\epsilon}(\text{다|아|어}) \square \hat{\chi}^{2\epsilon}(\text{보|으|}|\text{본|볼|봄|봄|봐|봤|자}) \rightarrow \hat{\chi} \square \hat{\chi}_1 \hat{\chi}_2$ ,       $\hat{\chi}_1 \epsilon (\text{난|단|란|잔}) \square \text{말} \hat{\chi}^{2\epsilon}(\text{씀|이}) \rightarrow \square \hat{\chi}_1 \text{말} \hat{\chi}_2$ ,  
 지  $\square X^\epsilon(\text{마느|마세|마셨|마시|마십|말겠|말더|말례|말았|말지|맙니|맙시}) \rightarrow \square \text{지} X$ ,      지  $\square (\sim X^\epsilon(\text{마는|마니|마})$

서|마신|마실|마심|만|말|말고|말기|말면|말아|말자|맑))→□지**X**,

**Algorithm 3.26** (Text cleansing for Turkish ebooks). *We normalize our text according to the following procedures:*

(1) *Do*( $\emptyset$ )  $\boxed{\text{U+002D}}$   $\boxed{\text{U+2014}}$  ( $\boxed{\text{U+0020}}$   $\boxed{\text{U+000A}}$ ) $_{m_0}$   $\boxed{[\mathbf{D}_{m_0}]} \boxed{(\text{U+0020} \text{ } \text{U+000A})_{m_0}} \rightarrow \emptyset, \hat{\chi}_1(\boxed{\text{U+2019}}) \hat{\chi}_2 \rightarrow \hat{\chi}_1 \# \hat{\chi}_2, I \rightarrow i, I \rightarrow i, BÖLÜM(\boxed{\phantom{0}})_m \mathbf{D}_m \rightarrow BÖLÜM(\boxed{\phantom{0}}).$

(2) Do  $\left( \square_m \boxed{U+000A}_m \boxed{-} \boxed{*} \boxed{_} \boxed{U+002D} \boxed{U+2013} \boxed{U+2014} \boxed{U+201E} \boxed{(} \boxed{)} \boxed{U+00AB} \boxed{U+00BB} \boxed{U+2026} \right) \rightarrow \square$ ,  
 $\left( \boxed{U+0027} \boxed{U+2018} \boxed{U+2019} \boxed{U+201C} \boxed{U+201D} \right) \rightarrow \square$ .

<sup>28</sup>The Basque word *eliza* “church” should not be confused with the proper name *Eliza*.

<sup>29</sup>Some Korean words are best disambiguated by their immediate neighbours, such as ~고□설~ “want to do” [45, p. 251] and ~(가)나□설~ “think it might” [45, p. 256]. This is the main reason why we reassign spaces to our text.

To obtain individual words from a normalized text, we split the text string at (□, □, □, ?, □, !, □, ;, □ U+0022) (for French, also apostrophe □' and hyphen □-). To measure gaps between words in a normalized text, we further substitute spaces and punctuation marks according to the rules: (□, □, :, ;) → □, (□, □, ?, !, □ U+0022) → □ × 2.

## Part III

# Protocols for word clustering

All the word clustering algorithms in this Part are available as *Mathematica* source codes (*English.nb* for word clustering and analysis of English texts, etc.), posted to Github

<https://github.com/yajun-zhou/linguae-naturalis-principia-mathematica>

To reproduce all the figures in this Part from these source codes, the readers need to download and/or purchase text corpora (not posted to our Github repository), following the instructions in Table S1.

Additionally, a separate source code *WikiQA.nb* is available, which implements our word clustering algorithm for English to the WikiQA dataset [44]. Based on the list of page titles in the WikiQA dataset, we extract *WikiSelecta.XML* (deposited as a zipped file *WikiSelecta.zip* in our Github folder) from the static dump of English Wikipedia dated Sept. 1, 2015 (file name: *enwiki-20150901-pages-articles-xml.bz2*). This subset of 1242 Wikipedia pages forms the only training source and knowledge base for our question-answering machine.

## 4 Approximate word clustering in English

In this section, we describe an automatic algorithm to cluster English words according to their morphologies. Due to the complicated linguistic history of English that derives from both Germanic heritage and Latin/Romance influence, the English language is rich in both inflectional complexity (such as vowel alternations in verb forms, a feature common to Germanic languages) and derivational complexity (such as many Latinate suffixes that distinguish words within the same word family). Our English word clustering algorithm thus involves methods that will be later used in the studies of Germanic (§5) and Romance (§6) languages.

### 4.1 Capitalization and stop words

Before presenting our algorithms in detail, we need to briefly discuss the capitalization problem in English and the list of English stop words.

English orthography requires capitalization of proper nouns and their derivatives.<sup>30</sup> It turns out that proper nouns and common nouns in English tend to have different etymological sources, and consequently, different morphological structures. Thus, it is helpful to single out “intrinsically capitalized words” used in an English text before performing word clustering algorithms. The precise identification of the capitalization status of a particular word is a non-trivial task: words that start any sentence are always capitalized; some authors may occasionally choose to use ALL CAPS for emphasis; some “intrinsically capitalized words” may share the same spelling with other ordinary words, up to capitalization (such as *Miss* and *miss*). For our text mining purposes, we adhere to the following heuristic determination of the capitalization status.

**Definition 4.1** (Heuristic capitalization). We say that a word is “heuristically capitalized” in a given text, if its all-lowercase form never appears in the same document. Otherwise, we consider the word’s spelling “heuristically uncapitalized”.<sup>31</sup> □

A small collection of English words perform important grammatical/logical functions, but transmit very little information on their own. These words need to be excluded in most natural language processing tasks, and are thus referred to as “stop words”. There is not a universal standard for the list of English stop words. We state our empirical list in the definition below.

<sup>30</sup> Among the major European languages, this capitalization rule perhaps only applies to English. In all the other European languages studied in this work, adjectives derived from proper nouns are spelt in lower case forms. In modern German and pre-1948 Danish, all nouns are capitalized. In Modern Turkish (which we regard as an Asian language), derivatives of proper nouns are also capitalized, as in English. However, we do not find it particularly helpful to treat Turkish proper nouns differently from other nouns in word clustering. Therefore, we will restrict Definition 4.1 to English only.

<sup>31</sup> According to this definition, a word’s spelling must be either “heuristically capitalized” or “heuristically uncapitalized”. If *Miss* and *miss* both appear in the same text, then they will be regarded as the same word, in the “heuristically uncapitalized” status of our classification scheme. On the other hand, if a text contains *Elizabeth* and *ELIZABETH*, but does not contain *elizabeth*, then we merge *Elizabeth* and *ELIZABETH* to the same word, in the “heuristically capitalized” status.

**Definition 4.2** (English stop words). If a word belongs to the following list<sup>32</sup>:

–, &c, a, about, above, across, after, afterward, afterwards, again, against, ago, ahead, albeit, all, almost, alone, along, already, also, although, always, am, amid, amidst, among, amongst, an, and, another, any, anybody, anybody's, anyhow, anyone, anything, anyway, anywhere, apart, are, aren't, around, as, aside, at, away, back, be, became, because, become, becomes, becoming, been, before, behind, being, below, beneath, beside, besides, between, beyond, both, but, by, can, can't, cannot, could, couldn't, did, didn't, do, does, doesn't, doing, don't, done, down, during, each, either, else, elsewhere, enough, etc, even, ever, every, everybody, everybody's, everyone, everyone's, everything, everywhere, except, few, for, from, front, full, further, furthermore, get, gets, getting, got, gotten, had, hadn't, hardly, has, hasn't, have, haven't, having, he, he'd, he'll, he's, hence, her, here, here's, hereabout, hereabouts, hereafter, hereby, herein, hereinafter, hereof, hereon, hereto, heretofore, hereunder, hereunto, hereupon, herewith, hers, herself, him, himself, his, hither, how, how's, however, howsoever, i, I, i'd, I'd, I'll, I'll, i'm, I'm, i've, I've, if, in, inside, instead, into, is, isn't, it, it's, its, itself, just, last, least, less, let, let's, letting, likewise, many, may, me, mere, merely, might, more, moreover, most, mostly, much, must, mustn't, my, myself, neither, never, next, no, nobody, nobody's, none, nor, not, nothing, now, nowhere, of, off, often, on, once, one, one's, ones, ones', oneself, only, or, other, others, otherwise, ought, our, ours, ourselves, out, outside, over, own, per, perhaps, quite, rather, same, several, shall, shan't, she, she'd, she'll, she's, should, shouldn't, since, so, some, somebody, somebody's, somehow, someone, someone's, something, sometimes, somewhat, somewhere, soon, sooner, still, such, than, that, that's, the, thee, their, theirs, them, themselves, then, there, there's, thereabout, thereabouts, thereafter, thereat, thereby, therefor, therefore, therefrom, therein, thereof, thereon, thereto, theretofore, thereunder, thereunto, thereupon, therewith, these, they, they'd, they'll, they're, they've, thine, this, those, thou, though, through, thus, thy, thyself, till, to, together, too, toward, towards, under, unless, until, up, upon, us, versus, very, via, was, wasn't, we, we'd, we'll, we're, we've, well, were, weren't, what, what's, whatever, whatsoever, when, whenever, when's, where, where's, whereabouts, whereas, whereat, whereby, wherefore, wherein, whereof, whereon, wheresoever, whereto, whereupon, wherever, wherewith, wherewithal, whether, which, whichever, whichever, while, who, who's, whoever, whole, whom, whomever, whomsoever, whose, whosoever, why, why's, whyever, whysoever, will, with, within, without, won't, would, wouldn't, yes, yet, you, you'd, you'll, you're, you've, your, yours, yourself, yourselves,

then we consider it an English stop word.<sup>33</sup> All the English stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list. □

## 4.2 Modified Porter stemming algorithm for English

The Porter stemming algorithm is a very useful word clustering tool for information retrieval in English [39] and many other European languages [40]. Through a series of string manipulations, the Porter stemming algorithm generates a token (known as “Porter stem”) for an arbitrary word in a given language, without consulting a dictionary. Words sharing the same token are identified with each other in the Porter stemming procedure. While this procedure correctly lumps together morphologically related words in many scenarios, it ignores morphological irregularities that prevail in the vocabularies of many European languages, including English.

In this work, we modify Porter’s stemming algorithm so that it works more accurately on realistic vocabulary lists generated from text mining tasks. Our modification has two aspects. First, instead of associating each word with a unique Porter stem, we construct the token in two stages: a longer stem called *effective spelling*, and a shorter stem called *essential root* (§4.2.1). Second, instead of using exact matching of tokens as a criterion for word clustering, we add some *admissible mutation* rules to accommodate to morphological irregularities, and build an *approximate clustering* algorithm thereupon (§4.2.2).

In English, there are about two hundred irregular verbs in common use. A limited subset of these verbs are highly irregular, and they need to be handled by explicitly written exceptions to rules during the stemming procedure (§4.2.1). The remaining

<sup>32</sup>Our list of English stop words differ from the built-in list in *Mathematica* v11.0 in the following aspects:

- (1) Our list includes – (two hyphens joined together), which accommodates to vintage encoding in certain ebooks;
- (2) Our list includes *shall* and *towards*, to accommodate to British English;
- (3) Our list includes some common function words: *afterward*, *afterwards*, *ago*, *ahead*, *albeit*, *amid*, *amidst*, *amongst*, *anyhow*, *anyway*, *apart*, *aside*, *away*, *beneath*, *beside*, *besides*, *beyond*, *else*, *elsewhere*, *everybody*, *everybody's*, *everyone*'s, *except*, *front*, *hardly*, *hence*, *hither*, *inside*, *instead*, *just*, *likewise*, *neither*, *oneself*, *otherwise*, *outside*, *somebody*, *somebody's*, *someone*'s, *sometimes*, *somewhat*, *soon*, *sooner*, *thereafter*, *therein*, *thereof*, *thereupon*, *unless*, *versus*, *via* and *yes*;
- (4) Our list includes *thee*, *thine*, *thou*, *thy* and *thyself*, to accommodate to Early Modern English (such as in the works of Shakespeare);
- (5) Our list excludes *find*, *first*, *give*, *go*, *interest*, *keep*, *made*, *noone*, *part*, *put*, *see*, *seem*, *seemed*, *seeming*, *seems*, *show*, *side* and *take*—while some of these words might appear disposable in certain rudimentary information retrieval tasks, we find it harmful to regard them as stop words for machine comprehension of English texts.

<sup>33</sup>In the list above (and also the substitution rules below), the apostrophe refers to either APOSTROPHE U+0027 or RIGHT SINGLE QUOTATION MARK U+2019, as is appropriate for the particular text mining task.

majority of these irregular verbs, however, do exhibit systematic vowel/consonant alternation patterns, such as *keep/kept*, *sleep/slept*, *sweep/swept*, *weep/wept*, *break/broke/broken*, *speak/spoke/spoken*, *steal/stole/stolen*. These systematic irregularities are reflexes of historical Germanic vowel mutations, which are also found in modern Danish, German and Dutch (§5). We will treat these vowel alternations with an *admissible mutation* mechanism (§4.2.2).

#### 4.2.1 Effective spelling and essential root

Before deducing analogs of the Porter stem as the *effective spelling* and the *essential root*, we need to regularize some English verbs (Algorithm 4.3) and normalize the vowel lengths in English spelling (Algorithm 4.4). It should be noted that our purpose is to deduce a token for an English word that is suitable for morphological studies, so the token itself does not have to resemble the word’s etymological stem (as was in the Porter stemming algorithm): in fact, we are going to exploit a wide assortment of special Latin characters (like the German ä, ö, ü, ß, the Spanish ñ, and various accented vowels appearing in Romance languages) as well as some Greek letters<sup>34</sup> to construct the effective spelling and essential root of an English word.

**Algorithm 4.3** (Regularization of English verbs). *For an English word  $\hat{\sigma}$ , its “regularized verb form”  $\text{RegVb}(\hat{\sigma})$  is deduced by the following four steps:*<sup>35</sup>

- (1) *Do*  $\mathbf{X}^e(a|i)bilit\mathbf{X} \rightarrow \mathbf{X}ble$ ,  $\hat{x}ful(ler|les|ness)\mathbf{X} \rightarrow \hat{x}$ , *feet* → *foot*, *heaven* →  $\eta eave\tilde{v}$ , *ificat* $\mathbf{X}$  → *ify*, *ili(s|z)(e|ing)* $\mathbf{X}$  → *ility*, *inabil* $\mathbf{X}$  → *unable*, *inalien* → *unalien*, *istic* $\mathbf{X}$  → *ist*, *positiv* → *poσitiv*, *reciprocal* → *reciprocity*, *toric* $\mathbf{X}$  → *tory*,  $\sim \mathbf{X}^e(\mathit{barb|hil})arit(ies|y) \rightarrow \mathbf{X}arious$ ,  $\sim \mathbf{X}^e(c|d|g|n)uit(ies|y) \rightarrow \mathbf{X}uous$ ,  $\sim alit(ies|y) \rightarrow al$ ,  $\sim arit(ies|y) \rightarrow ar$ ,  $\sim charit(ies|y) \rightarrow charitable$ ,  $\sim circuit(ies|y) \rightarrow circuitous$ ,  $\sim enmit(ies|y) \rightarrow enemy$ ,  $\sim femini(\emptyset|ni)t(ies|y) \rightarrow feminine$ ,  $\sim idit(ies|y) \rightarrow id$ ,  $\sim lar(l|y)s \rightarrow lar$ ,  $\sim necessit(ies|y) \rightarrow necessary$ ,  $\sim ntit(ies|y) \rightarrow ntify$ ,  $\sim simplicit(ies|y) \rightarrow simple$ ,  $\sim tricit(ies|y) \rightarrow tric$ ,  $\sim vanit(ies|y) \rightarrow vain$ ,  $\sim virgin(\emptyset|s) \rightarrow virginity$ .

(2) *In a single sweep, perform replacements according to the following substitution rules:*

- (2) In a single sweep, perform replacements according to the following substitution rules:

$\hat{x}^{\notin}(i)aris(m t)\mathbf{X}$	$\hat{x}^{\notin}(s t)ress$	$(Us Usa)$	$X^{\in}(\emptyset un)dece(it pt)\mathbf{X}\sim$
$\hat{x}ial$	$\hat{x}tor$	$American$	$Xdeceive$
$\mathbf{X}^{\in}(alchem narch)is(m t)\mathbf{X}$	$\mathbf{X}^{\in}(ap im multi)pl(y)(ie)\mathbf{X}$		$\mathbf{X}^{\in}(at g ic it iv nt)ious$
$\mathbf{Xy}$	$\mathbf{Xplication}$		$\mathbf{Xion}$
$\mathbf{X}^{\in}(b h)ang(ed ing)\sim$		$\mathbf{X}^{\in}(cent cloist integ neut sepulch spect)er$	
$\mathbf{X}ang$		$\mathbf{X}re$	
$\mathbf{X}^{\in}(con miscon per precon)cept\mathbf{X}\sim$	$a(c n)eous\mathbf{X}$	$ancestress$	$avuncular\mathbf{X}\sim$
$\mathbf{X}ceive$	$al$	$ancestor$	$uncle$
$children$	$choice$	$clar\sim$	$comic\sim$
$child$	$choose$	$clear$	$comedic$
$flexi(on ve)\mathbf{X}$	$flower$	$focus\mathbf{X}$	$France$
$flect$	$\varphi lower$	$focus$	$French$
$idiot$	$inguish$	$instab\sim$	$Ireland$
$idiot$	$inct$	$unstab$	$Irish$
$langu(id or)\mathbf{X}\sim$	$laughter\sim$	$leg\sim$	$lesson$
$languish$	$laugh$	$\lambda ey$	$\lambda ekcjon$
$menace$	$mentalit\mathbf{X}$	$Mexico$	$mistress$
$\mu e\tilde{n}ace$	$mental$	$Mexican$	$master$
$publish$	$rash\sim$	$rebellion\mathbf{X}$	$rece(ip pt)\mathbf{X}\sim$
$publication$	$\rho ash$	$rebel$	$receive$
$scienti(fic st)\mathbf{X}$	$Scotland$	$severed$	$soror\mathbf{X}\sim$
$science$	$Scottish$	$sever$	$sister$
$tender$	$tenet$	$tenta$	$theori\mathbf{X}\sim$
$te\tilde{v}der$	$tenet$	$\tau enta$	$theory$
$tician\mathbf{X}$			$tician\mathbf{X}$
$tow\tilde{v}$			$tow\tilde{v}$
$traged\sim$			$tragic$
$Ukraine$			$Ukrainian$
$usur$			$usup$
$Wales$			$Welsh$

<sup>34</sup>The Greek letter  $\beta$  should not be confused with German  $\beta$ .

<sup>35</sup>Not all these substitutions operate on verbs, yet most of them are conducive to correct clustering of certain verbs. For example, *missile* and *mission* should not be merged with *miss(Ø|ed)es|ing*, and *passion* should be separated from *pass(Ø|ed)es|ing*. The word family *attent(ion|ive)ly* was indeed etymologically related to the verb *attend*, but their modern meanings are well separated. To avoid such a conflation, while guaranteeing the correct grouping of *inten(d)t*, *exten(d)t* and so on, we need to alter the spelling *attend~*. A diligent reader may find the rationale behind the other entries in our substitution rules.

$weari_{\hat{\chi}^{\neq}}(n)X$	$weariX$	$woke$	$worr$	$chagrin(\emptyset s)$	$conic(\emptyset al)(\emptyset ly)$	$good$	$hung$	$made$	$radii$	$ran$
$\omega e a \rho i$	$\omega e a \rho i$	$wake$	$w o p r$	$chagrine$	$cone$	$\gamma o w \delta$	$hang$	$make$	$radius$	$run$
$\underline{sis}$	$\underline{ten}(\emptyset s)(thX)$			$\sim \hat{\chi}^{\neq}(a e i o u)ar(\emptyset ies ily ize y)$					$\sim \hat{\chi}^{\neq}(p)lus$	
$sister$	$10ten$			$\hat{x}ial$					$\hat{x}li$	
		$\sim(is iz)(ab able ation e ed ement er ing)(\emptyset ly s)$				$\sim X^{\epsilon}(atr fer prec)ocit(ies y)$				
		$ize$				$Xocious$				
$\sim acit(ies y)$	$\sim antal$	$\sim antas(ies v)$	$\sim ics$	$\sim istry$	$\sim monopol(ies y)$	$\sim nness$	$\sim radical(\emptyset s)$	$\sim uum$		
$acious$	$ance$	$antasize$	$ical$	$ister$	$monopolize$	$n$	$radicalize$	$ua$		

(3) *Do εearnest → εaapneστ,  $X^{\epsilon}(a|i|o|u)si(on|ve)X \rightarrow Xde$ , abdomen → abdomin, beautical → beauty, began → begin, caus $X^{\epsilon}(a|e|i) \rightarrow κaωζX$ , comedo → κομεδο, criti(c|que)(∅|al)(∅|ly) → criticize, deer → deer, hung → φaiμ, pretent → pretend, repeat $X \rightarrow repetition$ , savage → saβarye, serious → serious, sever(e|it) → σewere, suffer → swuffer, suspic → suspect, tent → tent, tranquil(∅|l)izX → tranquil, wearinessX → ωeapi, ~risen|rose → rize, ~χtic(∅|al)(∅|ly)()izX → χsis, ~s's → ∅.*

(4) *Replace*

$\sim$	$X^{\epsilon}(ab pre)sence$	$X^{\epsilon}(ad de)vi(c z)$	$X^{\epsilon}(col e)lid$	$X^{\epsilon}(de of)fen(c s)$	$X^{\epsilon}(di pro)vid$	
$\emptyset$	$Xsent$	$Xvis$	$Xlis$	$Xfend$	$Xvis$	
$X^{\epsilon}(dissatis lique petri putre rare satis stupe vitri)f(ie ied ies y ying)$						
			$Xfact$			
	$X^{\epsilon}(g i r s)onal(\emptyset i ly)X$	$X^{\epsilon}(k K)ing~$	$X^{\epsilon}(mpl rt tr)aint$	$attend~$	$aur$	$automation$
	$Xon$	$Xoenig$	$Xain$	$attenzd$	$aor$	$automasis$
$begg(ar ed)~$	$counten$	$deX^{\epsilon}(c r)id$	$duct(\emptyset ion)$	$dye$	$expence$	$explain~$
$beg$	$kouten$	$deXis$	$duce$	$dzye$	$expense$	$explan$
$forg(e o)t(\emptyset ten)X$	$former~$	$four$	$inal(\emptyset i ly)X$	$listen$	$mad$	$mention$
$vye\beta$	$φormer$	$φour$	$in$	$licen(c s)~$	$menzzion$	$migrant$
$mission$	$passion$	$pool$	$possess$	$licent$	$μφαδ$	$migrate$
$mit$	$patzion$	$powl$	$pozzezz$	$lisztzn$		$missle$
$sump$	$temper$	$thought$	$typ~$	$prett$	$prevail~$	$rouse$
$sume$	$temper$	$think$	$zyp$	$preet$	$prevent~$	$solde~$
$(mete meted metes meting metings)$						
$(ridden rode)$		$X^{\epsilon}(\emptyset in mis off out over p under)led$		$Xlain$	$X(b br c l r)ook$	
$ride$		$Xlead$		$Xlay$	$Xooked$	
$brought$	$caught$	$eat(\emptyset en ing s)$	$lad(\emptyset die dish hood)(\emptyset s)$	$meant$	$sought$	$taught$
$bring$	$catch$	$ate$	$λaδ$	$mean$	$seek$	$teach$
$(~sis siz)$		$~(did does doing done)$		$~(goes going gone went)$		$~bought$
$ses$		$do$		$go$		$buy$
$~ced(e ed es ing)$	$~ered$	$~gave$	$~X-in-law$	$~heard$	$~itten$	$~stood$
$cez\beta$	$er$	$give$	$λawX$	$hear$	$iten$	$stand$
				$rn$	$saw$	$vise$

The silent *e* in English spelling may affect the length of the immediately preceding vowel. Pairs of words like *hat/hate* and *rat/rate* are etymologically unrelated. Motivated by modern Dutch orthography (§5.3.1), we will introduce a heuristic procedure to mark the vowel lengths in English words explicitly. This procedure not only detects long vowels in CVCe patterns, but also words derived therefrom. (For example, the vowel *o* in *devote*, *devoting* and *devotion* should all be consistently marked as long.)

**Algorithm 4.4** (Normalization of English vowel lengths). *For a text string  $σ̂$ , its “normalized English vowel length form” NormVL( $σ̂$ ) is determined through the following procedure:*

- If  $σ̂$  comes from a word that is “heuristically capitalized” (Definition 4.1), then  $\text{NormVL}(σ̂) = zxσ̂'$  where  $σ̂'$  results from doing  $\sim ean(\emptyset|s) \rightarrow \emptyset$  on  $σ̂$ .<sup>36</sup>*

<sup>36</sup>Here, the prefix “zx” push the word down to the end of an alphabetized list later. In English, the ending  $\sim ean$  is found mainly in adjectives (and nouns) deriving from proper nouns. The procedure here ensures that words like *Argentine*(∅|an)(∅|s), *Boole*(∅|an) and *Euclid*(∅|ean) *Europe*(∅|an)(∅|s) are clustered properly.

(2) Otherwise, perform the following substitutions, in five sequential sweeps:

(2.1) Do  $\sim osit(ies|y) \rightarrow ous$ ,  $\sim \hat{\chi}\hat{\chi}oe(\emptyset|d|s) \rightarrow \hat{\chi}\hat{\chi}o$ .

(2.2) Replace<sup>37</sup>

$(lose loss)$	$\hat{\chi}tr$	$X^{\infty}(dis ex pre)tens$	$X^{\infty}(ex sus)pens$
$lost$	$\hat{\chi}ter$	$Xtend$	$Xpend$
$X^{\infty}(\hat{\chi}(a o))tion$	$au$	$capab\sim$	$deep(\emptyset eX)$
$Xtiv$	$aau$	$kapab$	$deep(\emptyset eX)$
$manner$	$ne$	$depth$	$earl(i y)\sim$
$maznner$	$ne$	$earzli$	$eed$
$mzere$	$oue$	$depth$	$fer$
$modz$	$ú$	$earzli$	$fine\sim$
$sudden\sim$	$summer$	$water\sim$	$freedom$
$swudden$	$suzmmer$	$ton\sim$	$gas\sim$
		$wazter$	$generat-$
		$woman$	$genero-$
		$idle$	$ll$
		$idl(est y)$	
		$genera$	
		$genus$	
		$gently$	
		$gentle$	
		$singly$	$ski(\emptyset ed ing)s$
		$single$	$szki$

(2.3) Replace every occurrence of double letter by a corresponding single capital letter (i.e.  $ee \rightarrow E$  etc.).

(2.4) Set  $C_1$  as any lowercase English letter<sup>38</sup> other than  $\{a, e, i, o, u, y\}$ ,  $V$  as any one among  $\{a, e, i, o, u\}$ , and  $C_2$  as any lowercase English letter other than  $\{a, e, i, o, u, w, y\}$ . Define  $V^+$  as the capital form of the vowel  $V$ . Replace<sup>39</sup>

$C_1 V^{\infty}((C_2(a a a a at e ing ion ish it iv ous)) ck gh nd))$	$ang(e ing)$
$C_1 V^+ X$	$Ange$

(2.5) Replace<sup>40</sup>

$A$	$E$	$\acute{e}$	$I$	$lez\beta$	$O$	$T$	$\acute{u}$	$U$	$\lambda$	$gOv\sim$	$inter\sim$	$be(Ter st)$
$aa$	$ee$	$e$	$i$	$zlez\beta$	$\acute{o}$	$tt$	$oo$	$u$	$ll$	$góv$	$jntr$	$\gamma\omega\delta$

before converting the string to lowercase.

With the preparations above, we can derive the effective spelling of an English word. The effective spelling (to be derived in Algorithm 4.5) is a conservative estimate for the (etymologically correct) word stem, which usually has longer string length than the essential root (to be derived in Algorithm 4.8).

**Algorithm 4.5** (English effective spelling). *For an English word  $\hat{\sigma}$ , its effective spelling EffSpell( $\hat{\sigma}$ ) is constructed in seven sequential steps:*

(1) Take NormVL(RegVb( $\hat{\sigma}$ )), and perform the following replacements in a single sweep:

$\hat{\chi}aid$	$ck$	$clamat$	$emn$	$enten$	$family$	$gral$	$íoón$	$iv$	$la$	$ld$	$lóó$	$ng$	$oi$	$ou$	$ply$
$\hat{\chi}ayed$	$c$	$claimat$	$μn$	$entzen$	$families$	$graation$	$ion$	$ív$	$lôa$	$llz$	$lo$	$zñ$	$zö$	$i$	$plie$
$rd$	$sh$	$tim$	$ui$	$ye$	$\sim'(\emptyset s)(\emptyset ')$	$\sim ant$	$\sim ble(\emptyset d s)$	$\sim hie(r st)$	$\sim ook$	$\sim s$	$-$	$(hyphen)$			
$rzd$	$\check{s}$	$ztim$	$zü$	$yme$	$\emptyset$	$ance$	$bly$	$hy$	$aaken$	$\emptyset$	$\emptyset$				

(2) Replace

$\hat{\chi}ied\sim$	$C_1^{\infty}(c w)a\sim$	$ba\sim$	$C_2^{\infty}(s t w)X^{\infty}(e o \acute{o}\sim)$	$dóó\sim$	$f^{N_1^{\infty}(a e)\sim}$	$f^{N_2^{\infty}(i i)\sim}$	$móó\sim$	$dad(\emptyset y)$	$m(o u)m$
$\hat{\chi}eed$	$C_1\beta a$	$\beta a$	$C_2zX$	$\delta o$	$\varphi V_1$	$ffV_2$	$\mu o$	$\varphi a\theta er$	$mother$

(3) Do  $\sim f(\emptyset|e) \rightarrow ve$ ,  $\sim xe \rightarrow x$ ,  $\sim \hat{\chi}\hat{\chi}\hat{\chi}ly \rightarrow \hat{\chi}\hat{\chi}\hat{\chi}$  (i.e. remove final “ly” if it is preceded by at least three letters).

(4) Do  $h(i|i) \rightarrow h\bar{i}$ ,  $lu \rightarrow lvu$ ,  $y \rightarrow ie$ .

(5) Do  $\sim \hat{\chi}\hat{\chi}\hat{\chi}(tic|tion) \rightarrow \hat{\chi}\hat{\chi}\hat{\chi}t$ ,  $\sim \hat{\chi}\hat{\chi}\hat{\chi}(age|ed|ful|ical|izñ|izñ|ment|nez\beta) \rightarrow \hat{\chi}\hat{\chi}\hat{\chi}$ .

<sup>37</sup>As we do  $deep(\emptyset|eX) \rightarrow depth$  to the words *deep*, *deeper*, *deepest*, *deeply*, we obtain *depth*, *depth*, *depth*, *depthly*. (In Mathematica, the command `StringReplace[{"deep", "deeper", "deepest", "deeply"}, {"deep"~~""|("e"~~_) :> "depth"}]` results in `{"depth", "depth", "depth", "depthly"}`.) Here, the wildcard symbol  $X \in \{r, st\}$  is deleted after replacement, because  $X$  does not reappear in the destination string  $\hat{\sigma}' = depth$  (cf. Definition 3.6).

<sup>38</sup>A letter is an English letter if and only if it is one of the 26 members in the ISO basic Latin alphabet.

<sup>39</sup>For readers’ benefit, we point out that the corresponding Mathematica substitution rules are `{(C1: Except["e"|"i"|"o"|"u"|"y"], CharacterRange["b", "z"]])~(V: "a"|"e"|"i"|"o"|"u"~~(X: (Except["e"|"i"|"o"|"u"|"w"|"x"|"y"], CharacterRange["b", "z"]))~~"ag"|"al"|"al"|"at"|"e"|"ing"|"ish"|"it"|"ous") | "ck"|"gñ"|"nd") :> C1 <> ToUpperCase[V] <> X, "ang"~~"e"|"ing" :> "Ange"}`. As our substitutions disallow overlaps, the letter *e* in patterns like *ndering* (as in *plundering*, *pondering*) will not be marked as a long vowel.

<sup>40</sup>Note that during a single sweep, the longer matches take priority over shorter matches. Therefore, the result `NormVL(government) = góvernment` does not contain *óó*.

## (6) Replace

$(eft iev \hat{iev})$	$(o óó)^{\mathbf{X}^e}(p r)$	$ch$	$ea$	$ir$	$ll$	$mi\sim$	$ph$	$th$	$\sim\hat{x}\notin(d)ent(\emptyset ive)$	$\sim\mathbf{X}^e(pe s)ion$	$\sim\hat{x}\hat{x}\hat{x}i\check{s}$	$\sim lt$	$\sim own$
$\ddot{a}\ddot{a}ve$	$\delta\mathbf{X}$	$\check{c}$	$\ddot{a}\ddot{a}$	$\check{ir}$	$\lambda$	$\mu i$	$\varphi\varphi$	$\theta$	$\hat{x}eend$	$\mathbf{X}iv$	$\hat{x}\hat{x}\hat{x}$	$l$	$ow$

(7) Replace<sup>41</sup>

$beg$	$d\ddot{a}\ddot{a}(d \theta)$	$\hat{x}_1\mathbf{X}_1\hat{x}_2(e i)e$	$\underline{zx}\mathbf{X}_3$
$bbeg$	$die$	$\hat{x}_1\mathbf{X}_1\hat{x}_2ae$ , if $\hat{x}_1\mathbf{X}_1 = zx\mathbf{X}_2$ ; $\hat{x}_1\mathbf{X}_1\hat{x}_2$ otherwise.	$zx\widetilde{\mathbf{X}}_3$ , where $\widetilde{\mathbf{X}}_3$ results from doing $\sim ie \rightarrow y$ on $\mathbf{X}_3$ .

Here, “ $\hat{x}_1\mathbf{X}_1 = zx\mathbf{X}_2$ ” means that the string  $\hat{x}_1\mathbf{X}_1$  begins with two letters “zx”.

Up to the current stage, we have introduced many non-English letters into the tokens assigned to English words. Some of these non-English letters will be regarded as extensions to the set of English vowels.

**Definition 4.6** (English Vowel Extensions). Hereafter in §4, the symbol  $\mathbf{V}^*$  stands for any member from the list  $\{a, \ddot{a}, e, \dot{e}, i, \dot{i}, \hat{i}, \ddot{i}, o, \acute{o}, \ddot{o}, \ddot{o}, u, \ddot{u}\}$ , the so-called English vowel extensions. In line with the multiplicity notations introduced in Definition 3.3, the symbol  $\mathbf{V}_m^*$  stands for a text string formed by consecutive appearance of one or more (not necessarily identical or distinct) members from the set of English vowel extensions.

Dual to the notations above, the symbol  $\mathbf{C}^*$  stands for any character that does not belong to the list  $\{a, \ddot{a}, e, \dot{e}, i, \dot{i}, \hat{i}, \ddot{i}, o, \acute{o}, \ddot{o}, u, \ddot{u}\}$ , and  $\mathbf{C}_{m_0}^*$  stands for a text string formed by consecutive appearance of zero or more (not necessarily identical or distinct) characters that do not belong to the same list.  $\square$

As mentioned before, the conservative effective spelling algorithm assigns a long token to an English word, while the aggressive essential root algorithm will produce a short token, removing as many final letters as “appropriate”. To ensure “appropriateness” during the short token generation, we need to explicitly label some part of an English word as unremovable. The string length of this unremovable part is called the English protected range, as defined below.

**Definition 4.7** (English protected range). Let  $\hat{\sigma}$  be the effective spelling of an English word, its protected range  $\text{ProtRg}(\hat{\sigma})$  is an integer determined as follows:<sup>42</sup>

- Look for the string pattern  $(\emptyset|bel|cinter|co|cô|dee|di|fô|ha|ma|pre|pro|próó|r\mathbf{V}_m^*|ster|su|te|\mathbf{V}_m^*)\mathbf{C}_{m_0}^*\mathbf{V}_m^* \sim$  in the string  $\hat{\sigma}$ ;
- If the string pattern above is found, the last position occupied by such a string defines  $\text{ProtRg}(\hat{\sigma})$ ; otherwise, set  $\text{ProtRg}(\hat{\sigma}) = 0$ .  $\square$

*Example 4.7.1.* Let  $\hat{\sigma} = TNT$ , then  $\text{EffSpell}(\hat{\sigma}) = zxTNT$ , and  $\text{ProtRg}(\text{EffSpell}(\hat{\sigma})) = 0$ ; let  $\hat{\sigma} = trinitrotoluene$  (the IUPAC name for TNT), then  $\text{EffSpell}(\hat{\sigma}) = trinitrotoluuen$ , and  $\text{ProtRg}(\text{EffSpell}(\hat{\sigma})) = 3$ .

*Example 4.7.2.* Let  $\hat{\sigma} = DDT$ , then  $\text{EffSpell}(\hat{\sigma}) = zxDDT$ , and  $\text{ProtRg}(\text{EffSpell}(\hat{\sigma})) = 0$ ; let  $\hat{\sigma} = 1,1'-(2,2,2-trichloroethane-1,1-diyl)bis(4-chlorobenzene)$  (the IUPAC name for DDT), then

$$\text{EffSpell}(\hat{\sigma}) = 1,1(2,2,2tričlôrieθaane1,1diiel)bi(4člôróóbenzeen),$$

and  $\text{ProtRg}(\text{EffSpell}(\hat{\sigma})) = 39$ . Here is how the effective spelling looks like with its first 39 characters underlined:

$$\underline{1,1(2,2,2tričlôrieθaane1,1diiel)bi(4člôróóbenzeen)}.$$

We note that the string pattern of interest here is counted from an internal word boundary (the numerical digit 4).

**Algorithm 4.8** (English essential root). Let  $\hat{\sigma}$  be the effective spelling of an English word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:

- (1) Break down  $\hat{\sigma} = \hat{\sigma}_1\hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}^{[\text{ProtRg}(\hat{\sigma})]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma})$  is equal to the protected range of  $\hat{\sigma}$ .
- (2) On  $\hat{\sigma}_2$ , perform the following substitutions in a sequel:<sup>43</sup>

<sup>41</sup>Here, the status of heuristic capitalization affects our final treatments. If a word is heuristically capitalized, its final *e* (not preceded by *e* or *i*) will be altered to *ae*; if a word is heuristically uncapitalized, its final *e* (not preceded by *e* or *i*) will be removed. Furthermore, If a word is heuristically capitalized, its final *ie* will be changed to *y*.

<sup>42</sup>Roughly speaking, we want to protect an English word up to its first vowel (or vowel cluster) counting from an internal or external word boundary, bypassing some common prefixes.

<sup>43</sup>In other words, the core algorithm for essential root extraction runs as follows: keep the last “strong” vowel *a*, *i*, *o* or *u* in non-final position, plus one subsequent letter; delete final *a*; erase the final appearance of *e* and all the letters thereafter.

- (2.1) Do  $\sim \hat{\chi}^{\epsilon}(a|i|o|u)\mathbf{X}^{\epsilon}\overline{(a|i|o|u)}_m \rightarrow \hat{\chi}\mathbf{X}_1$ , where  $\mathbf{X}_1$  is the first character in  $\mathbf{X}$ .  
 (2.2) Do  $\sim a \rightarrow \emptyset$ .  
 (2.3) Do  $\sim e\overline{(e)}_{m_0} \rightarrow \emptyset$ .

The result after these three steps of operations is called  $\hat{\sigma}'_2$ .

- (3) Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ .  
 (4) Do  $\sim pe\lambda \rightarrow pel$ ,  $reebel$   $\rightarrow reebel\lambda$ .

*Example 4.8.1.* Let  $\hat{\sigma}$  be the IUPAC name for TNT, then  $\text{EssRoot}(\text{EffSpell}(\hat{\sigma})) = \text{tríinitrotolu}$ . Let  $\hat{\sigma}$  be the IUPAC name for DDT, then  $\text{EssRoot}(\text{EffSpell}(\hat{\sigma})) = 1,1(2,2,2\text{tričlóríeθaanel}, 1\text{diel})bi(4\text{člóróbenze})$ . Note, in particular, that in the latter example, the last internal word boundary (closing parenthesis) is preserved in the essential root.

*Example 4.8.2.* Here are more examples showing how the functions  $\text{RegVb}$ ,  $\text{NormVL}$ ,  $\text{EffSpell}$ ,  $\text{EssRoot}$  act on some etymologically related words.

$\hat{\sigma}$	$\text{RegVb}(\hat{\sigma})$	$\text{NormVL}(\text{RegVb}(\hat{\sigma}))$	$\text{EffSpell}(\hat{\sigma})$	$\text{EssRoot}(\text{EffSpell}(\hat{\sigma}))$
<i>environment</i>	<i>environment</i>	<i>environment</i>	<i>enviřon</i>	<i>enviřon</i>
<i>environmental</i>	<i>environmental</i>	<i>environmental</i>	<i>enviřonm</i>	<i>enviřonm</i>
<i>environmentalist</i>	<i>environmentalist</i>	<i>environmentalist</i>	<i>enviřonmentalist</i>	<i>enviřonm</i>
<i>protect</i>	<i>protect</i>	<i>próóTECT</i>	<i>próóTECT</i>	<i>próóTECT</i>
<i>protection</i>	<i>protection</i>	<i>próóTECTION</i>	<i>próóTECT</i>	<i>próóTECT</i>
<i>protector</i>	<i>protector</i>	<i>próóTector</i>	<i>próóTECTôR</i>	<i>próóTECTôR</i>
<i>protectorate</i>	<i>protectorate</i>	<i>próóTECTÓRáte</i>	<i>próóTECTÔRat</i>	<i>próóTECTÔRat</i>

We note that the vowel *e* in *protect* had not been deleted in the construction of essential root. In fact, it remains in the unremovable  $\hat{\sigma}_1$  in Algorithm 4.8(1).

#### 4.2.2 Admissible mutation and approximate clustering

Unlike the Porter stemming algorithm, we shall not use the exact match of tokens as the sole criterion for word clustering. Instead, we shall sort words by their tokens, and compare consecutive neighbors in such an alphabetized list. Roughly speaking, our clustering algorithm (see Algorithm 4.15) runs as follows: if the tokens of two neighbors in such a list fail certain “approximate matching criteria”, we add a demarcation line between them; after comparisons between all the consecutive neighbors, we split the alphabetized list into word clusters.

As noted before, some English irregular verbs exhibit systematic patterns of vowel alternations, such as *drink/drank/drunk* and *sing/sang/sung*. However, the tokens for the separate forms of an irregular verb may not sit next to each other in an alphabetized list. To resolve this issue in token sorting order, we need to blot certain vowels (more precisely, the extended vowels prescribed by Definition 4.6) occurring in our token, by the following algorithm.

**Algorithm 4.9** (English vowel blotting). *For a token string  $\hat{\sigma}$ , its blotted form  $\text{BlotV}_1(\hat{\sigma})$  is constructed as follows:*

- If the first appearance of pattern  $\mathbf{V}_m^*$  occurs neither word-initially nor word-finally (as measured by external word boundaries of  $\hat{\sigma}$ ), then replace the aforementioned pattern with a single letter “a”.
- Otherwise, leave the string  $\hat{\sigma}$  intact.

*Example 4.9.1.* As a continuation of Example 4.8.2, one can append one more column  $\text{BlotV}_1(\text{EssRoot}(\text{EffSpell}(\hat{\sigma})))$  with entries *environ*, *enviřonm*, *enviřonm*, *práTECT*, *práTECT*, *práTECTôR*, *práTECTÔRat*.

*Example 4.9.2.* A few more examples:

$\hat{\sigma}$	<i>automatic</i>	<i>bee</i>	<i>find</i>	<i>join</i>	<i>solo</i>
$\text{BlotV}_1(\text{EssRoot}(\text{EffSpell}(\hat{\sigma})))$	<i>aaútóómat</i>	<i>bee</i>	<i>ffand</i>	<i>jzan</i>	<i>szalo</i>

Our “approximate matching criteria” are encoded in a bivariate Boolean-valued function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$ , which tests whether string  $\hat{\beta}$  is a “legitimate heir” to string  $\hat{\alpha}$ . Before stating this “heredity test function” in its entirety (Algorithm 4.14), we shall describe a “simple heredity test function” in Algorithm 4.10, which judges similarities between two strings, without worrying about the problems associated with irregular verbs.

**Algorithm 4.10** (Simple heredity test). Let  $\hat{\alpha}'$  be the result of doing  $\sim e \rightarrow \emptyset$ ,  $\hat{\chi} lie(er|st) \rightarrow \hat{\chi}$  on  $\hat{\alpha}$ . The Boolean-valued function SimpHrdTest( $\hat{\alpha}, \hat{\beta}$ ) returns **TRUE** if the lowercase form of  $\hat{\alpha}$  contains at least one instance of  $\mathbf{V}^*$  (Definition 4.6) **AND** at least one of the following 12 conditions holds:<sup>44</sup>

- (i)  $\hat{\alpha} = \hat{\beta}$ ;
- (ii)  $\hat{\alpha} = \hat{\beta}^\circ$ , where  $\hat{\beta}^\circ$  results from doing  $\sim_{\hat{\chi}\#}(a|e|i|o|u)ster \rightarrow \hat{\chi}$  on  $\hat{\beta}$ .
- (iii)  $\hat{\alpha} = \hat{\beta}^*$ , where  $\hat{\beta}^*$  results from doing  $\hat{\chi} lie(er|st) \rightarrow \hat{\chi}$  on  $\hat{\beta}$ ;
- (iv)  $\hat{\alpha} = \hat{\beta}^\dagger$ , where  $\hat{\beta}^\dagger$  results from doing  $\sim ie \rightarrow i$  on  $\hat{\beta}$ ;
- (v)  $\hat{\alpha} = \hat{\beta}^\ddagger$ , where  $\hat{\beta}^\ddagger$  results from doing  $\sim us \rightarrow u$  on  $\hat{\beta}$ ;
- (vi) Appending the last letter of string  $\hat{\alpha}$  to itself, one obtains  $\hat{\beta}$ ;
- (vii)  $\hat{\beta} = \hat{\alpha}'i$ ;
- (viii)  $\hat{\beta} = \hat{\alpha}iv$ ;
- (ix)  $\hat{\beta} = \hat{\alpha}n$ ;
- (x)  $\hat{\beta} = (\hat{\alpha}|\hat{\alpha}')\check{skip}$ ;
- (xi)  $\hat{\beta} = \hat{\alpha}\theta$ ;
- (xii)  $\ell(\hat{\beta}) - 1 > \ell(\hat{\alpha}') \geq \frac{\ell(\hat{\beta})}{2}$  **AND**  $\hat{\alpha} = \hat{\beta}^{[\ell(\hat{\alpha}')]}(\emptyset|e)$  **AND**  $\hat{\alpha}'\mathbf{X}_1(gu|i\hat{o}|ow|z\beta)\mathbf{X}_2 \neq \hat{\beta}$  **AND**  $\hat{\beta}^{[\ell(\hat{\alpha}')+1]} = (a|e|h|i|o|ó|ó|r|u)$ . (See Definition 3.1 for the notations  $\hat{\beta}^{[n]}$  and  $\hat{\beta}^{[n]}$ .)

When SimpHrdTest( $\hat{\alpha}, \hat{\beta}$ ) = **FALSE**, it is still possible that the two underlying words are etymologically related. To detect affinity between words in such scenarios, we need to invoke the Needleman–Wunsch (NW) and Smith–Waterman (SW) algorithms for more detailed string comparisons (cf. Definition 3.7).

**Algorithm 4.11** (Roots and suffixes by NW and SW). For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the function values RootNW( $\hat{\alpha}, \hat{\beta}$ ), SuffixNW( $\hat{\alpha}, \hat{\beta}$ ) and NW\*( $\hat{\alpha}, \hat{\beta}$ ) are determined through the following procedure:

- Use the Needleman–Wunsch algorithm to align the sequences as NW( $\hat{\alpha}, \hat{\beta}$ ).
- If NW( $\hat{\alpha}, \hat{\beta}$ ) ends with a mismatch (shown in brackets in Examples 3.7.1 and 3.7.2), use this mismatch to define SuffixNW( $\hat{\alpha}, \hat{\beta}$ ); otherwise, define SuffixNW( $\hat{\alpha}, \hat{\beta}$ ) =  $[\emptyset, \emptyset]$ .
- If NW( $\hat{\alpha}, \hat{\beta}$ ) starts with a matching string, define RootNW( $\hat{\alpha}, \hat{\beta}$ ) as the concatenation of all the (not necessarily contiguous) matching string portions in NW( $\hat{\alpha}, \hat{\beta}$ ), and define NW\*( $\hat{\alpha}, \hat{\beta}$ ) by dropping the first matching string from NW( $\hat{\alpha}, \hat{\beta}$ ) as well as its contribution (if non-void) to SuffixNW( $\hat{\alpha}, \hat{\beta}$ ); otherwise, define RootNW( $\hat{\alpha}, \hat{\beta}$ ) =  $[\#\%] = [\text{U+0023} | \text{U+0025}]$  (a nonsensical string containing two characters) and define NW\*( $\hat{\alpha}, \hat{\beta}$ ) by deleting from NW( $\hat{\alpha}, \hat{\beta}$ ) the contribution (if non-void) to SuffixNW( $\hat{\alpha}, \hat{\beta}$ );

The function values RootSW( $\hat{\alpha}, \hat{\beta}$ ), SuffixSW( $\hat{\alpha}, \hat{\beta}$ ) and SW\*( $\hat{\alpha}, \hat{\beta}$ ) are determined similarly, through the SW function.

**Algorithm 4.12** (Admissible suffix mismatch). For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function

$$\text{AdmSM}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta}))$$

returns **TRUE** if the lowercase form of RootNW( $\hat{\alpha}, \hat{\beta}$ ) contains at least one instance of  $\mathbf{V}^*$  (Definition 4.6) **AND** NW\*( $\hat{\alpha}, \hat{\beta}$ ) is void **AND** SuffixNW( $\hat{\alpha}, \hat{\beta}$ )  $\neq [(\text{nd}|ow|z\beta), (\text{nd}|ow)]$ <sup>45</sup> **AND** at least one of the following four conditions holds:

- (i) SuffixNW( $\hat{\alpha}, \hat{\beta}$ ) =  $[i, n]$ ;
- (ii) SuffixNW( $\hat{\alpha}, \hat{\beta}$ ) =  $[i, u]$ ;
- (iii) SuffixNW( $\hat{\alpha}, \hat{\beta}$ ) =  $[\emptyset, s]$  **AND**  $\mathcal{Q}(\text{RootNW}(\hat{\alpha}, \hat{\beta})) = i$ ;
- (iv) SuffixNW( $\hat{\alpha}, \hat{\beta}$ ) =  $[(\emptyset|i)((a|e|h)\mathbf{X}), (en|i|iz|iz\mathbf{X}|(i|ia|io|i\hat{o}|iv|m|r|\hat{\chi}')\hat{\chi}\mathbf{X})]$  where  $\hat{\chi}' = \text{RootNW}(\hat{\alpha}, \hat{\beta})^{[\ell(\text{RootNW}(\hat{\alpha}, \hat{\beta})}]}$  is the last character in RootNW( $\hat{\alpha}, \hat{\beta}$ ).

<sup>44</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical **OR**.

<sup>45</sup>Hereafter,  $[\hat{\sigma}_1, \hat{\sigma}_2] \neq [\hat{\tau}_1, \hat{\tau}_2]$  means  $(\hat{\sigma}_1 \neq \hat{\tau}_1 \text{ AND } \hat{\sigma}_2 \neq \hat{\tau}_2)$ , while  $[\hat{\sigma}_1, \hat{\sigma}_2] = [\hat{\tau}_1, \hat{\tau}_2]$  means  $(\hat{\sigma}_1 = \hat{\tau}_1 \text{ AND } \hat{\sigma}_2 = \hat{\tau}_2)$ .

Similarly, one can evaluate another Boolean-valued function

$$\text{AdmSM}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta}))$$

by trading all the occurrences of NW in the statements above with SW.

**Algorithm 4.13** (Admissible vowel alternation). *For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function*

$$\text{AdmVA}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta}))$$

*returns TRUE if  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = ([\emptyset](e|i|n)\mathbf{X})$ ,  $(\emptyset|(e|i|n|t|z)\mathbf{X})$  AND  $\text{NW}^*(\hat{\alpha}, \hat{\beta})$  consists of one mutation bracket and one non-void matching string AND  $\text{RootNW}(\hat{\alpha}, \hat{\beta}) \neq \mathbf{X}_1(i|u)\mathbf{X}_2$  AND at least one of the following three conditions holds:*

- (i)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = ([a, e][a, i][a, u][ää, öö][ää, öö][e, \emptyset][e, o][ee, oo][ee, öö][i, u][i, öö][i, öö][oo, öö])\mathbf{X}$
- (ii)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = [\emptyset, \hat{\chi}']\mathbf{X}$  where  $\hat{\chi}' = \text{RootNW}(\hat{\alpha}, \hat{\beta})^{\ell(\text{RootNW}(\hat{\alpha}, \hat{\beta}))}$  is the last character in  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$ ;
- (iii)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = [\emptyset, e]e$ .

Similarly, one can evaluate  $\text{AdmVA}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta}))$ .

**Algorithm 4.14** (Heredity test function). *Let  $\hat{\alpha}'$  be the result of performing a substitution  $\sim e \rightarrow \emptyset$  on  $\hat{\alpha}$ . For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  returns TRUE if at least one of the following three conditions holds:*

- (i)  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta}) = \text{TRUE}$ ;
- (ii)  $\ell(\text{RootNW}(\hat{\alpha}, \hat{\beta})) \geq \frac{1}{2} \max\{\ell(\hat{\alpha}'), \ell(\hat{\beta})\}$  AND  $(\text{AdmSM}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta})) = \text{TRUE}$  OR  $\text{AdmVA}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta})) = \text{TRUE}$ );
- (iii)  $\ell(\text{RootSW}(\hat{\alpha}, \hat{\beta})) \geq \frac{1}{2} \max\{\ell(\hat{\alpha}'), \ell(\hat{\beta})\}$  AND  $(\text{AdmSM}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta})) = \text{TRUE}$  OR  $\text{AdmVA}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta})) = \text{TRUE}$ ).

**Algorithm 4.15** (Approximate clustering of English words). *The approximate clustering of a list of English words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  is completed in two stages:*

- (1) *We sort the list  $\{(\hat{\alpha}_1, \text{EffSpell}(\hat{\alpha}_1)), \dots, (\hat{\alpha}_N, \text{EffSpell}(\hat{\alpha}_N))\}$  alphabetically according to the second component (effective spelling) of each entry. If two consecutive neighbors in the alphabetized list  $\{(\hat{\alpha}_{(1)}, \text{EffSpell}(\hat{\alpha}_{(1)})), \dots, (\hat{\alpha}_{(N)}, \text{EffSpell}(\hat{\alpha}_{(N)}))\}$  satisfy  $\text{HrdTest}(\text{EffSpell}(\hat{\alpha}_{(n)}), \text{EffSpell}(\hat{\alpha}_{(n+1)})) = \text{FALSE}$ , where  $n \in \mathbb{Z} \cap [1, N]$ , then add a demarcation line between these two entries. In this way, the alphabetized list is divided into separate groups of words tagged with their effective spellings:  $\{g_1 = \{(\hat{\alpha}_{(1,1)}, \text{EffSpell}(\hat{\alpha}_{(1,1)})), \dots, (\hat{\alpha}_{(1,n_1)}, \text{EffSpell}(\hat{\alpha}_{(1,n_1)}))\}, \dots, g_M = \{(\hat{\alpha}_{(M,1)}, \text{EffSpell}(\hat{\alpha}_{(M,1)})), \dots, (\hat{\alpha}_{(M,n_M)}, \text{EffSpell}(\hat{\alpha}_{(M,n_M)}))\}\}$ , where each sublist also preserves alphabetic order. In particular, we have  $\hat{\alpha}_{(1,1)} = \hat{\alpha}_{(1)}$  and  $\hat{\alpha}_{(M,n_M)} = \hat{\alpha}_{(N)}$ .*
- (2) *For each group of words  $g_m = \{(\hat{\alpha}_{(m,1)}, \text{EffSpell}(\hat{\alpha}_{(m,1)})), \dots, (\hat{\alpha}_{(m,n_m)}, \text{EffSpell}(\hat{\alpha}_{(m,n_m)}))\}$  where  $m \in \mathbb{Z} \cap [1, M]$ , we augment it into a tagged entry  $G_m = (g_m, \text{EffSpell}(\hat{\alpha}_{(m,1)}), \text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})), \text{BlotV}_1(\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))))$ . The list  $\{G_1, \dots, G_M\}$  is sorted alphabetically, with respect to  $\text{BlotV}_1(\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})))$  (with highest priority),  $\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))$  (with medium priority), and  $\text{EffSpell}(\hat{\alpha}_{(m,1)})$  (with lowest priority). If two consecutive neighbors in the alphabetized list  $\{G_{(1)} = (g_{(1)}, \hat{\gamma}'_{(1)}, \hat{\gamma}''_{(1)}, \hat{\gamma}'''_{(1)}), \dots, G_{(M)} = (g_{(M)}, \hat{\gamma}'_{(M)}, \hat{\gamma}''_{(M)}, \hat{\gamma}'''_{(M)})\}$  satisfy*

$$\text{HrdTest}(\hat{\gamma}''_{(m)}, \hat{\gamma}''_{(m+1)}) = \text{FALSE}$$

AND

$$\text{HrdTest}(\hat{\gamma}'_{(m)}, \hat{\gamma}'_{(m+1)}) = \text{FALSE}$$

where  $m \in \mathbb{Z} \cap [1, M]$ , then add a demarcation line between these two entries. In this manner, the alphabetized list  $\{G_{(1)}, \dots, G_{(M)}\}$  is divided into separate groups of tagged clusters  $\{\Gamma_1 = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, \Gamma_K = \{G_{(K,1)}, \dots, G_{(1,m_K)}\}\}$ . Finally, the output list of word clusters  $\{\check{\Gamma}_1, \dots, \check{\Gamma}_K\}$  is constructed by discarding all the tags (effective spellings, essential roots, vowel blotted forms) from each  $\Gamma_k$ , where  $k \in \mathbb{Z} \cap [1, K]$ .

*Example 4.15.1.* To partially demonstrate the capabilities of our clustering algorithm, we pick the following list of English words:

abolish, abolished, abolishes, abolishing, abolishment, abolition, abolitionism, abolitionist, abolitionists, admonish, admonishable, admonisher, admonishes, admonishing, admonishingly, admonishment, admonition, admonitory, aid, aided, aids, America, American, analyses, analysis, analytic, analytically, analyze, analyzed, ant, antenna, antennae, anterior, anxieties, anxiety, anxious, anxiously, arise, arisen, arose, automatic, automatically, bad, bat, beauties, beautiful, beautifully, beauty, bed, bedding, beds, beg, began, beggar, begged, begin, beginning, begun, belie, belied, belief, beliefs, belies, believe, believed, believes, believing, bereave, bereft, best, bestow, bestowed, bet, better, bid, bit, bitten, break, breaking, breaks, Britain, British, broke, broken, brow, brows, browse, browsed, browses, browsing, brutal, brutality, brutalize, brutally, brute, brutish, child, childhood, children, children's, child's, conquer, conquered, conquest, cycle, cylinder, dance, danced, dances, dancing, dead, deal, dealt, death, destroy, destroyed, destruction, die, died, dies, differ, difference, different, difficult, difficulty, dosage, dose, drain, drainage, dried, dries, dry, dryly, dye, dyed, dyeing, dyes, dying, elect, elected, election, elections, electoral, emphases, emphasis, emphasize, emphasized, emphasizes, emphatic, employ, employed, employee, employer, employment, England, English, enjoy, enjoyed, enjoying, enjoyment, enjoys, environment, environmental, environmentalist, establish, established, establishes, establishing, establishment, estimate, estimated, estimates, estimation, fall, fallen, falling, falls, fame, family, famous, fat, feel, feeling, feels, fell, fellow, fellows, fellowship, felt, fight, fighting, fill, find, fit, fold, follow, fond, fought, found, fret, fretting, full, fund, funding, gentleman, gentlemen, German, Germany, good, happen, happened, happening, happens, happier, happiest, happily, happiness, happy, hat, hate, hated, hating, hatred, hats, heavy, heft, held, hell, hill, hillside, hit, hits, hitting, hold, increase, increased, increases, increasing, incredible, incredulity, incredulous, infer, inference, inferior, inferiority, inferred, infers, insist, insisted, insistence, integral, integrally, integrate, integration, interfere, interference, interfering, introduce, introduced, introduces, introduction, introductory, lack, lacking, ladies, ladies', lady, lady's, law, lawyer, leave, leaves, leaving, left, lick, lie, lied, lies, life, life's, lift, lived, lives, living, lock, love, loved, lovely, loves, loving, low, lower, lowest, lowing, lowly, luck, luckily, lucky, lying, maid, maiden, maids, male, man, manage, manipulate, mankind, manner, manners, man's, mansion, mansions, marriage, married, marries, marrow, marry, marrying, men, mend, men's, merry, mile, mind, mole, mule, natural, naturally, nature, nature's, paid, pain, painful, painfully, pay, paying, payment, payments, pays, plan, planned, planning, plans, plant, plantation, play, played, player, players, playing, plays, presidency, president, presidential, presidents, prince, princes, princess, princesses, protect, protection, protector, protectoral, protectorate, ran, rang, range, ranged, ranges, ranging, rank, rant, rate, rated, rates, rating, real, realization, realize, refer, reference, referral, rid, ridden, ride, ring, ringing, rings, road, rode, rude, ruin, ruined, ruins, run, rung, running, runs, Russell, Russia, Russian, sad, sadden, saddened, sadly, said, sang, sat, say, saying, says, secede, seceded, secedes, seceding, secession, secessionist, secessionists, secessions, sell, selling, send, sending, sent, sentence, sentences, sentiment, sentimental, sentiments, set, sets, setting, sing, singing, sit, sits, sitting, slain, slave, slave-holder, slaveholding, slavery, slavery-restricting, slaves, slay, slow, slowly, sold, solution, solutions, solve, solved, solves, son, song, sons, sore, sorely, sorrow, sorry, sort, speak, speaker, speaking, speaks, speech, speeches, spoke, spoken, spokesman, strata, stratum, strife, strifes, strive, strived, striven, strove, sun, sung, swam, swear, swell, swelled, swelling, swells, swim, swimming, swims, swollen, swoon, swooned, swoons, sword, sworn, swum, tall, tame, tamed, tax, taxes, teeth, tell, telling, thank, thanking, thanks, theft, thief, thieve, thieved, thieves, thin, thing, things, think, thinking, thrift, thrifty, thrive, thrived, thriven, thrives, throne, time, times, timing, told, tooth, traffic, trafficking, transfer, transference, transferral, want, wear, weave, weep, weeping, weeps, weft, wept, wet, wife, wife's, winter, wipe, wiped, wipes, wit, wives, woman, woman's, women, women's, word, words, wore, world, worn, woven.

Applying Algorithm 4.15 to the list above, we arrive at the following result

{abolish, abolished, abolishes, abolishing, abolishment, abolition, abolitionism, abolitionist, abolitionists}, {admonish, admonishable, admonisher, admonishes, admonishing, admonishingly, admonishment}, {admonition, admonitory}, {aid, aided, aids}, {America, American}, {analyses, analysis, analytic, analytically, analyze, analyzed}, {ant}, {antenna, antennae}, {anterior}, {anxieties, anxiety, anxious, anxiously}, {arise, arisen, arose}, {automatic, automatically}, {bad}, {bat}, {beauties, beautiful, beautifully, beauty}, {bed, bedding, beds}, {beg, begged}, {began, begin, beginning, begun}, {beggar}, {belie, belied, belies}, {belief, beliefs, believe, believed, believes, believing}, {bereave, bereft}, {best, better, good}, {bestow, bestowed}, {bet}, {bid}, {bit}, {bitten}, {break, breaking, breaks, broke, broken}, {Britain, British}, {brow, brows}, {browse, browsed, browses, browsing}, {brutal, brutality, brutalize, brutally, brute, brutish}, {child, child's, childhood, children, children's}, {conquer, conquered}, {conquest}, {cycle}, {cylinder}, {dance, danced, dances, dancing}, {dead, death, die, died, dies, dying}, {deal, dealt}, {destroy, destroyed, destruction}, {differ, difference, different}, {difficult, difficulty}, {dosage, dose}, {drain, drainage}, {dried, dries, dry, dryly}, {dye, dyed, dyeing, dyes}, {elect, elected, election, elections, electoral}, {emphases, emphasis, emphasize, emphasized, emphasizes, emphatic}, {employ, employed, employee, employer, employment}, {England, English}, {enjoy, enjoyed, enjoying, enjoyment, enjoys}, {environ-

ment, environmental, environmentalist}, {establish, established, establishes, establishing, establishment}, {estimate, estimated, estimates, estimation}, {fall, fallen, falling, falls, fell}, {fame, famous}, {family}, {fat}, {feel, feeling, feels, felt}, {fellow, fellows, fellowship}, {fight, fighting, fought}, {fill}, {find, found}, {fit}, {fold}, {follow}, {fond}, {fret, fretting}, {full}, {fund, funding}, {gentleman, gentlemen}, {German, Germany}, {happen, happened, happening, happens}, {happier, happiest, happily, happiness, happy}, {hat, hats}, {hate, hated, hating, hatred}, {heavy, heft}, {held, hold}, {hell}, {hill}, {hillside}, {hit, hits, hitting}, {increase, increased, increases, increasing}, {incredible}, {incredulity, incredulous}, {infer, inference, inferred, infers}, {inferior, inferiority}, {insist, insisted, insistence}, {integral, integrally, integrate, integration}, {interfere, interference, interfering}, {introduce, introduced, introduces, introduction, introductory}, {lack, lacking}, {ladies, ladies'}, lady, lady's}, {law}, {lawyer}, {leave, leaves, leaving, left}, {lick}, {lie, lied, lies, lying}, {life, life's, lived, lives, living}, {lift}, {lock}, {love, loved, lovely, loves, loving}, {low, lower, lowest, lowing, lowly}, {luck, luckily, lucky}, {maid, maiden, maids}, {male}, {man, man's, men, men's}, {manage}, {manipulate}, {mankind}, {manner, manners}, {mansion, mansions}, {marriage, married, marries, marry, marrying}, {marrow}, {mend}, {merry}, {mile}, {mind}, {mole}, {mule}, {natural, naturally, nature, nature's}, {paid, pay, paying, payment, payments, pays}, {pain, painful, painfully}, {plan, planned, planning, plans}, {plant, plantation}, {play, played, player, players, playing, plays}, {presidency, president, presidential, presidents}, {prince, princes, princess, princesses}, {protect, protection, protector, protectoral, protectorate}, {ran, run, running, runs}, {rang, ring, ringing, rings, rung}, {range, ranged, ranges, ranging}, {rank}, {rant}, {rate, rated, rates, rating}, {real, realization, realize}, {refer, reference, referral}, {rid}, {ridden, ride, rode}, {road}, {rude}, {ruin, ruined, ruins}, {Russell, Russia, Russian}, {sad, sadden, saddened, sadly}, {said, say, saying, says}, {sang, sing, singing, sung}, {sat, sit, sits, sitting}, {secede, seceded, secedes, seceding, secession, secessionist, secessionists, secessions}, {sell, selling, sold}, {send, sending, sent}, {sentence, sentences}, {sentiment, sentimental, sentiments}, {set, sets, setting}, {slain, slay}, {slave, slave-holder, slave-holding, slavery, slavery-restricting, slaves}, {slow, slowly}, {solution, solutions, solve, solved, solves}, {son, sons}, {song}, {sore, sorely, sorry}, {sorrow}, {sort}, {speak, speaker, speaking, speaks, spoke, spoken, spokesman}, {speech, speeches}, {strata, stratum}, {strife, strifes, strive, strived, striven, strove}, {sun}, {swam, swim, swimming, swims, swum}, {swear, sworn}, {swell, swelled, swelling, swells, swollen}, {swoon, swooned, swoons}, {sword}, {tall}, {tame, tamed}, {tax, taxes}, {teeth, tooth}, {tell, telling, told}, {thank, thanking, thanks}, {theft, thief, thieve, thieved, thieves}, {thin}, {thing, things}, {think, thinking}, {thrift, thrifty}, {thrive, thrived, thriven, thrives}, {time, times, timing}, {traffic, trafficking}, {transfer, transference, transferral}, {want}, {wear, wore, worn}, {weave, weft, woven}, {weep, weeping, weeps, wept}, {wet}, {wife, wife's, wives}, {winter}, {wipe, wiped, wipes}, {wit}, {woman, woman's, women, women's}, {word, words}, {world},

which contains only very few errors. For such a list of test words, the time cost of our algorithm is only twice as long as conventional Porter stemming<sup>46</sup>, the latter of which yields the following clustering result:

{abolish, abolished, abolishes, abolishing, abolition}, {abolition}, {abolitionism}, {abolitionist, abolitionists}, {admonish, admonishable, admonisher, admonishes, admonishing, admonishment}, {admonishingly}, {admonition}, {admonitory}, {aid, aided, aids}, {America}, {American}, {analyses}, {analysis}, {analytic, analytically}, {analyze, analyzed}, {ant}, {antenna, antennae}, {anterior}, {anxieties, anxiety}, {anxious}, {anxiously}, {arise}, {arisen}, {arose}, {automatic, automatically}, {bad}, {bat}, {beauties, beautiful, beauty}, {beautifully}, {bed, bedding, beds}, {beg, begged}, {began}, {beggar}, {begin, beginning}, {begun}, {belie, belied, belies}, {belief, beliefs}, {believe, believed, believes, believing}, {bereave}, {bereft}, {best}, {bestow, bestowed}, {bet}, {better}, {bid}, {bit}, {bitten}, {break, breaking, breaks}, {Britain}, {British}, {broke}, {broken}, {brow, brows}, {browse, browsed, browses, browsing}, {brutal, brutality, brutalize, brutally}, {brute}, {brutish}, {child}, {childhood}, {children}, {children's}, {child's}, {conquer, conquered}, {conquest}, {cycle}, {cylinder}, {dance, danced, dances, dancing}, {dead}, {deal}, {dealt}, {death}, {destroy, destroyed}, {destruction}, {die}, {died, dies}, {differ, difference, different}, {difficult}, {difficulty}, {dosage}, {dose}, {drain}, {drainage}, {dried, dries}, {dry}, {dryly}, {dye, dyeing, dyes}, {dyed, dying}, {elect, elected, election, elections}, {electoral}, {emphasize, emphasize, emphasized, emphasizes}, {emphasis}, {emphatic}, {employ, employed}, {employee}, {employer, employment}, {England}, {English}, {enjoy, enjoyed, enjoying, enjoys}, {enjoyment}, {environment}, {environmental}, {environmentalist}, {establish, established, establishes, establishing, establishment}, {estimate, estimated, estimates, estimation}, {fall, falling, falls}, {fallen}, {fame}, {family}, {famous}, {fat}, {feel, feeling, feels}, {fell}, {fellow, fellows}, {fellowship}, {felt}, {fight, fighting}, {fill}, {find}, {fit}, {fold}, {follow}, {fond}, {fought}, {found}, {fret, fretting}, {full}, {fund, funding}, {gentleman}, {gentlemen}, {German}, {Germany}, {good}, {happen, happened, happening, happens}, {happier}, {happiest}, {happily}, {happiness, happy}, {hat, hats}, {hate, hated, hating}, {hatred}, {heavy}, {heft}, {held}, {hell},

<sup>46</sup>We use the Porter stemming functionality built in *Mathematica* v11.0 for this numerical experiment. Certain variations on the Porter stemming algorithm may have better performances in word clustering.

{hill}, {hillside}, {hit, hits, hitting}, {hold}, {increase, increased, increases, increasing}, {incredible}, {incredulity, incredulous}, {infer, inference, inferred, infers}, {inferior, inferiority}, {insist, insisted, insistence}, {integral, integrally, integrate, integration}, {interfere, interference}, {interfering}, {introduce, introduced, introduces}, {introduction}, {introductory}, {lack, lacking}, {ladies, lady}, {ladies'}, {lady's}, {law}, {lawyer}, {leave, leaves, leaving}, {left}, {lick}, {lie}, {lied, lies}, {life}, {life's}, {lift}, {lived, lives, living}, {lock}, {love, loved, lovely, loves, loving}, {low, lowing}, {lower}, {lowest}, {lowly}, {luck}, {luckily}, {lucky}, {lying}, {maid, maids}, {maiden}, {male}, {man}, {manage}, {manipulate}, {mankind}, {manner, manners}, {man's}, {mansion, mansions}, {marriage}, {married, marries, marry, marrying}, {marrow}, {men}, {mend}, {men's}, {merry}, {mile}, {mind}, {mole}, {mule}, {natural, naturally, nature}, {nature's}, {paid}, {pain, painful}, {painfully}, {pay, paying, pays}, {payment, payments}, {plan, planned, planning, plans}, {plant}, {plantation}, {play, played, playing, plays}, {player, players}, {presidency, president, presidents}, {presidential}, {prince, princes}, {princess, princesses}, {protect, protection}, {protector, protectoral, protectorate}, {ran}, {rang, range, ranged, ranges, ranging}, {rank}, {rant}, {rate, rated, rates, rating}, {real}, {realization, realize}, {refer, reference}, {referral}, {rid}, {ridden}, {ride}, {ring, ringing, rings}, {road}, {rode}, {rude}, {ruin, ruined, ruins}, {run, running, runs}, {rung}, {Russell}, {Russia}, {Russian}, {sad}, {sadden, saddened}, {sadly}, {said}, {sang}, {sat}, {say, saying, says}, {secede, seceded, secedes, seceding}, {secession, secessions}, {secessionist, secessionists}, {sell, selling}, {send, sending}, {sent}, {sentence, sentences}, {sentiment, sentimental, sentiments}, {set, sets, setting}, {sing, singing}, {sit, sits, sitting}, {slain}, {slave, slaves}, {slave-holder, slaveholding}, {slavery}, {slavery-restricting}, {slay}, {slow}, {slowly}, {sold}, {solution, solutions}, {solve, solved, solves}, {son, sons}, {song}, {sore, sorely}, {sorrow}, {sorry}, {sort}, {speak, speaking, speaks}, {speaker}, {speech, speeches}, {spoke}, {spoken}, {spokesman}, {strata}, {stratum}, {strife, strifes}, {strive, strived}, {striven}, {strove}, {sun}, {sung}, {swam}, {swear}, {swell, swelled, swelling, swells}, {swim, swimming, swims}, {swollen}, {swoon, swooned, swoons}, {sword}, {sworn}, {swum}, {tall}, {tame, tamed}, {tax, taxes}, {teeth}, {tell, telling}, {thank, thanking, thanks}, {theft}, {thief}, {thieve, thieved, thieves}, {thin}, {thing, things}, {think, thinking}, {thrift}, {thrive, thrived, thrives}, {thriven}, {throve}, {time, times, timing}, {told}, {tooth}, {traffic}, {trafficking}, {transfer, transference}, {transferral}, {want}, {wear}, {weave}, {weep, weeping, weeps}, {west}, {wept}, {wet}, {wife}, {wife's}, {winter}, {wipe, wiped, wipes}, {wit}, {wives}, {woman}, {woman's}, {women}, {women's}, {word, words}, {wore}, {world}, {worn}, {woven},

containing more errors than our output.

*Example 4.15.2.* In Fig. S3, we apply our modified Porter stemming algorithm to topic mining from three English masterpieces (see Table S1 for provenances), and compare the results with ground truths.

According to our stochastic model for recurrence time statistics, word patterns should stay either around the critical line of Poissonian banality (colored gray in Fig. S3) or below it (colored red in Fig. S3). The pattern *Eliza(Ø|beth)(Ø|'*s) in Fig. S3a is a prominent counterexample to this general rule. Our explanation for this is the aliasing of “Elizabeth” by pronouns and nicknames.

Furthermore, if a language lacks gender distinction in third-person pronouns, then its speaker might sometimes spell out a specific person’s name to avoid confusion. This practice also complicates the statistical behavior of *Eliza(Ø|beth)(Ø|'*s) in translations of *Pride and Prejudice*. In Fig. S3e, we see that in four versions (Hungarian, Basque, Korean and Turkish), our statistical algorithm fails to identify the equivalent of *Eliza(Ø|beth)(Ø|'*s) as a topic, and the corresponding word counts also deviate significantly from the English original. These four languages, together with Finnish, do not have separate words for “he” and “she”.<sup>47</sup> The use of pronouns (or possessive suffixes, or personal verb markers, or explicit names) in these five languages can be demonstrated by the following excerpt (see Table S1 for provenances) from Chapter 43 of *Pride and Prejudice*:

#### English

*Mrs. Gardiner was standing a little behind; and on her pausing, he asked her if she would do him the honour of introducing him to her friends.*

#### Finnish

*Gardinerit olivat jääneet vähän taemaksi; ja kun Elizabethilta puhe juutui, pyysi Darcy tulla esitetyksi hänen ystävilleen.*

#### Hungarian

*Mrs. Gardiner kissé hátrább maradt, és amikor Elizabeth elhallgatott, Darcy megkérte, legyen szíves bemutatni őt ismerőseinek.*

<sup>47</sup>This statement is true for colloquial Korean, where 그 means both “he” and “she”. In written Korean texts that are translations of certain European languages, the word 그녀 is sometimes used for “she”. This explains why the Korean word count statistics conform to the English original for *Eliza(Ø|beth)(Ø|'*s) better than Hungarian, Basque and Turkish.

### Basque

Gardiner anderea atzeratxoago zegoen; eta *Elizabeth* isildu zelarik, *Darcy* jaunak galde egin *zion* ea ez al *zion* egingo lagun *horien* aurrean *bera* aurkezteko ohorea.

### Korean

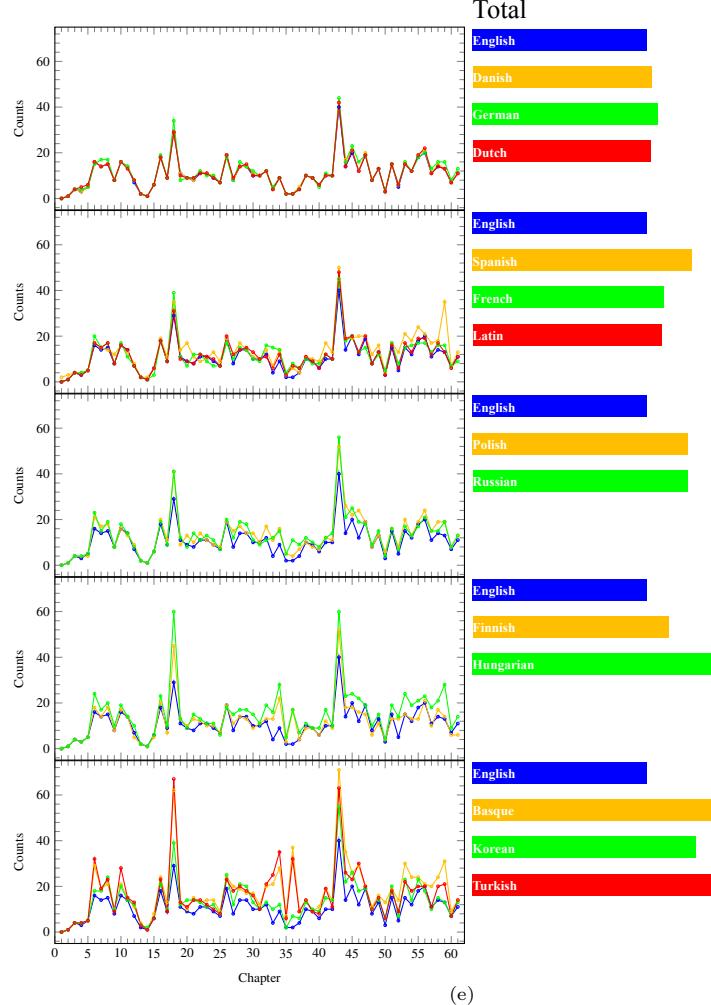
가디너 부인은 조금 뒤에 서 있었다. 엘리자베스가 말을 멈추자 다시 씨는 그녀에게 친구 분들을 소개하는 영광을 베풀어 주지 않겠냐고 물었다.

### Turkish

*Mrs. Gardiner biraz arkada duruyordu; Elizabeth* 'in duraklamasıyla Mr. *Darcy*, arkadaşlarıyla tanışma şerefine erişmeyi rica etti.

Here, in the parallel texts given above, the magenta parts point to Elizabeth herself, either in name or by coreference; the blue parts refer to Mr. Darcy. Agglutinative case prefixes/suffixes (other than those personal markers) are not shown in color. To avoid potential misinterpretations of pronominal coreferences, all the non-Indo-European translations above explicitly mentioned the names Elizabeth and Darcy as antecedents, and the subsequent pronouns (or personal markers) referred to these antecedents, according to their order of appearance in the translated texts.

ELIZABETH	MR DARCY	BENNET MRS	AAE	BINGLEY JANE	LOOKS	COLLINS	LYDIA	LIZZIE	ROOM LADY	DAY
SEE	WICKHAM	INT MISS COME	CHARLOTTE	KITTY KNOW	GODSON	CATHERINE	TAKES	SISTER	THLNK	FMGHT
GARDINER	EXT	CAROLINE	MARKE	LEAV	DRAWING PLEASE	MAN HEAR	LIKE	GIRLS CUT MARRIED	CARRIAGE	TURNS DOOR
FAMILY	LITTLE	NIGHT STAB	LOVE	FOUND LETTER	SIR REALY	GIVE STANDIN	SURE	HOP GREAT	DAUGHTER	TRL MARY CONY
BOW	MOTHER	THREE	THANK	FITZWILLIAM	FATHER COUNTRY	PEMBERLEY	INDRED KIN	WAY PEOPLE	BALL PASSION	STORY
WRITING	STAY	END JOIN	UNCLE	LATE PELL	VILLAGE	MERTON	FIGHTIN	TRYIN	MOMENT WORD	DINTE
PNT	THEIR	EMBARRASSMENT	ESTATE	EMBARRASSMENT	PAPA LOND	HANSDOME	DRESSED	BLAINE	MILITA NAME	CLOSE
WHITE	SHOP	HUNS福德	ENGAGED	ROB	PRESENT	PAINTER	PAINT	MISS	OBSECT	LUCAS COULD
ACCOMPLISHED	POST	PRESENT	V	SOLDIERS	RUNED	RAIN	PLACE	FEAR	RED CORD	REID
BREAKFAST	MEET	CRYS	ATRIBL	POST	POST	POST	POST	POST	POST	PULL FENCE
(c)	ACCOMPLISHED	POST	POST	POST	POST	POST	POST	POST	POST	LAST LIZZ
CHAPTER	INQUIRY	NEIGHBOURHOOD	HURST	PHILLIPS	NIECE	ABSENCE	FANCY	FORTNIGHT	READY	TON
(b)	SISTER	FRIEND	FRIEND	FRIEND	LOS	PERSUAD	EARTLY	LATE	SMALL	REALLY AIR
(d)	SIDE	LOSE	PERSUAD	EARTLY	SMALL	NIGHT	MISTAKES	AIR DIES	WORK	COLL



SAY SEE<sup>w</sup> THOUGHT<sup>w</sup> MR GO<sup>nt</sup> LOOKED<sup>w</sup> LIKE<sup>w</sup> COME<sup>w</sup> KNOW<sup>w</sup> TAKE<sup>883ter</sup>  
 FEELING<sup>w</sup> ROCHESTER<sup>w</sup> MARE<sup>w</sup> LITTLE JANE MISS<sup>w</sup> LONG<sup>w</sup> LIFE<sup>w</sup> SIR EYES DAY<sup>w</sup> LEAV<sup>w</sup>  
 HEARD<sup>w</sup> GIVE<sup>w</sup> SEEMED<sup>w</sup> HAND<sup>w</sup> TIME<sup>w</sup> MRS ASKED LOVE<sup>w</sup> FOUND ROOM<sup>w</sup> WISHED FACE JOHN<sup>w</sup> TWO HOUSE<sup>w</sup> PASSED<sup>w</sup> WORDS  
 STOOD<sup>w</sup> TURNED<sup>w</sup> DEAF<sup>w</sup> FIRST ANSWERED HEART DOOR RASE<sup>w</sup> NEAR<sup>w</sup> LADIES<sup>w</sup> SAT TELL<sup>w</sup> NIGHT STRANGE<sup>w</sup> PUT<sup>w</sup> WAY  
 PRESENT<sup>w</sup> HOUR MAN RETURNED<sup>w</sup> MARRIED<sup>w</sup> THINGS PLACE FIRE<sup>w</sup> CALLED FAIRFAX<sup>w</sup> ST HEAD MASTER<sup>w</sup> ADELE<sup>w</sup> NATURE<sup>w</sup>  
 GREAT CHILD<sup>w</sup> VOICE SPEAK<sup>w</sup> BESSIE<sup>w</sup> MIND READ<sup>w</sup> MOMENT MORNING GENTLEMAN NEW<sup>w</sup> FEAR<sup>w</sup> DRAW<sup>w</sup> ROUND GIRLS LIGHT<sup>w</sup>  
 OLD REED<sup>w</sup> SURE<sup>w</sup> FAR PART<sup>w</sup> WOMAN DARK<sup>w</sup> OH INDEED APPEARED<sup>w</sup> CONTINUED<sup>w</sup> YEARS TALK<sup>w</sup> CERTAIN PLEASE<sup>w</sup> STRONG<sup>w</sup> MEANS<sup>w</sup> HALL FELL<sup>w</sup>  
 BROUGHT<sup>w</sup> SMILE<sup>w</sup> THORNFIELD HOPE COLD<sup>w</sup> EYRE POINT<sup>w</sup> YOUNG<sup>w</sup> MINUTES GOD FOLLOWED CARE<sup>w</sup> MARY<sup>w</sup> BELIEVE<sup>w</sup> SERVANT<sup>w</sup> ENTERED QUIET<sup>w</sup>  
 NAME BED EVENING LAY<sup>w</sup> WALK<sup>w</sup> KEEP<sup>w</sup> ARMS WINDOW<sup>w</sup> SIDE RUN<sup>w</sup> KIND<sup>w</sup> LOW<sup>w</sup> BLACK THREE<sup>w</sup> RIGHT REAL<sup>w</sup> DEEP<sup>w</sup> WANT<sup>w</sup>  
 SENSE<sup>w</sup> SEE<sup>w</sup> HARD<sup>w</sup> FIELD<sup>w</sup> CREED<sup>w</sup> HAPPY<sup>w</sup> HOME FORGET<sup>w</sup> EXPECTED<sup>w</sup> MOVED<sup>w</sup> DRESSED<sup>w</sup> DARE<sup>w</sup> STAY<sup>w</sup> FIST<sup>w</sup> WORLD WIFE TOLD WATCH<sup>w</sup> STEP<sup>w</sup> SISTER<sup>w</sup>  
 QUESTION SHOWED FORMED<sup>w</sup> REST SUPPOSE<sup>w</sup> PLEASURE WILD<sup>w</sup> USED WHITE SEN<sup>w</sup> OPENED WORK LISTENED<sup>w</sup> BEGAN SCHOOL WEND<sup>w</sup> DOUBT SORT DIANA<sup>w</sup>  
 FAST TABLE LARGE<sup>w</sup> MINE BEAR<sup>w</sup> SOUND<sup>w</sup> CONSIDERED<sup>w</sup> OPEN LIPS TEACHER<sup>w</sup> REMEMBER ACT<sup>w</sup> FRIENDS<sup>w</sup> CHANGE REPA<sup>w</sup> BOOK RECEIVED<sup>w</sup> TR<sup>w</sup> EXPRESS<sup>w</sup> PAN<sup>w</sup>  
 SOFT WANTED WAIT<sup>w</sup> ORDER<sup>w</sup> SPOKE PERSON<sup>w</sup> CARRIED<sup>w</sup> MEET<sup>w</sup> HAIR SPIRIT DIRECT<sup>w</sup> GLANCE<sup>w</sup> MATTER HIGH<sup>w</sup> POOR COURSE WONDER OFFERED<sup>w</sup> CLEAR<sup>w</sup> NURSEY SET CHAIR<sup>w</sup>  
 WATER FINE PEOPLE LAUGH<sup>w</sup> STATE<sup>w</sup> FATHER BROKE GRACE<sup>w</sup> HELEN<sup>w</sup> BEAUTIFUL<sup>w</sup> SILENCE SEAT<sup>w</sup> POWER<sup>w</sup> HELP<sup>w</sup> APPROACHED PALE CANDLE<sup>w</sup>  
 AIR PROCESSED<sup>w</sup> NEED<sup>w</sup> BEN<sup>w</sup> TEMPLE<sup>w</sup> MASON<sup>w</sup> ENGLISH CLOSE<sup>w</sup> LATE<sup>w</sup> WEEK IDEAL END TEARS<sup>w</sup> MOTHER<sup>w</sup> BUSINESS TRUE<sup>w</sup> TONE GRIP<sup>w</sup> STOPPED DISTANCE<sup>w</sup>  
 BROCKLEHURST<sup>w</sup> RESOLVED<sup>w</sup> CONTENT<sup>w</sup> SWEET<sup>w</sup> BLOW<sup>w</sup> WALL WHATEVER<sup>w</sup> GAZED<sup>w</sup> CHARACTER<sup>w</sup> SUFFERED<sup>w</sup> DRIVING<sup>w</sup> CONVERSATION TROUBLE PLAIN<sup>w</sup> BROTHER<sup>w</sup> RIVERS<sup>w</sup>  
 LOWWOOD<sup>w</sup> BURN<sup>w</sup> PORT<sup>w</sup> MA<sup>w</sup> AM FEATURES CHAMBER BRIGHT KISSED<sup>w</sup> CR<sup>w</sup> DREAM<sup>w</sup> SHUT FAMILY CURTAIN WANDERED<sup>w</sup> SUDEN PAUSED DIFFERENT<sup>w</sup>  
 SOFT HARD INTEREST<sup>w</sup> EXCITED<sup>w</sup> MUCH<sup>w</sup> LOST EXISTENCE SMALL ROAD GRAVE FLOWER SOUL GROUND<sup>w</sup> EARL CHAPTER SECOND<sup>w</sup> CID<sup>w</sup> PREPARED<sup>w</sup> DISCOVERED<sup>w</sup> UNDERSTAND<sup>w</sup> NEED USUAL<sup>w</sup>  
 SLEEP SCARCELY MONT<sup>w</sup> MANNER FILLED EXCUSE<sup>w</sup> EAS CASE<sup>w</sup> CROSS<sup>w</sup> TASTE PICTURE<sup>w</sup> LIP<sup>w</sup> SILENT GATHERED<sup>w</sup> WOOD<sup>w</sup> FORCES GLAD TRAVELED<sup>w</sup> COVERED UTTERED<sup>w</sup> POSSESS<sup>w</sup> JUDGEMENT<sup>w</sup> DELIGHT<sup>w</sup>  
 EASY HUMAN WHISPERS HATE<sup>w</sup> REPEATED<sup>w</sup> ADDRESS<sup>w</sup> REACHES<sup>w</sup> FRENCH EFFORT STRUCK GATESHEAD ADDED RESPECTED<sup>w</sup> SHANG<sup>w</sup> RING<sup>w</sup> LEARN<sup>w</sup> COMFORT<sup>w</sup> STRENG<sup>w</sup> SAVE PLAY<sup>w</sup>  
 GOVERNESS<sup>w</sup> TREES TALL PROBABLY<sup>w</sup> CLOUD BREAD DISTING<sup>w</sup> SUCCESS<sup>w</sup> DRIVE<sup>w</sup> VAN<sup>w</sup> P<sup>w</sup> EXT<sup>w</sup> PROMISE NOTICE LETTER AFTERNOON TO-NIGHT SKY CHURCH BOOND<sup>w</sup>  
 GOVERNESS<sup>w</sup> PRAYER<sup>w</sup> EXPERIENCE DREAD<sup>w</sup> DESIRE CHARM FAIR SURPRISE SCENE GENERAL<sup>w</sup> HANNAH<sup>w</sup> RELATIVES<sup>w</sup> NOTES INTENDED<sup>w</sup> REASON<sup>w</sup> RAIN PUPIL<sup>w</sup> PASSION<sup>w</sup> FRAME DIFFICULT<sup>w</sup> DEMAND<sup>w</sup>  
 MILES GLASS FIELD<sup>w</sup> CHEEK PAST HANDSOME THANK<sup>w</sup> HABITS<sup>w</sup> SA<sup>w</sup> HORSE<sup>w</sup> REGARD<sup>w</sup> GATE FREE<sup>w</sup> PRE<sup>w</sup> FORWARD CIRCUMSTANCE<sup>w</sup> AID<sup>w</sup> TEN MOUNTAINS<sup>w</sup> ARRANGED<sup>w</sup> WAB<sup>w</sup> TRACE<sup>w</sup> FIXED FIGURE DIM<sup>w</sup>  
 BREATH FANCY ALLOWED<sup>w</sup> SITTING FIVE CLOSED BLOOD WAKE<sup>w</sup> REFLECTED<sup>w</sup> AFFECTION OBJECT LOCATED<sup>w</sup> HAPPENED EXPLO<sup>w</sup> EIGHT<sup>w</sup> TRUTH RICH<sup>w</sup> CITY<sup>w</sup> MOOR<sup>w</sup> HUSBAND HILLS DAUGHTER<sup>w</sup>  
 SING ENJOYED<sup>w</sup> SANK OCCASION<sup>w</sup> LONELY EXAMINED<sup>w</sup> DEAR TRED<sup>w</sup> FRESH ABSOLUTELY SHOULD<sup>w</sup> BIRD AFRAID VISIT<sup>w</sup> SICK<sup>w</sup> QUICK<sup>w</sup> STONE<sup>w</sup> KITCHEN FAITHFUL<sup>w</sup> BREAKFAST<sup>w</sup> TASK<sup>w</sup>  
 SUN AGE FOUR DEAL BAD SPENSE OBLIGED<sup>w</sup> DEPARTED ATTENTION ADVANCED TWENTY RECALLED<sup>w</sup> BITTER SUBJECT PURE KNEE FINGER COUNTEEN<sup>w</sup> SCHOOLROOM MED GALLERY SIGN<sup>w</sup> METAS<sup>w</sup> STUD<sup>w</sup> REMARK<sup>w</sup> PROVE LINE<sup>w</sup>  
 LEADS GUARD<sup>w</sup> CURLS<sup>w</sup> COOL COMPLETE WITHDRAW REMOVED<sup>w</sup> INFORM EQUAL<sup>w</sup> CAPBY BLW<sup>w</sup> ARRIVED SETTLED POOLE<sup>w</sup> HEAVEN BURNS<sup>w</sup> ACCOUNT SIGHT NARROW BLUE STRATEG<sup>w</sup> SLIP<sup>w</sup> SATISFY<sup>w</sup> PERIOD<sup>w</sup>  
 EVIDENTLY ACCOMPAN<sup>w</sup> SECRET SHEET DANGER<sup>w</sup> IL<sup>w</sup> GOLD<sup>w</sup> ELIZA<sup>w</sup> EARTH CLEAN<sup>w</sup> BREAD<sup>w</sup> TEA READY MONEY GARDEN FOREHEAD FLESH DINNER<sup>w</sup> WRONG SYMPATH<sup>w</sup> BED<sup>w</sup> QUIT<sup>w</sup> MARK<sup>w</sup> FREQUENT ENDURE<sup>w</sup> OP<sup>w</sup>  
 CHANCE<sup>w</sup> BONNET TOWN OLIVER<sup>w</sup> ESHTON<sup>w</sup> BREAK<sup>w</sup> BELL HEARTH PRACT<sup>w</sup> THROW KEEN<sup>w</sup> G<sup>w</sup> FORGIVE<sup>w</sup> FIR<sup>w</sup> ATTEMPT<sup>w</sup> VIEW UNCLE<sup>w</sup> STATE<sup>w</sup> POSSIBLY MURK<sup>w</sup> GRIEV<sup>w</sup> DESCENDING CLASS<sup>w</sup>  
 FAULT<sup>w</sup> DUTY<sup>w</sup> CUT<sup>w</sup> MILLCOTE<sup>w</sup> GREEN FUTURE<sup>w</sup> DEPEND<sup>w</sup> WE<sup>w</sup> SUIT<sup>w</sup> HEALTH<sup>w</sup> BANK<sup>w</sup> WICKED<sup>w</sup> SWEEP<sup>w</sup> RUSH<sup>w</sup> REQUIRED<sup>w</sup> CONDUCT<sup>w</sup> COMPANION<sup>w</sup> CLAUD WINTER SHAPE PERFECT<sup>w</sup> PARTY<sup>w</sup> MAD<sup>w</sup> HUR<sup>w</sup> GLOOM<sup>w</sup>  
 FINISHED<sup>w</sup> FETCH<sup>w</sup> DISPOSE<sup>w</sup> UPSTAIRS DECIDE<sup>w</sup> TREAT<sup>w</sup> SEPARATE<sup>w</sup> IMPRESSION<sup>w</sup> GLEAM<sup>w</sup> FOLD DEAR<sup>w</sup> ADMIRE<sup>w</sup> VISION<sup>w</sup> REQUIRES<sup>w</sup> EFFECT<sup>w</sup> CREATURE CONNECT<sup>w</sup> CLOTHES<sup>w</sup> ADMIT<sup>w</sup> TRUST<sup>w</sup> STEIN<sup>w</sup> SPREAD<sup>w</sup> INSTANCE FEVER<sup>w</sup> DOG<sup>w</sup>  
 CONFIDENCE WET TOMORROW THEN LANGUAGE COMPANY BLIND STOREY<sup>w</sup> RENDER<sup>w</sup> WAYSIDE SECURE<sup>w</sup> HID<sup>w</sup> FURNITURE EMPLOYEE<sup>w</sup> ACCORDING<sup>w</sup> YIELD<sup>w</sup> TIDE SUMMERS<sup>w</sup> SOLAN<sup>w</sup> SHADE<sup>w</sup> EXACT<sup>w</sup> ENERGY<sup>w</sup> EAGER<sup>w</sup> CHRISTIAN<sup>w</sup> WIFE WEALTH THOUSAND<sup>w</sup>  
 BELIEVED PRETTY POSITION OBEYED INTERVAL IMMEDIATE DRAWING ROOM<sup>w</sup> SOPHIE PARLOUR<sup>w</sup> LEAH<sup>w</sup> SINGER<sup>w</sup> SUGGEST<sup>w</sup> COMPARE<sup>w</sup> BRIDE<sup>w</sup> WORTH<sup>w</sup> SHADOW<sup>w</sup> SEVER<sup>w</sup> FRON<sup>w</sup> PECCULAR<sup>w</sup> JOY<sup>w</sup> EVENT<sup>w</sup> CURIO<sup>w</sup>  
 CONSEQUENT<sup>w</sup> COMPREHEND COMMENCE<sup>w</sup> TOP BEHIND VEIL<sup>w</sup> THEM SOCIETY PATH MORTON LIBRARY LIBERTY IMPOSSIBLE EARLY SHAD<sup>w</sup> CLOAK<sup>w</sup> SIX COMMUNICAT<sup>w</sup> RECOLLECT<sup>w</sup> CAST ACCEPT<sup>w</sup> THICK SHOCK<sup>w</sup> PROVIDED<sup>w</sup> PAPER<sup>w</sup> MENTION<sup>w</sup>  
 EARNEST<sup>w</sup> COUSING SERVICE<sup>w</sup> PROSPECT<sup>w</sup> PRINCIPLE<sup>w</sup> INFLUENCE HUNGRY FAIR<sup>w</sup> CLERGYMAN<sup>w</sup> APARTMENT ACQUAINT<sup>w</sup> SOLO<sup>w</sup> SIGHT<sup>w</sup> PILOT<sup>w</sup> INSTRUMENT<sup>w</sup> FAVOUR<sup>w</sup> ELDER<sup>w</sup> WEAK<sup>w</sup> STRENG<sup>w</sup> STRIKE<sup>w</sup> PIECE<sup>w</sup> MYSTER<sup>w</sup>  
 MISSIONARY<sup>w</sup> CLO<sup>w</sup> FOOD ACCOMPLISH<sup>w</sup> SILK<sup>w</sup> PLAN<sup>w</sup> ESPECIALLY EDWARD<sup>w</sup> CHARGE BRILLIANT ABSENC<sup>w</sup> MOUTH MILLER<sup>w</sup> LIT GREY COMPLY<sup>w</sup> REPORT<sup>w</sup> DECLAR<sup>w</sup> SINCER<sup>w</sup> LABOUR<sup>w</sup>  
 CAUSES<sup>w</sup> CAUSE<sup>w</sup> SACRIFICE REQUEST<sup>w</sup> HAD<sup>w</sup> PURSU<sup>w</sup> PEACE<sup>w</sup> BRAINS<sup>w</sup> DEGREE BRAIN<sup>w</sup> BLANCHE<sup>w</sup> ANGEL STAIRCASE HENY<sup>w</sup> AUNT<sup>w</sup> AH TOL<sup>w</sup> BAND<sup>w</sup> PRACT<sup>w</sup> SHARE<sup>w</sup> REB<sup>w</sup> PUZZLED MEAN<sup>w</sup> INCLIN<sup>w</sup> CRIM<sup>w</sup> FAIL<sup>w</sup> PERFORM<sup>w</sup>  
 GRAPHIC FRAM<sup>w</sup> FRAM<sup>w</sup> WEARY<sup>w</sup> SLOWLY SINGLE<sup>w</sup> SHRED<sup>w</sup> SENTIMENT<sup>w</sup> MATERIAL<sup>w</sup> PURPOSE<sup>w</sup> LOC<sup>w</sup> CLO<sup>w</sup> KEY<sup>w</sup> DROM<sup>w</sup> HEATH<sup>w</sup> EAST<sup>w</sup> CUP<sup>w</sup> CONSENT<sup>w</sup> CASTEL<sup>w</sup> ASKED<sup>w</sup> ROOF LYNN JANET FULL COUNTRY COACH<sup>w</sup> SOAR<sup>w</sup> BENEFIT<sup>w</sup> HOM<sup>w</sup> JOURNEY<sup>w</sup>  
 (a')  
 JANE ROCHESTER MRS JOHN ST FAIRFAX<sup>w</sup> CONT'D LOOK<sup>w</sup> SEE<sup>w</sup> INT DAY THORNFIELD GO<sup>w</sup> COME ADELE<sup>w</sup>  
 REED<sup>w</sup> LIKE ROOM<sup>w</sup> MISS TAKE<sup>w</sup> EYES DOOR<sup>w</sup> SIR HAND<sup>w</sup> NIGHT HOUSE TURNS<sup>w</sup> KNOW<sup>w</sup> LOVE<sup>w</sup> HEAR<sup>w</sup> EYRE DEAF<sup>w</sup> GOOD<sup>w</sup> MR LIFE<sup>w</sup> MARY HELEN<sup>w</sup> FEND<sup>w</sup> SKY<sup>w</sup> OPEN<sup>w</sup> DIANA<sup>w</sup>  
 BLANCHE<sup>w</sup> BESSIE<sup>w</sup> LAUGH<sup>w</sup>  
 MASON<sup>w</sup> MAN DARK STAND<sup>w</sup> PLEASE<sup>w</sup> GOD<sup>w</sup> GIVE<sup>w</sup> TELL<sup>w</sup> RED GRACE<sup>w</sup> EVENING BRIGGS<sup>w</sup> STARS<sup>w</sup> FULL<sup>w</sup> BURNS<sup>w</sup> WINDOW<sup>w</sup> CHILD<sup>w</sup> OH WATCH<sup>w</sup> WOMAN<sup>w</sup> HEAD<sup>w</sup> CLO<sup>w</sup>  
 BROCKLEHURST<sup>w</sup> SIDE DREN<sup>w</sup> SMILE<sup>w</sup> MARI<sup>w</sup> WORK<sup>w</sup> BERTHA<sup>w</sup> NEED<sup>w</sup> O HORSE<sup>w</sup> FLOOR<sup>w</sup> CUT<sup>w</sup> GESTURE<sup>w</sup> BEAUTY<sup>w</sup> V PASS<sup>w</sup> SIGHT<sup>w</sup> RAIN<sup>w</sup> CRY<sup>w</sup> THE HOME SLOW<sup>w</sup> SOUL<sup>w</sup> MEAN<sup>w</sup> LEV<sup>w</sup> STONES<sup>w</sup> WORD<sup>w</sup>  
 VOM<sup>w</sup> WEAT<sup>w</sup> HUM<sup>w</sup> SOUL<sup>w</sup> MOTHER<sup>w</sup> BROTHER<sup>w</sup> CURTAINS<sup>w</sup> SOUND<sup>w</sup> SMALL POOLE LOWWOOD<sup>w</sup> SET<sup>w</sup> THRO<sup>w</sup> LADS<sup>w</sup> BEAUTY<sup>w</sup> V PASS<sup>w</sup> SIGHT<sup>w</sup> RAIN<sup>w</sup> CRY<sup>w</sup> THE HOME SLOW<sup>w</sup> SOUL<sup>w</sup> MEAN<sup>w</sup> LEV<sup>w</sup> STONES<sup>w</sup> WORD<sup>w</sup>  
 HANNAH<sup>w</sup> KISSES<sup>w</sup> LEAH<sup>w</sup> EXPRESSION STRONG<sup>w</sup> BOOK<sup>w</sup> AUNT<sup>w</sup> WALL<sup>w</sup> SCREAM<sup>w</sup> LAND<sup>w</sup> POND<sup>w</sup> LIGH<sup>w</sup> STARS<sup>w</sup> PART<sup>w</sup> MARTHA<sup>w</sup> HATE<sup>w</sup> FIGURE COACH<sup>w</sup> TOUCH<sup>w</sup> RESIST<sup>w</sup> STRETCH<sup>w</sup> EMBRACE<sup>w</sup> CAR<sup>w</sup> ADO<sup>w</sup> UNCLE<sup>w</sup> STRONG<sup>w</sup> HOLD<sup>w</sup> YEARS LAUGH<sup>w</sup>  
 (b')  
 GEORGIANA DISTING<sup>w</sup> SUCCEED<sup>w</sup> AFTERNOON<sup>w</sup> TO-NIGHT<sup>w</sup> SAID<sup>w</sup> SEE<sup>w</sup> MARE<sup>w</sup> STRONG<sup>w</sup> ENTERED<sup>w</sup> MINE BEAR<sup>w</sup> EXPRESS<sup>w</sup> OFFER<sup>w</sup>  
 GENERAL<sup>w</sup> DIFFICULT<sup>w</sup> CIRCUMSTANCES AID<sup>w</sup> OCCASION<sup>w</sup> BLUE STREED<sup>w</sup> ACCOMPAN<sup>w</sup> ELIZA<sup>w</sup> ENDRE<sup>w</sup> SET<sup>w</sup> NEED<sup>w</sup> FREE<sup>w</sup> FORWARD<sup>w</sup> BREATH<sup>w</sup> KITCHEN<sup>w</sup> RED<sup>w</sup>  
 OLIVER<sup>w</sup> BELL MILLCOTE DENT<sup>w</sup> DISPOSED<sup>w</sup> SEPARATE<sup>w</sup> SPREAD<sup>w</sup> CONFIDENCE WET LANGUAGE STOREY<sup>w</sup>  
 (c')  
 (d')  
 Fig. S3. Text mining in English. (Continued) A similar service on Charlotte Brontë's *Jane Eyre*. Moira Buffini's screenplay adaptation is chosen as the ground truth for topic extraction.

**SPECIES FORMATION** NATURAL DIFFERENCES CASES VARIETIES SELECTION  
 LIFE GREAT VARIATIONS ORGANS PLANTS PARTS ANIMALS MODIFICATION TWO SEE  
 PRODUCE PON DISTINCT CHARACTERS GENERA CONDITIONS LENGTH GENERALLY NUMBER PERIOD CERTAIN  
 LARGE STRUCTURE GIVEN FACTS GROUPS EXISTING FOUND INDIVIDUALS BELIEVE CHANGE DESCENDENTS FIRST COMMON  
 IMPORTANT NEW TIME HIGHEST CLASSES INSTANCE BESIDE CLOSELY STATE RELATION INHABITANTS DEGREE  
 CROSSED SHOW ACTION BREEDS KNOWN PERFECT WIDELY PRESENT NEARLY BIRDS SLIGHT PROBABLY ISLANDS  
 DEVELOPED DOMESTICATION USEFUL CONSIDER REMARKS FAR APPEAR TENDENCY HABITS SUCCESSIVE MARE DIFFICULTY  
**HYBRIDS** MANNER MR FLOWERS FERTILITY VIEW EXTINCTION SEEDS EFFECTS INHERITANCE SUPPOSED INSECTS

Doubt COUNTRY STERILITY CAUSE SEEMS LIKE KINDS FORMERLY INTERMEDIATE MEANS PLACE INSTINCTS  
 INCREASE WORLD SMALL WORKS NAMELY AMERICA ALLIED EXTREMELY LITTLE OCCUR MAN COME GEOLOGICAL THEORY PARENTS BEINGS ADAPTED  
 SIMILAR POINTS SLOWLY THREE SOUTHERN AREAS SUBJECT LOOK GENERATIONS TAKE EXTEN LOWE EARLY YOUNG ACCORDING ORDER LATE OLD  
 SINGLE FAMILIES POWER BELONGING FOLLOWING PRINCIPLE OCCASIONALLY ADMITTED RESPECT AMOUNT STRONG OBSERVER EUROPE CONTINUOUS ADVANTAGE LAWS  
 PRESERVED OFFSPRING PROCESS MALES SEXUAL FAVOURABLE HAND NECESSARY RANKED NORTHERN UNDERSTAND POLLEN RANGE  
 POSSIBLE YEARS SEPARATED LATELY STAGES CONCLUDING SIZE CREATE DISTANT REASON CLIMATE SPECIAL TRUE PROGENITOR LAND DIVERSITY SAY  
 ORIGINAL IMPROVED COLOURED OBJECTIONS COMPARED SIDES LEAD CONTINENT PLAINLY ORGANISATION AGE BEAR GROWTH GRADUATIONS CALLED CONSEQUENTLY  
 EXPLANATION RARE ANCIENT ANALOGOUS WATER EGGS ACCUMULATIVES REMAINS DIRECT SEA RESULT MAMMALS GO STAND THINN COMPETITION PECULIAR CHAPTER VARY

**TREES NEST** BEES STRUGGLE NUMEROUS RESEMBLANCE POSSESSIVE MIGRATION SYSTEM RACES MEMBERS LINES EYES EVIDENCE CONNECTED  
 REGIONS COMPLEX SECONDARY WINGS FISHES RUDIMENTARY ACQUIRED RULE RECENTLY MARKED DESTRUCTION LATENT PHYSICAL DAY THROUGHOUT UNITED CONFINED FEMALEs  
 FITTED SIMPLE DEPEND PURPOSE ASSURE BRANCHES RENDERED MOUNTAINS BENEFICIAL INTERVALS CLEARLY BODY INDEPENDENTLY APPARENTLY FEET PIGEONS AFFINITIES LARVAE  
 HEAD CHANCE COURSE SUFFICIENT INCLUDE FOOD DUE SERVE LINKS AUTHORITY SUDDENLY DISCOVERED VALUE STOCK BEAK PASSAGE FREQUENTLY REPRODUCTIVE FOSSILS ORDINARY  
 DR REPRESENT MIND LEGS GLACIAL EXPERIENCING WONDERFUL RESEMBLE WAY RECORD FREELY THICKNESS EXCEPTION ACCOUNT ABLE CARRIER DEPOSITOR AFFECTION CAREFULLY SHORT  
 EMBRYO INFER CONSTITUTION BEAUTIFUL ESPECIALLY EXTERMINATE COMPLETELY HORSE SURVIVE SURPRISE PROPORTION STRIKING LESSER STEPS UNDERGO DETERMINES WILD JUDGES  
 PROVED FINALLY DISCUSSED TRANSITIONAL ELEMENTS EASILY REGARD EXPOSE EXPECT DISTRIBUTION ASSUME FUNCTION ADVANCED LOST TRANSPORT RETAINED SPACE SERIES REAL AUSTRALIA ADULT PLAY  
**CELLS** RACES MEMBERS LINES EYES EVIDENCE CONNECTED

PREY STRICTLY SPECIFIC DISEASE DISAPPEAR COLLECTOR YIELD CRUSTACEANS REDUCED RAISED UNIVERSAL SURFACE PALAEONTOLOGISTS DEEP TREES INTERCROSSING DEFINITE CORRELATION ASKED  
 BEDS WHOLLY RISE BIRCHEN ARISE TRANSMITTED THOUSAND SHAPE DOGS OWING GEOGRAPHICAL EXCLUSIVELY ARRANGEMENT TYPE FLYING EXPLAIN EARTH ATTRIBUTED ARCHIPELAGO ACTUALLY TERRESTRIAL FORCE BRIGHT DIVISION

SPREAD OPPOSITE OCEAN BAT WAX PROFESSOR TWENTY EXTRAORDINARY OCEANIC DOMINANT CATTLE ILLUSTRATIS PREVENT UNIFORM ROUND EXTERNAL CURIOUS CONSTRUCTORS ATTEMPTED  
 MOVEMENT APPLIED SUPPLANT SEVERE LIMIT GRADUATED UTTERLY TAIL SCALE JAWS RATE OPEN ALLOWED ENTER ENDURE INVARIABLE SHELLS ISOLATOR BRITAN TEMPERATE DIAGRAM FLAMES EXPRESSO

TREAT BASE TURN SURROUNDING STIGMA ROCKS POSITION AVERAGE MATURE INJURIOUS HOLD PARALLEL MONGRELS MAIN HOMOLOGOUS TERTIARY NEVERTHELESS  
 MARINE FUEL SENSE SIMULTANEOUSLY QUESTION ENORMOUS ENGLAND END CONVERTED COMMUNITY ABSOLUTELY VAST MILES F PROFIT SUPPORT PARTICULAR IGNORANT AGENCY STRANGE LIGHT FOUR STAGED COMPARISON INFINITE CHECK ARGUMENT ANSWER

GARTNER FRUIT AFRICA RIGHT DRAW EXAMINED DESCRIBED TEETH SIR PURE OBVIOUS MOUTH MIVART HAVING LEVEL INDIA FUTURE EXACTLY ELEMENT SOLELY REQUIRED LASTLY  
 CENTRAL RAPID DUCK ABSURD ABDUDE DE SEIZZ INFORMATIVE GROUND STAMENS NATIVE MUTUAL HOME ABSENCE USELESS FOREGOING ANTS VERTEBRATES IMPLANT MASTERS

THROWS METAMORPHOSIS ARRIVED SPHERES AL QUICKLY OFFER ECONOMY COMMENCEMENT BARREN INHABITANT ROCK-PIGEON FAUNAS CONSTANT COLD TRILING SIX FRESH-WATER SUBSIDE REPRESENTATIVE  
 ARCTIC ACCIDENTAL REACH IMMIGRATES VEGETABLE TRACE TELE INFLUENCE DETAIL SHORES QUADRUPEDS LIMBS INTERESTING HEMISPHERE ABRUPT SHEEP HOOKER BEGIN EXCERPTS REPORTS ALTERED WHALE

PROTECTIVE CULTIVATED BOTANISTS UNCONSCIOUS REMEMBER PISTILS MULTITUDE INSISTED HIVE-BEE HEAD FEATHERS DEVIATIONS INNUMERABLE FURNISHED CAPABLE REPEATEDLY SPECIESS NEUTERS  
 ENDOWED EMBRYOLOGICAL DISSEMBLAR DISPERALS DEGRADATION CONCERN BUILDING ATTENDED REPTILES RED NEED MATTER MADEIRA INEXPLICABLE FRESH DURATION DOCTRINE

SUFFER PROGRESS MEASURED KEEP COMB SUBORDINATE SET REVERSION IMPOSSIBLE GRAFTER FLOATER ZEALAND FLIGHT DIVERGENCE WEATH MOTHER TUMBLER DECIDE REAPPEAR

**ILLEGITIMATE LAMELLAE** HABITUALLY ENTRE DR DIVERSITY CONVINCED CIRRIPEDES AVAILABLE TRUTH SEDIMENT INCIDENT GOVERNORS STUDS HAPPEN DEFINED UNUSUAL UNTHINKABLE THIRD  
 SAMPLE HORNS COVERED WELL-MARKED STICK SO-CALLED SERVICE PARENT-SPECIES FACILITY AIR PROVIDED SPEAK RETURN INTRODUCE IMAGINE NERVES FEED CAVES WALL TEN SPECIMENS RECIPROCAL QUALITY PERMANENT

MULLER HORNS HATCHED ESSENTIAL EIGHT CAPACITY BEETLES ALLEGED WHITE WEIGHT SECRETORY REJECT READERS FALSE EXCRETE DECREASE AFFORD REARED MASTERS DISPUTED ATTACHED REGULARLY PAIR LOCAL INDEPENSABLE IMPROBABLE HEIGHT

EQUATORIAL SKIN LONG-CONTINUED LAPSE FIVE CANDOLE BIRTH RECEIVE POLYPODIFLIC MANIFEST SPOT RUDIMENTARY MONSTROUS METAMORPHOSIS LONG ELECTRIC DISPERSAL CAVE BRITAIN BROUGHT AREA APPENDAGES ANALOGIES

ENDOWED EMBRYOLOGICAL DISSEMBLAR DISPERALS DEGRADATION CONCERN BUILDING ATTENDED REPTILES RED NEED MATTER MADEIRA INEXPLICABLE FRESH DURATION DOCTRINE

(a'')

**VARIETIES PLANTS SPECIES FORMATION ORGANS ANIMALS**  
 NATURAL MR INSECTS FLOWERS DEVELOPMENT SELECTION GENERALLY  
 CROSSED CHARACTERS PARTS PRODUCTS DOMESTICATION CLASSIFICATION SEEDS  
 BIRDS STRUCTURE GROUPS INSTINCTS DR FOSSIL VARIATION PROF RELATIONS HYBRIDS  
 ROCKS LIFE COLOUR HABITS CHAPTER DISTRIBUTION CHANGES SUCCESSION GREAT STERILITY  
 GLACIERS GENES YOUNG LOWER CALLED INCLUDING CONDITIONS BODY MODIFICATION EXTINCTION USUALLY  
 ORIGIN IMPORTANCE EYES CRUSTACEANS WINGS EFFECTS DIVISION ORDER GEOLOGICAL FISHES GIVE FERTILITY SHELLS  
 SEXUAL PERIOD FEET BEES RUDIMENTARY ISLANDS WATER APPEARING MEANS CERTAIN AFFINITIES DESCENDENTS INHABITANTS  
 CAUSES IMPERFECT FLORA DISTINCT BELONGING TERM SMALL SIR EGGS COMMON TWO SYSTEM M SECONDARY COMPOSED  
 STRIPED LAWS NEW FRESH-WATER FIRST DIFFERENT BLIND PERFECT STATE REPRODUCTION USED SEE INDIVIDUALS HOMOLOGOUS EARLIEST  
 DESTROYED TEETH PRESENT JAWS EXISTENCE EARTH STRUGGLE PECULIAR MAMMALS EXAMPLES BEINGS APPLIED HIGHEST RANGE CORRELATION  
 CONCLUSIONS STAGE SKIN POWER OLDEST MAN LEAVES FURNISHED DENUDATION MADEIRA WIDELY OBJECTIONS CENTRE ALLIED SUDDEN  
 PALAEONTOLOGY ORGANISMS ON MALE LARGE JOINTED GROWTH CELLS AUSTRALIA ANCIENT AGE ACCLIMATISATION ZEALAND UNITED MAMMALIA INTERMEDIATE  
 INCREASE NERVOUS FILIES TYPES RESEMBLANCE QUADRUPEDS POLLEN PISTIL LAND-SHELLS INHERITANCE GILLS EMBRYOLOGICAL CLOSE CLIMATE CASE ATTACHED  
 AMERICA TREES TIME THEORY LIKE FACTS CATTLE TRANSPORTING POZOZO NESTS MEMBRANES KNOWN BEARING TAIL RATE NEUTER NEARLY HORN HAIR FAVOURABLE EMBRYO  
 DOG DIFFICULTIES CRUSTACEA BREEDS WORLD UMBELLIFERAE THROUGHOUT TERTIARY SILURIAN MATERIAL HEAD GRAFTS ANTS ACTION SEPARATE POSSESSIONS VEGETABLE FRUIT  
 ADAPTED ACQUIRED SOUTH SIMILAR GEOPOLITICAL PORTION NUMBER NAME MONGRELS METAMORPHOSIS LONG LABOUR INTELLIGENCE FALL SURE SEASON OWN MILLION MARSUPIALS KINGDOM HERMAPHRODITES  
 W SURFACE SUMMARY ST SEGMENTS MOUTH J GEOGRAPHICAL BEETLES ARCHIPELAGO PASSAGE SCALP TRANSITION THICKNESS STUDY REMAINS REDUCED RATIO MARE LAURENTIAN FILAMENT  
 DIFFERENCES DEPOSITION EXTRAVERTED VERTEBRATES SUPPORTED SUBJECT STALK SLIGHT POSSIBLE PETALS PARENT PAIR MOTHER LAND INDIAN HORSES EUROPE COUNTRY BEAK BASE AFFECTING  
 Tarsi STIGMA STATES STAMENS SINGLE SAID MUTUAL MINUTE MARSUPIALS MALAY LIMBS FIGHTING FAB CUCKOO CAPABLE C ABSENCE UNDERGOING PUBLISH PROCESSES MODELED LEGS INTERCROSSING FUNCTION  
 ENVELOPE CONSIST HABIT BONE WOLF STEM SIDES SERVING SERIES LINES KINGDOM IMPROVEMENT HERMAPHRODITE FOUND FEMALE EPOCH CONTAINING COLLECTBNS BUTTERFLIES ARTICULATE APPARENTLY  
 WAX PROFESSOR TWENTY EXTRAORDINARY UNIFORM ROUND CURIOUS ATTEMPTED TEMPERATE DIAGRAM  
 LOOK RANKED UNDERSTAND PROGENITOR PLAINLY THINK MIGRATION EVIDENCE MARKED  
 BRANCHES RENDERED INTERVALS DISCOVERED MIND WONDERFUL FREELY EXCEPTION LESSER FINALLY EXPOSED  
 ASSUMED PLAY KEPT REMOTE DISTRICTS PREY YIELD RAISED THREW ASKED ARISE TRANSMITTED THOUSAND SHAPE

AVERAGE END VAST MILES F STRANGE INFINITE ANSWER AFRICA RIGHT PURE SEIZE INFORMATIVE NATIVE USELESS IMPLANT SUBSIDENCE  
 HUNDRED QUICKLY BARBES ROCK-PIGEON SIX DERIVED ACCIDENTAL IMMIGRATES TELL INTERESTING BEGIN  
 REPRESENTATIVE BOTANISTS REMEMBER MULTITUDE INSISTED SPECIES SUFFER KEEP COMB SET FLOTTER DECIDE

**ILLEGITIMATE** DR DRY TRUTH ENEMIES FIXED ADJOINING TEN SPECIMENS QUALITY PERMANENT  
 WELL-MARKED PARENT-SPECIES FACILITY SPEAK DIAGNOSIS BOND SERIOUS QUARTER PERISH HATCHED EIGHT WEIGHT REJECT EXCRETE  
 DECREASE REARED INDEPENSABLE IMPROBABLE HEIGHT EQUATORIAL LONG-CONTINUED SUR SEASON MILLION BUILD INEXPLICABLE

(c'')  
 GREAT GIVE IN FACTS INSTANCE SHOW CONSIDER REMARK FAR SEEMS NAMELY  
 ACCORDING NECESSARY CONCLUSION REASON CONSEQUENTLY EXPLANATION NUMEROUS LATTER CLEARLY INDEPENDENTLY APPARENTLY COURSE WAY ACCOUNT  
 EDITORIAL STRIKING PHOTO EVIDENCE HYPOTHESIS CHIMERIC WITHIN BIOCENOSIS DIVERSITY INFLUENCE DETAILS BEING CAPTURED VICTIMS  
 ENABLED LIFE EXAMINED OBVIOUS MAINTAIN EMINENTLY LASTLY ASSURED AGED INDEXED FOREGOING OFFERS UNDERTAKEN CIRCUMSTANCES ENTERTAINED  
 EXTREME ADMIRABLE RAPIDLY ULTIMATELY SAFER DECODE SUSPECTS STRUCK SERVICE DIMINISHES ESSENTIAL ALIENS FAIRLY OTHERS CONCERNED NEED

(d'')  
 (b'')  
 Fig. S3. Text mining in English. (Continued) (a'') Word patterns  $W_i$  in Charles Darwin's *Origin of Species* (including 15 chapters in the main text, but excluding the synopsis—the table of contents, the introduction, the glossary, and the index), sorted by descending  $n_{ii} \geq 20$ . (b'') Content word clusters  $W_i$  found in the synopsis of the aforementioned book, sorted by descending  $n_{ii} \geq 3$ . Font sizes are proportional to  $\sqrt{n_{ii}}$ . Temporal structures of the word patterns are ignored. (c'') False positives: automatically extracted topical patterns (colored green or red in panel a'') that are not mentioned anywhere in the synopsis. (d'') False negatives: automatically identified non-topic patterns (gray in panel a'') that are found in panel b''.  
 (c'')  
 (d'')

*Example 4.15.3.* We use exactly the same clustering algorithm to automatically generate words related to a particular question, in our WikiQA experiments. In Table S11 below, we list all the questions whose top-scoring answers (according to our algorithm) contain at least one sentence that has been officially labeled as correct by the WikiQA team. All the tabulated entries take the form of

[WikiQA-Q#]: [WikiQA question]	[Reference Wikipedia page]
[Sentence(s) bearing the highest score in our algorithm.]	

Here, while tabulating our results, we make no attempt to correct the spelling or capitalization in the question (see however, Algorithm 3.11), and the punctuation marks in the candidate answer(s) are kept in the same form as the original WikiQA dataset, such as

26: how did anne frank die	Anne Frank
----------------------------	------------

Anne Frank and her sister, Margot , were eventually transferred to the **Bergen-Belsen concentration camp** , where they **died of typhus** in March **1945**.

In the example above, the words related to the question (including the place of death “Bergen-Belsen concentration camp”, cause of death “typhus” and year of death “1945”, all discovered by our automated text mining) are highlighted in red.

Some (85 out of 580 total) question numbers in Table S11 are prefixed with  $\sharp$ , because the corresponding questions are classified as “quantitative”, by Algorithm 3.12. Intuitively speaking, one might answer such “quantitative” questions better by focusing on candidate sentences that mention numbers, dates and so on. However, in our current work, we have made no effort to treat “quantitative” questions with special screening procedures for the answers.

**Table S11. List of WikiQA hits**

0: HOW AFRICAN AMERICANS WERE IMMIGRATED TO THE US	African immigration to the United States
As such, <b>African immigrants</b> are to be distinguished from <b>African American</b> people, the latter of whom are descendants of mostly West and <b>Central Africans</b> who were involuntarily brought to the United States by means of the historic Atlantic slave trade .	
1: how are glacier caves formed?	Glacier cave
A <b>glacier cave</b> is a <b>cave formed</b> within the ice of a <b>glacier</b> .	
$\sharp$ 16: how much is 1 tablespoon of water	Tablespoon
In the USA one <b>tablespoon</b> (measurement unit) is approximately 15 mL; the capacity of an actual <b>tablespoon</b> (dining utensil) ranges from 7 mL to 14 mL.	
In countries where a <b>tablespoon</b> is a serving spoon, the nearest equivalent to the US <b>tablespoon</b> is either the dessert spoon or the soup spoon .	
A <b>tablespoonful</b> , nominally the capacity of one <b>tablespoon</b> , is commonly used as a measure of volume in cooking .	
18: how a rocket engine works	Rocket engine
A <b>rocket engine</b> , or simply “ <b>rocket</b> ”, is a <b>jet engine</b> that uses only stored propellant mass for <b>forming</b> its high speed propulsive <b>jet</b> .	
26: how did anne frank die	Anne Frank
Anne Frank and her sister, Margot , were eventually transferred to the <b>Bergen-Belsen concentration camp</b> , where they <b>died of typhus</b> in March <b>1945</b> .	
$\sharp$ 31: how old was monica lewinsky during the affair	Monica Lewinsky
<b>Monica Samille Lewinsky</b> (born July 23, 1973) is an American woman with whom United States President Bill Clinton admitted to having had an “improper relationship” while she worked at the White House in 1995 and 1996.	
33: how are antibodies used in	antibody
An <b>antibody</b> (Ab), also known as an immunoglobulin (Ig), is a large Y-shaped <b>protein produced</b> by B-cells that is <b>used</b> by the immune system to <b>identify</b> and neutralize foreign objects such as bacteria and viruses .	
$\sharp$ 45: how old is kirk douglas, the actor?	Kirk Douglas
<b>Kirk Douglas</b> (born Issur Danielovitch, ; December 9, 1916) is an American <b>stage</b> and film <b>actor</b> , film producer and author.	
$\sharp$ 50: how long was richard nixon a president	Richard Nixon
<b>Richard Milhouse Nixon</b> (January 9, 1913 – April 22, 1994) was the 37th <b>President</b> of the United States , <b>serving</b> from 1969 to 1974, when he became the only <b>president</b> to resign the office.	
$\sharp$ 64: How long was Mickie James with WWE?	Mickie James
<b>James</b> appeared in <b>World Wrestling Entertainment (WWE)</b> in October 2005 and was placed in a storyline with <b>Trish Stratus</b> , in which <b>James'</b> gimmick was that of <b>Stratus'</b> biggest fan turned obsessed stalker, an angle which lasted almost a year.	
77: how are the # of electrons in each shell determined	Electron shell
Each <b>shell</b> can <b>contain</b> only a fixed number of <b>electrons</b> : The 1st <b>shell</b> can hold up to two <b>electrons</b> , the 2nd <b>shell</b> can hold up to eight <b>electrons</b> , the 3rd <b>shell</b> can hold up to 18, and 4th <b>shell</b> can hold up to 32 and so on.	
$\sharp$ 94: how old is beatrice author	Bea Arthur
<b>Beatrice "Bea" Arthur</b> (May 13, 1922 – April 25, 2009) was an American actress, comedian, and singer whose career spanned several decades.	
98: how are public schools funded	state school
<b>State schools</b> (also known as <b>public schools</b> or <b>government schools</b> ) generally <b>refer</b> to primary or secondary <b>schools</b> mandated for or offered to all children without charge paid for, in whole or in part, by taxation .	
102: how does interlibrary loan work	Interlibrary loan
<b>Interlibrary loan</b> (abbreviated ILL, and sometimes called interloan, document delivery, or document supply) is a service whereby a user of one library can borrow books or receive photocopies of documents that are owned by another library. In many cases, nominal fees accompany <b>interlibrary loan</b> services.	
104: what did mia hamm do his work	Mia Hamm
Mariel Margaret "Mia" Hamm (born March 17, 1972) is a retired American professional soccer player.	
105: what bacteria grow on macconekey agar	MacConekey agar
<b>MacConkey agar</b> is a culture medium designed to <b>grow</b> Gram-negative <b>bacteria</b> and differentiate them for lactose fermentation .	
110: how do forensic auditors examine financial reporting	Financial audit
The purpose of an audit is provide and objective independent <b>examination</b> of the <b>financial statements</b> , which <b>increases</b> the value and credibility of the <b>financial statements</b> produced by management, thus <b>increase</b> user confidence in the <b>financial statement</b> , reduce investor <b>risk</b> and consequently reduce the cost of capital of the preparer of the <b>financial statements</b> .	
127: What committees are joint committees	Joint committee
<b>A Joint Committee</b> is a term in politics that is used to refer to a <b>committee</b> made up of members of both chambers of a bicameral legislature.	
$\sharp$ 130: How many states and territories are within India?	States and territories of India
<b>India</b> is a federal union of <b>states</b> comprising twenty-eight <b>states</b> and seven union <b>territories</b> .	
140: how is single malt scotch made	Single malt Scotch
<b>Single Malt Scotch</b> is <b>single malt</b> whisky <b>made</b> in Scotland using a pot still distillation process at a <b>single</b> distillery , with <b>malted</b> barley as the only grain ingredient.	
164: what county in texas is conroe located in	Conroe, Texas
<b>Conroe</b> is the seat of Montgomery <b>County</b> and falls within the metropolitan area.	
178: how does a dredge work?	Dredging
A <b>dredger</b> (or “ <b>dredge</b> ” as is the general usage in the Americas) is any device, <b>machine</b> , or vessel that is used to <b>excavate</b> and remove material from the bottom of a body of water.	
$\sharp$ 181: how many world series did curt schilling have	Curt Schilling
He helped lead the Philadelphia <b>Phillies</b> to the <b>World Series</b> in and <b>won</b> <b>World Series</b> championships in with the <b>Arizona Diamondbacks</b> and in and with the Boston Red Sox .	
188: what area code is 479	Area code 479

331: what day is st. patricks day	<i>Saint Patrick's Day</i>
Saint Patrick's Day or the Feast of Saint Patrick (, "the Day of the Festival of Patrick") is a cultural and religious holiday celebrated on 17 March.	
334: what cheese is made from goat's milk	<i>Goat cheese</i>
<b>Goat cheese</b> , or chèvre (from the French word for goat), is cheese made out of the milk of goats .	
338: what country is belize in	<i>Belize</i>
Belize , is a country located on the northeastern coast of Central America.	
342: how does a cat purr	<i>Purr</i>
However, using a strict definition of purring that continuous sound production must alternate between pulmonic egressive and ingressive airstream (and usually go on for minutes), Peters (2002), in an exhaustive review of the scientific literature, reached the conclusion that until then only 'purring cats' (Felidae) and two species of genets , Genetta tigrina, and most likely also Genetta genetta, had been documented to purr.	
346: what county is St. Elizabeth MO in	<i>St. Elizabeth, Missouri</i>
St. Elizabeth is a village in Miller County , Missouri , United States .	
348: what county is wilton ca in	<i>Wilton, California</i>
Wilton is a census-designated place (CDP) in Sacramento County , California , United States .	
#354: how many gold gloves does barry larkin have	<i>Barry Larkin</i>
Larkin is considered one of the top players of his era, winning nine Silver Slugger awards and three Gold Glove awards .	
361: how did women's role change during the war	<i>Women's roles in the World Wars</i>
Whether it was on the front front or the front-lines, for civilian or enlisted women, the World Wars started a new era for women's opportunities to contribute in war and be recognized for efforts outside of the home.	
366: how does lsd impact the human body	<i>Lysergic acid diethylamide</i>
Lysergic acid diethylamide, abbreviated LSD or LSD-25, also known as lysergide ( INN ) and colloquially as acid, is a semisynthetic psychedelic drug of the ergoline family, well known for its psychological effects which can include altered thinking processes, closed and open eye visuals, synesthesia , an altered sense of time and spiritual experiences , as well as for its key role in 1960s counterculture .	
Currently, a number of organizations—including the Beckley Foundation , MAPS , Heffter Research Institute and the Albert Hofmann Foundation—exist to fund, encourage and coordinate research into the medicinal and spiritual uses of LSD and related psychedelics.	
374: what county is orono maine in	<i>Orono, Maine</i>
Orono is a town in Penobscot County , Maine , United States .	
377: what it is a pilot study	<i>Pilot experiment</i>
A pilot experiment, also called a pilot study , is a small scale preliminary study conducted in order to evaluate feasibility, time, cost, adverse events, and effect size (statistical variability) in an attempt to predict an appropriate sample size and improve upon the study design prior to performance of a full-scale research project.	
#383: how many humps on a camel	<i>Camel</i>
The two surviving species of camel are the dromedary , or one-humped camel , which is native to the Middle East and the Horn of Africa , and the Bactrian , or two-humped camel , which inhabits Central Asia .	
384: what can be powered by wind	<i>Wind power</i>
Wind power is the conversion of wind energy into a useful form of energy, such as using wind turbines to make electrical power , windmills for mechanical power, wind pumps for water pumping or drainage , or sails to propel ships.	
389: what chili wants wiki	<i>What Chilli Wants</i>
What Chilli Wants is an American reality series on VH1 starring Chilli , one-third of the Grammy Award -winning R&B trio TLC .	
398: what division is boise state football	<i>Boise State Broncos football</i>
The Boise State Broncos football program represents Boise State University in college football and compete in the Football Bowl Subdivision (FBS) of Division I as a member of the Mountain West Conference .	
#409: how much caffeine is in a shot of espresso	<i>Espresso</i>
Espresso has more caffeine per unit volume than most beverages, but the usual serving size is smaller—a typical 60 mL (2 US fluid ounce ) of espresso has 80 to 150 mg of caffeine , little less than the 95 to 200 mg of a standard 240 mL (8 US fluid ounces ) cup of drip-brewed coffee.	
417: what does add my two cents mean	<i>My two cents</i>
"My two cents" (2¢) and its longer version "put my two cents in" is an United States (US) idiom expression, taken from the original English idiom expression: to put in "my two pennies worth" or "my tuppence worth."	
447: what does bruce jenner do	<i>Bruce Jenner</i>
William Bruce Jenner (born October 28, 1949) is a former U.S. track and field athlete , motivational speaker , socialite , television personality and businessman .	
448: what does alkali do to liquids?	<i>Alkali</i>
In chemistry , an alkali ( ; from Arabic : al-qaly , الْقَلْيٰ ) is a basic , ionic salt of an alkali metal or alkaline earth metal element .	
This broad use of the term is likely to have come about because alkalis were the first bases known to obey the Arrhenius definition of a base and are still among the more common bases.	
Some authors also define an alkali as a base that dissolves in water .	
450: what area code is 217	<i>Area code 217</i>
Area code 217 is the North American telephone area code for much of western and central Illinois .	
462: what does gloria in excelsis deo mean	<i>Gloria in Excelsis Deo</i>
"Gloria in excelsis Deo" ( Latin for "Glory to God in the highest") is a hymn known also as the Greater Doxology (as distinguished from the "Minor Doxology" or Gloria Patri ) and the Angelic Hymn.	
#463: how many grams in a troy ounce of gold	<i>Troy weight</i>
The troy ounce is 480 grains , compared with the avoirdupois ounce , which is 437½ grains.	
#473: how much does united states spend on health care	<i>Health care in the United States</i>
According to the World Health Organization (WHO), the United States spent more on health care per capita (\$7,146, and more on health care as percentage of its GDP (15.2%), than any other nation in 2008.	
#488: how many died in hiroshima and nagasaki	<i>Atomic bombings of Hiroshima and Nagasaki</i>
Within the first two to four months of the bombings, the acute effects killed 90,000–166,000 people in Hiroshima and 60,000–80,000 in Nagasaki , with roughly half of the deaths in each city occurring on the first day.	
490: what county is Holly Ridge nc in?	<i>Holly Ridge, North Carolina</i>
Holly Ridge is a town in Onslow County , North Carolina , United States .	
503: what kind of books does debbie macomber writes	<i>Debbie Macomber</i>
Debbie Macomber (born October 22, 1948 in Yakima, Washington ) is a best-selling American author of over 150 romance novels and contemporary women's fiction.	
504: what county is catonsville md in	<i>Catonsville, Maryland</i>
Catonsville is the home of the University of Maryland, Baltimore County (UMBC), a public research university with over 12,000 students.	
509: what kind of literature did john steinbeck writing	<i>John Steinbeck</i>
As the author of twenty-seven books , including sixteen novels, six non-fiction books , and five collections of short stories, Steinbeck received the Nobel Prize for Literature in 1962.	
510: whatever happened clint walker	<i>Clint Walker</i>
Norman Eugene Walker , known as Clint Walker (born May 30, 1927), is a retired American actor .	
512: what does the family leave act	<i>Family and Medical Leave Act of 1993</i>
The Family and Medical Leave Act of 1993 (FMLA) is a United States federal law requiring covered employers to provide employees job-protected and unpaid leave for qualified medical and family reasons .	
515: How is the pothole formed	<i>Pothole</i>
A pothole (sometimes called a kettle and known in parts of the Western United States as a chuckhole) is a type of disruption in the surface of a roadway where a portion of the road material has broken away, leaving a hole.	
517: what does estee lauder do	<i>Estée Lauder Companies</i>
Estée Lauder Companies, Inc. is a manufacturer and marketer of prestige skincare, makeup, fragrance and hair care products.	
534: what creates sonic boom	<i>Sonic boom</i>
A sonic boom is the sound associated with the shock waves created by an object traveling through the air faster than the speed of sound .	
551: what city is oregon state university in	<i>Oregon State University</i>
Oregon State University (OSU) is a coeducational , public research university located in Corvallis , Oregon , United States .	
#561: how many bones are in the skeletal system is composed of 306 bones	<i>Human skeleton</i>
Humans are born with over 270 bones , some of which fuse together into a longitudinal axis, the axial skeleton , to which the appendicular skeleton is attached.	
564: what county is north myrtle beach in SC	<i>North Myrtle Beach, South Carolina</i>
North Myrtle Beach is a coastal resort city in Horry County , South Carolina , United States .	
#568: how much of our universe does plasma make up	<i>Plasma (physics)</i>
In the universe , plasma is the most common state of matter for ordinary matter , most of which is in the rarefied intergalactic plasma (particularly intracluster medium ) and in stars.	
571: what does barefoot and pregnant mean	<i>Barefoot and pregnant</i>
The phrase "barefoot and pregnant" was probably first used sometime before 1950. The only way to keep a woman happy," he said, "is to keep her barefoot and pregnant."	
"Barefoot and pregnant" is a phrase most commonly associated with the controversial idea that women should not work outside the home and should have many children during their reproductive years.	
572: what county is oakhurst, nj in	<i>Oakhurst, New Jersey</i>
Oakhurst is a census-designated place and unincorporated community within Ocean Township , in Monmouth County , New Jersey , United States .	
574: what did lawrence joshua chamberlain do?	<i>Joshua Chamberlain</i>
Joshua Lawrence Chamberlain (September 8, 1828 – February 24, 1914), born as Lawrence Joshua Chamberlain , was an American college professor from the State of Maine , who volunteered during the American Civil War to join the Union Army .	
575: what circuit court is maryland	<i>Maryland Circuit Courts</i>
The Circuit Courts of Maryland are the state trial courts of general jurisdiction in Maryland.	
577: what causes a deficiency in adenosine deaminase	<i>Adenosine deaminase deficiency</i>
Adenosine deaminase deficiency , also called ADA deficiency or ADA-SCID, is an autosomal recessive metabolic disorder that causes immunodeficiency .	
582: how is human height measured	<i>Human height</i>
Human height is the distance from the bottom of the feet to the top of the head in a human body , standing erect.	
587: what does a timing belt do	<i>Timing belt</i>
Timing belt (camshaft), a toothed belt used to drive the camshaft(s) within an internal combustion engine A timing belt is a non-slipping mechanical drive belt and the term may refer to either:	
604: what channel is shopnbc on	<i>ShopNBC</i>
ShopNBC is an American broadcast and cable home shopping network, owned and operated by ValueVision Media , which is in turn 30% owned by GE Equity and NBC Universal .	
606: how post and lintels are used	<i>Post and lintel</i>
Post and lintel , "prop and lintel" or "trabeated" is a simple construction method using a lintel , header, or architrave as the horizontal member over a building void supported at its ends by two vertical columns , or .	
616: what county is cambrria wi in	<i>Cambrria, Wisconsin</i>
Cambrria is a village in Columbia County , Wisconsin , United States .	
#622: HOW MANY STRIPES ARE ON THE AMERICAN FLAG	<i>Flag of the United States</i>
The national flag of the United States of America , often simply referred to as the American flag , consists of thirteen equal horizontal stripes of red (top and bottom) alternating with white , with a blue rectangle in the canton (referred to specifically as the "union") bearing fifty small, white, five-pointed stars arranged in nine offset horizontal rows of six stars (top and bottom) alternating with rows of five stars.	
625: what do biologists do	<i>Biologist</i>
Biologists involved in basic research attempt to discover underlying mechanisms that govern how organisms work.	
Biologists involved in applied research attempt to develop or improve medical, industrial or agricultural processes.	
A biologist is a scientist who studies living organisms and their relationship to their environment.	
626: what does base jumping stand for	<i>BASE jumping</i>
BASE jumping , also sometimes written as B.A.S.E. jumping , is an activity where participants jump from fixed objects and use a parachute to break their fall.	
631: what country is the largest stalagmite	<i>Stalagmite</i>
The largest stalagmite in the world is high and is located in the cave of Cueva Martin Inferno, Cuba.	
633: what does Gringo mean	<i>Gringo</i>
Roger Axtell, a travel etiquette expert, notes that "[t]he word gringo is not necessarily a bad word.	
Gringo (, , ) is a slang Spanish and Portuguese word used in Ibero-America , to denote foreigners, often from the United States .	
638: what does hair testing show	<i>Drug test</i>
A drug test is a technical analysis of a biological specimen – for example urine, hair, blood, sweat, or oral fluid / saliva – to determine the presence or absence of specified parent drugs or their metabolites .	
644: what hormones produce thyroid	<i>Thyroid hormone</i>
The thyroid hormones , triiodothyronine (T3) and thyroxine (T4), are tyrosine -based hormones produced by the thyroid gland that are primarily responsible for regulation of metabolism.	
645: what does automatic paper feeder on printers mean	<i>Automatic Document Feeder</i>
In multifunction or all-in-one printers , fax machines , photocopiers and scanners , an automatic document feeder or ADF is a feature which takes several pages and feeds the paper one page at a time into a scanner or copier, allowing the user to scan , and thereby copy , print , or fax , multiple-page documents without having to manually replace each page.	
#660: how old was a child pedophile crime	<i>Pedophilia</i>
As a medical diagnosis, pedophilia or paedophilia is a psychiatric disorder in persons 16 years of age or older typically characterized by a primary or exclusive sexual interest toward prepubescent children (generally age 11 years or younger , though specific diagnosis criteria for the disorder extends the cut-off point for prepubescence to age 13).	
#662: how many people were killed in the holocaust	<i>The Holocaust</i>
Over one million Jewish children were killed in the Holocaust , as were approximately two million Jewish women and three million Jewish men.	
679: what channel is letterman on	<i>Late Show with David Letterman</i>
Late Show with David Letterman is an American late-night talk show hosted by David Letterman on CBS .	
683: what does karma mean in buddhism	<i>Karma in Buddhism</i>
Karma ( Sanskrit , also karman , Pāli : Kamma) means "action" or "doing"; whatever one does, says, or thinks is a karma.	
687: what kind a tilapia	<i>Tilapia</i>
Tilapia ( ) is the common name for nearly a hundred species of cichlid fish from the tilapiine cichlid tribe .	
#695: how many grape farms in united states	<i>Agriculture in the United States</i>
As of the last census of agriculture in 2007, there were 2.2 million farms , covering an area of , an average of per farm .	
696: what freezes faster? hot or cold water?	<i>Mpemba effect</i>
The Mpemba effect, named after Tanzanian student Erasto Mpemba , is the assertion that warmer water can freeze faster than colder water .	
698: what causes a small bowel obstruction	<i>Bowel obstruction</i>
Bowel obstruction (or intestinal obstruction) is a mechanical or functional obstruction of the intestines, preventing the normal transit of the products of digestion.	
#699: how many qfc stores are there	<i>QFC</i>
Quality Food Centers (QFC) is a supermarket chain based in Bellevue, Washington , with 64 stores in the Puget Sound region of Washington state and in the Portland, Oregon metropolitan area.	
705: What Causes Brain Freeze	<i>Ice-cream headache</i>
It is caused by having something cold touch the roof of the mouth ( palate ), and is believed to result from a nerve response causing rapid constriction and swelling of blood vessels or a "referring" of pain from the roof of the mouth to the head.	
708: what does leeroy jenkins mean	<i>Leeroy Jenkins</i>
Leeroy Jenkins , sometimes misspelled Leroy Jenkins and often elongated with numerous additional letters, is an Internet meme named for a player character created by Ben Schatz in Blizzard Entertainment 's MMORPG , World of Warcraft .	
#722: how many episodes of Lost were there	<i>List of Lost episodes</i>
A total of 121 episodes of Lost were produced, the last of which aired on May 23, 2010.	
ABC announced that Lost would end after six seasons, having produced a total of 121 episodes.	
#742: how much does a gold bar weigh	<i>Gold bar</i>
The standard gold bar held as gold reserves by central banks and traded among bullion dealers is the 400-troy-ounce (12.4 kg or 438.9 ounces) Good Delivery gold bar .	
743: what county is bolingbrook il in?	<i>Bolingbrook, Illinois</i>
Bolingbrook is a large village in Will and DuPage Counties in the U.S. state of Illinois .	
#745: how many kids does archie manning have	<i>Archie Manning</i>
He is the father of current Denver Broncos quarterback Peyton Manning , current New York Giants starting quarterback Eli Manning , and former Ole Miss receiver Cooper Manning .	

#746: what does an advocacy website promote?	<i>Advocacy</i>
Research is beginning to explore how <b>advocacy</b> groups in the U.S. and Canada are using social media to facilitate civic engagement and collective action.	
Lobbying (often by lobby groups) is a form of <b>advocacy</b> where a direct approach is made to legislators on an issue which plays a significant role in modern politics.	
<b>Advocacy</b> is a political process by an individual or group which aims to influence public-policy and resource allocation decisions within political, economic, and social systems and institutions.	
<b>Advocacy</b> can include many activities that a person or organization undertakes including media campaigns, public speaking, commissioning and publishing research or polls or the filing of an <i>amicus brief</i> .	
#750: how many amendments in us	<i>United States Constitution</i>
The Constitution has been <b>amended</b> seventeen <i>additional</i> times (for a total of <b>twenty-seven</b> amendments).	
751: what city is george washington university	<i>George Washington University</i>
The <b>George Washington University</b> (GW, GWU, or <b>George Washington</b> ) is a comprehensive private , coeducational research <b>university</b> located in <b>Washington, D.C.</b>	
#755: How many Muslims live in the United Kingdom?	<i>Islam in the United Kingdom</i>
The vast majority of <b>Muslims</b> in the <b>United Kingdom</b> live in England and Wales : of 1,591,000 <b>Muslims</b> recorded at the 2001 Census, 1,536,015 were <b>living</b> in England and Wales , where they formed 3% of the <b>population</b> in 2001; 42,557 were <b>living</b> in Scotland , forming 0.84% of the <b>population</b> ; and 1,943 were <b>living</b> in Northern Ireland .	
#763: how long to take two jima	<i>Battle of Iwo Jima</i>
The Battle of <b>Iwo Jima</b> (19 February – 26 March 1945), or <b>Operation Detachment</b> , was a major battle in which the <b>United States</b> Armed Forces fought for and captured the island of <b>Iwo Jima</b> from the Japanese Empire .	
764: What a Margarita contains	<i>Margarita</i>
The <b>margarita</b> is a Mexican cocktail consisting of tequila mixed with Cointreau or similar orange -flavoured liqueur and lime or lemon juice , often served with salt on the glass rim.	
The drink is served shaken with ice (on the rocks), blended with ice (frozen margarita), or without ice (straight up).	
766: what color is burgundy	<i>Burgundy (color)</i>
<b>Burgundy</b> is a dark red <b>color</b> associated with the <b>Burgundy</b> wine of the same name, which in turn is named after the <b>Burgundy</b> region of France .	
773: what day is the feast of st joseph's?	<i>Saint Joseph's Day</i>
Saint <b>Joseph's Day</b> , March 19, the <b>Feast of St. Joseph</b> is in Western Christianity the principal <b>feast day</b> of Saint Joseph , Spouse of the Blessed Virgin Mary .	
791: what does a laboratory in a gynecologist office consist of	<i>Medical laboratory</i>
A <b>medical laboratory</b> or <b>clinical laboratory</b> is a <b>laboratory</b> where tests are done on clinical specimens in order to get information about the health of a patient as pertaining to the diagnosis, treatment, and prevention of disease.	
812: what countries allow gays to openly serve in the military	<i>Sexual orientation and military service</i>
Nations that permit <b>gay</b> people to <b>serve openly</b> in the <b>military</b> include the 4 of the 5 members of the UN Security Council (United States, United Kingdom, France, and Russia), the Republic of China (Taiwan), Australia , <b>Israel</b> , <b>South Africa</b> , Argentina , and all NATO members excluding Turkey .	
813: How Works Diaphragm Pump	<i>Diaphragm pump</i>
A <b>diaphragm pump</b> (also known as a Membrane <b>pump</b> , Air Operated Double <b>Diaphragm Pump</b> (AODD) or Pneumatic <b>Diaphragm Pump</b> ) is a positive displacement <b>pump</b> that uses a combination of the reciprocating action of a rubber , thermoplastic or teflon <b>diaphragm</b> and suitable valves either side of the <b>diaphragm</b> ( check valve , butterfly valves, flap valves , or any other form of shut-off valves) to <b>pump</b> a fluid .	
#815: how many percent is a basis point	<i>Basis point</i>
The relationship between percentage changes and <b>basis points</b> can be summarized as follows: 1 percentage <b>point</b> change = 100 <b>basis points</b> , and 0.01 percentage <b>points</b> = 1 <b>basis point</b> .	
819: what genre is bloody beetroots	<i>The Bloody Beetroots</i>
The <b>Bloody Beetroots</b> is well-known for the black Venom mask he wears during performances.	
The <b>Bloody Beetroots</b> is the pseudonym of Sir Bob Cornelius Rifo, the Italian electro house and dance-punk music producer , DJ and photographer .	
"The <b>Bloody Beetroots</b> DJ set" contains Sir Bob Cornelius Rifo and Tommy Teia.	
#838: how many asian Indians live in usa	<i>Indian American</i>
<b>Indian Americans</b> are citizens of the United States of <b>Indian</b> ancestry and comprise about 3.18 million <b>people</b> , or ~1.0% of the U.S. population , the country's third largest self-reported <b>Asian</b> ancestral group after Chinese <b>Americans</b> and Filipino <b>Americans</b> according to <b>American</b> Community Survey of 2010 data.	
843: what caused the world war 2	<i>Causes of World War II</i>
The main <b>causes</b> of <b>World War II</b> were nationalistic issues, unresolved issues, and resentments resulting from <b>World War I</b> and the interval period in Europe , in addition to the effects of the Great Depression in the 1930s.	
867: how is root beer made?	<i>Root beer</i>
<b>Root beer</b> is a carbonated , sweetened beverage , originally <b>made</b> using the <b>root</b> of a sassafras plant (or the bark of a sassafras tree) as the primary flavor.	
872: what does informal logic mean	<i>Informal logic</i>
<b>Informal logic</b> , intuitively, refers to the principles of <b>logic</b> and <b>logical</b> thought outside of a formal setting.	
874: what country is turkey in	<i>Turkey</i>
Turkey ( ), officially the Republic of <b>Turkey</b> , is a transcontinental <b>country</b> , <b>located</b> mostly on Anatolia in Western Asia and on East Thrace in Southeastern Europe .	
876: what county is Augusta,GA located in?	<i>Augusta, Georgia</i>
<b>Augusta</b> is the principal city of the <b>Augusta – Richmond County</b> Metropolitan Statistical Area , which as of 2010 had an estimated population of 556,877, making it both the second-largest city and the second-largest metro area in the <b>state</b> after <b>Atlanta</b> .	
893: what does a cutter do	<i>Cutter (baseball)</i>
In baseball , a <b>cutter</b> , or <b>cut</b> fastball, is a type of fastball which breaks slightly toward the pitcher 's glove side as it reaches home plate .	
When a batter is able to hit a <b>cutter</b> pitch, it often results in soft contact and an easy out, due to the pitch's movement keeping the ball away from the bat's sweet spot .	
The <b>cutter</b> is typically 2–5 mph slower than a pitcher's four-seam fastball .	
Some pitchers use a <b>cutter</b> as a way to prevent hitters from expecting their regular fastballs .	
In 2010, the average pitch classified as a <b>cutter</b> by PITCHf/x thrown by a right-handed pitcher was 88.6 mph; the average four-seamer was 92.1 mph .	
An animated diagram of a <b>cutter</b>	
A common technique used to throw a <b>cutter</b> is to use a four-seam fastball grip with the baseball set slightly off center in the hand .	
895: what it takes aerosmith album	<i>What It Takes (song)</i>
"What It <b>Takes</b> " is a power ballad by American hard rock band <b>Aerosmith</b> .	
898: what county is erie colorado	<i>Erie, Colorado</i>
<b>Erie</b> is a Statutory Town in Boulder and Weld <b>counties</b> in the U.S. <b>state</b> of <b>Colorado</b> .	
902: what classes are considered humanities	<i>Humanities</i>
The <b>humanities</b> that are also regarded as social <b>sciences</b> include history , anthropology , area <b>studies</b> , communication studies , cultural studies , law , economics and linguistics .	
910: what food is in afghan	<i>Afghan cuisine</i>
Accompanying these staples are dairy products ( yogurt and whey ), various <b>nuts</b> , and native vegetables, as well as fresh and dried <b>fruit</b> ; <b>Afghanistan</b> is well known for its grapes .	
922: what kind of cut is tri tip	<i>Tri-tip</i>
The <b>tri-tip</b> is a <b>cut</b> of beef from the bottom sirloin primal <b>cut</b> .	
923: what county is willmar mn in?	<i>Willmar, Minnesota</i>
<b>Willmar</b> is a city in, and the <b>county</b> seat of, Kandiyohi <b>County</b> , Minnesota , United States .	
#943: how many innings makes an official game	<i>Official game</i>
Since most professional baseball <b>games</b> are nine <b>innings</b> long, the fifth <b>inning</b> is used as the threshold for an <b>official game</b>	
963: how tennessee became a state	<i>Tennessee</i>
<b>Tennessee</b> was the <b>last state</b> to leave the <b>Union</b> and join the <b>Confederacy</b> at the outbreak of the U.S. Civil War in 1861, and the <b>first state</b> to be readmitted to the <b>Union</b> at the end of the war.	
967: what county is coatesville indiana located in	<i>Coatesville, Indiana</i>
<b>Coatesville</b> is a town in Clay Township, Hendricks <b>County</b> , Indiana , United States .	
968: how does Delaware support its claim to being the first state?	<i>Delaware</i>
<b>Delaware</b> was one of the 13 <b>colonies</b> participating in the <b>American</b> Revolution and on December 7, 1787, became the <b>first state</b> to ratify the Constitution of the United States , thereby becoming known as The <b>First State</b> .	
973: what does fidelity do	<i>Fidelity Investments</i>

1120: who sings i am a man of constant sorrow <b>"Man of Constant Sorrow"</b> (also known as "I Am A Man of Constant Sorrow") is a traditional American folk song first recorded by Dick Burnett , a partially blind fiddler from Kentucky .	<i>Man of Constant Sorrow</i>
1123: who kill franz ferdinand ww1 On 28 June 1914, Archduke <b>Franz Ferdinand</b> of Austria , heir presumptive to the Austro-Hungarian throne, and his wife, Sophie, Duchess of Hohenberg , were shot <b>dead</b> in Sarajevo , by Gavrilo Princip , one of a group of six Bosnian Serb assassins coordinated by Danilo Ilic .	<i>Assassination of Archduke Franz Ferdinand of Austria</i>
1140: what is a nms message Multimedia <b>Messaging Service (MMS)</b> is a standard way to send <b>messages</b> that include multimedia content to and from mobile phones .	<i>Multimedia Messaging Service</i>
1150: who sang proud mary <b>"Proud Mary"</b> is a rock song written by American singer-songwriter and multi-instrumentalist John Fogerty , and <b>recorded</b> by his band Creedence Clearwater Revival .	<i>Proud Mary</i>
1154: where do crocodiles live <b>Crocodiles</b> (subfamily <i>Crocodylinae</i> ) or <b>true crocodiles</b> are large aquatic tetrapods that <b>live</b> throughout the tropics in Africa , Asia , the Americas and Australia .	<i>Crocodiles</i>
1157: what relates to erosion Water and wind <b>erosion</b> are now the two primary causes of land degradation ; combined, they are responsible for 84% of degraded acreage, making excessive <b>erosion</b> one of the most significant global <b>environmental</b> problems.	<i>Erosion</i>
1159: where in oregon is albany <b>Albany</b> is the 11th largest city in the U.S. state of <b>Oregon</b> , and is the <b>county</b> seat of <b>Linn County</b> .	<i>Albany, Oregon</i>
#1168: when did abraham lincoln write the emancipation proclamation?? <b>The Emancipation Proclamation</b> was an <b>order</b> issued to all segments of the <b>Executive</b> branch (including the Army and Navy) of the United States by President <b>Abraham Lincoln</b> on January 1, 1863, during the American Civil War .	<i>Emancipation Proclamation</i>
1172: what part of the pre-world war 1 arms race was the most intense? The <b>causes</b> of World War I , which <b>began</b> in central Europe in late <b>July</b> 1914 and finished in 1918, included many factors, such as the conflicts and hostility of the four decades leading up to the <b>war</b> .	<i>Causes of World War I</i>
#1182: when president nixon resigns Richard Milhouse <b>Nixon</b> (January 9, 1913 – April 22, 1994) was the 37th <b>President</b> of the United States , <b>serving</b> from 1969 to 1974 , when he became the <b>only president</b> to <b>resign</b> the office.	<i>Richard Nixon</i>
1184: where fourth of july came from Independence Day, commonly known as the <b>Fourth of July</b> , is a federal <b>holiday</b> in the <b>United States</b> commemorating the adoption of the Declaration of Independence on <b>July 4</b> , 1776, declaring independence from the Kingdom of Great Britain .	<i>Independence Day (United States)</i>
1194: what people used mayan numeral system <b>Maya numerals</b> are a vigesimal ( base - twenty ) <b>numerical system</b> used by the Pre-Columbian Maya civilization .	<i>Maya numerals</i>
1196: what is a newsgroup message <b>A Usenet newsgroup</b> is a usually within the <b>Usenet</b> system, for <b>messages</b> posted from many users in different locations.	<i>Usenet newsgroup</i>
1204: who accompanied King louis the VII of France on the second crusade The <b>Second Crusade</b> was announced by <b>Pope Eugene III</b> , and was the <b>first</b> of the <b>crusades</b> to be led by <b>European</b> kings, namely <b>Louis VII of France</b> and <b>Conrad III of Germany</b> , with help from a number of other <b>European</b> nobles.	<i>Second Crusade</i>
#1207: when did gary moore die Robert William <b>Gary Moore</b> (4 April 1952 – 6 February 2011), was a Northern Irish musician, most widely recognised as a singer and guitarist.	<i>Gary Moore</i>
1209: what is a hosting company on a website A <b>web hosting service</b> is a type of <b>Internet hosting service</b> that allows individuals and organizations to make their <b>website</b> accessible via the World Wide Web .	<i>Web hosting service</i>
1218: who owns exxon mobil It is a direct descendant of John D. Rockefeller 's Standard Oil company, and was formed on November 30, 1999, by the <b>merger</b> of <b>Exxon</b> and <b>Mobil</b> .	<i>ExxonMobil</i>
1222: what is a rock quarry <b>A quarry</b> is a type of open-pit mine from which <b>rock</b> or minerals are extracted.	<i>Quarry</i>
1223: what kind of school is MIT The <b>MIT Sloan School of Management</b> (also known as <b>MIT</b> Sloan or <b>Sloan</b> ) is the business <b>school</b> of the <b>Massachusetts Institute of Technology</b> , in Cambridge , <b>Massachusetts</b> , USA .	<i>MIT Sloan School of Management</i>
1225: what are banana plugs for A <b>banana connector</b> (commonly <b>banana plug</b> for the male , <b>banana socket</b> or <b>banana jack</b> for the female ) is a single-wire (one conductor ) electrical connector used for joining wires to equipment .	<i>Banana connector</i>
#1229: when did the trojan war take place The ancient <b>Greeks</b> thought that the <b>Trojan War</b> was a <b>historical event</b> that had <b>taken place</b> in the 13th or 12th century BC , and believed that Troy was <b>located</b> in modern-day Turkey near the Dardanelles .	<i>Trojan War</i>
1230: what is a letterbox movie The term refers to the shape of a <b>letter</b> box , a slot in a wall or door through which mail is delivered, being rectangular and wider than it is high.	<i>Letterboxing (filming)</i>
1236: what are the busiest airports in the world The definition of the <b>world's busiest airport</b> has been specified by the <b>Airports</b> Council International in Geneva, Switzerland .	<i>World's busiest airport</i>
1248: where the streets have no name filming location The song was notably performed on a Los Angeles <b>rooftop</b> for the <b>filming</b> of its music <b>video</b> , which won a Grammy Award for Best Performance Music <b>Video</b> . Recently the song has been used by the NFL's Baltimore Ravens as their entrance song in Super Bowl XLVII	<i>Where the Streets Have No Name</i>
#1243: when did the civil war start and where The <b>American Civil War</b> (ACW), also known as the <b>War between the States</b> or simply the <b>Civil War</b> (see naming ), was a <b>civil war</b> fought from 1861 to 1865 between the United States (the "Union" or the "North") and several <b>Southern slave states</b> that declared their secession and formed the Confederate States of America (the "Confederacy" or the "South").	<i>American Civil War</i>
1244: What are the busiest airports in the world The definition of the <b>world's busiest airport</b> has been specified by the <b>Airports</b> Council International in Geneva, Switzerland .	<i>World's busiest airport</i>
1248: where the streets have no name filming location The song was notably performed on a Los Angeles <b>rooftop</b> for the <b>filming</b> of its music <b>video</b> , which won a Grammy Award for Best Performance Music <b>Video</b> . Recently the song has been used by the NFL's Baltimore Ravens as their entrance song in Super Bowl XLVII	<i>Where the Streets Have No Name</i>
1253: where did the mayflower land <b>The Mayflower</b> was the ship that in 1620 transported 102 English Pilgrims , including a core group of Separatists , to New England .	<i>Mayflower</i>
1257: what languages are spoken in south africa The English version of the <b>South African</b> constitution refers to the <b>languages</b> by the names in those <b>languages</b> : isiZulu , isiXhosa , Afrikaans , Sepedi (referring to Northern Sotho) , Setswana , English , Sesotho (referring to Southern Sotho) , Xitsonga , Siswati , Tshivenda and isiNdebele (referring to <b>Southern</b> Ndebele).	<i>Languages of South Africa</i>
1262: Who controlled Alaska before US? The name " <b>Alaska</b> " (Аляска) was already introduced in the <b>Russian</b> colonial period , when it was used only for the peninsula and is derived from the Aleut alaxsxaq, meaning "the mainland" or, more literally, "the object towards which the action of the sea is directed".	<i>Alaska</i>
1268: who sang black velvet <b>"Black Velvet"</b> is a blues verse with a rock chorus written by Canadian songwriters Christopher Ward and David Tyson , recorded by Canadian <b>singer</b> songwriter Alannah Myles .	<i>Black Velvet (song)</i>
1274: who produced loyal to the game? Released in the United States December 14, 2004 (December 12 in the United Kingdom ) , <b>Loyal to the Game</b> was <b>produced</b> by Eminem .	<i>Loyal to the Game</i>
1275: WHAT ARE HERITABLE TRAITS <b>Heritability</b> of a <b>trait</b> within a population is the proportion of observable differences in a <b>trait</b> between individuals within a population that is due to genetic differences.	<i>Heritability</i>
1284: who does afge represent The American Federation of Government Employees (AFGE) is an American labor union <b>representing</b> over 650,000 employees of the federal government , about 5,000 employees of the District of Columbia , and a few hundred private sector employees, mostly in and around federal facilities.	<i>American Federation of Government Employees</i>
1292: who wrote puff the magic dragon <b>"Puff, the Magic Dragon"</b> is a song <b>written</b> by Leonard Lipton and Peter Yarrow , and made popular by Yarrow's group Peter, Paul and Mary in a 1963 recording	<i>Puff, the Magic Dragon</i>
1301: what is a vetting process <b>Vetting</b> is the <b>process</b> of performing a background check on someone before offering them employment, conferring an award, etc.	<i>Vetting</i>
#1303: when did the free soilers party start? The <b>Free Soil Party</b> was a short-lived political <b>party</b> in the United States active in the 1848 and 1852 presidential elections, and in some state elections.	<i>Free Soil Party</i>
1308: where does the expression "knocking on wood" come from <b>Knocking on wood</b> , or touch <b>wood</b> , refers to the apotropaic tradition in western folklore of literally touching <b>knocking on wood</b> , or merely stating that you are doing or intend same, in order to avoid " tempting fate " after making a favourable observation, a boast, or declaration concerning one's own death.	<i>Knocking on wood</i>
1310: who sang the nights lights went out "The Night the Lights Went Out in Georgia" is a Southern Gothic song written by songwriter Bobby Russell and performed in 1972 by his then-wife Vicki Lawrence .	<i>The Night the Lights Went Out in Georgia</i>
1315: what is a synthetic conduit A <b>nerve guidance conduit</b> (also referred to as an <b>artificial nerve conduit</b> or <b>artificial nerve graft</b> , as opposed to an <b>autograft</b> ) is an <b>artificial</b> means of guiding axonal regrowth to facilitate <b>nerve</b> regeneration and is one of several clinical <b>treatments</b> for <b>nerve</b> injuries .	<i>Nerve guidance conduit</i>
1316: what is a wwii theater The <b>European Theatre</b> of World War II, also known as the <b>European War</b> , was a huge area of heavy fighting across <b>Europe</b> from Germany's invasion of Poland on <b>September</b> 1, 1939 until the end of the <b>war</b> with the German unconditional surrender on May 8, 1945 (V-E Day ).	<i>European Theatre of World War II</i>
1318: what is amoxicillin for? <b>Amoxicillin</b> is susceptible to degradation by $\beta$ -lactamase -producing bacteria, which are resistant to a broad spectrum of $\beta$ -lactam antibiotics, such as penicillin .	<i>Amoxicillin</i>
1321: what state was the civil war in The <b>American Civil War</b> (ACW), also known as the <b>War between the States</b> or simply the <b>Civil War</b> (see naming ), was a <b>civil war</b> fought from 1861 to 1865 between the United States (the "Union" or the "North") and several <b>Southern slave states</b> that declared their <b>secession</b> and formed the <b>Confederate States of America</b> (the "Confederacy" or the "South").	<i>American Civil War</i>
1322: where did the early humans live? According to the Recent African Ancestry theory , modern <b>humans</b> evolved in Africa possibly from <b>Homo heidelbergensis</b> , <b>Homo rhodesiensis</b> or <b>Homo</b> antecessor and migrated out of the continent some 50,000 to 100,000 years ago, replacing local populations of <b>Homo erectus</b> , <b>Homo denisova</b> , <b>Homo floresiensis</b> and <b>Homo neanderthalensis</b> .	<i>Human evolution</i>
1325: who invaded north africa during ww2 Operation Torch (initially called Operation Gymnast) was the British – American <b>invasion</b> of French <b>North Africa</b> in World War II during the <b>North African Campaign</b> , started on 8 November 1942.	<i>Operation Torch</i>
#1326: WHAT YEARS WAS THE C5 VETTE PRODUCED The <b>Chevrolet Corvette (C5)</b> is a sports car <b>produced</b> by the <b>Chevrolet</b> division of General Motors for the 1997 through 2004 model years .	<i>Chevrolet Corvette (C5)</i>
1328: what U.S. President's head has been featured on the nickel (five-cent coin) since 1938? <b>Nickel (United States coin)</b> The <b>Buffalo nickel</b> was <b>introduced</b> in 1913 as part of a drive to increase the beauty of American <b>coinage</b> , in 1938, the <b>Jefferson nickel</b> followed.	<i>Nickel (United States coin)</i>
1329: who sings stand by me "Stand by Me" is a song originally performed by Ben E. King and written by King, Jerry Leiber , and Mike Stoller , inspired by the spiritual "Lord Stand by Me," plus two lines rooted in Psalms 46:2–3.	<i>Stand by Me (song)</i>
#1331: When did F15s first fly The <b>Eagle first flew</b> in July 1972 , and entered service in 1976.	<i>McDonnell Douglas F-15 Eagle</i>
1336: what are the arb medications Angiotensin receptor blocker (antagonist), a <b>medication</b> for treating high blood pressure	<i>ARB</i>
1341: what are SATA Power Connector Serial ATA (SATA) is a computer bus interface that <b>connects</b> host bus adapters to mass storage <b>devices</b> such as hard disk <b>drives</b> and optical <b>drives</b> .	<i>Serial ATA</i>
1344: What kind of poem is "This Is Just To Say" "This Is Just To Say" (1934) is a famous imagist <b>poem</b> by William Carlos Williams .	<i>This Is Just To Say</i>
1347: who played dumbledore in harry potter Dumbledore is <b>portrayed</b> by Richard <b>Harris</b> in the film adaptions of <b>Harry Potter</b> and the Philosopher's Stone and Harry Potter and the Chamber of Secrets .	<i>Albus Dumbledore</i>
1348: what produces calcitonin <b>Calcitonin</b> (also known as thyrocalcitonin) is a 32-amino <b>acid</b> linear polypeptide hormone that is <b>produced</b> in humans primarily by the parafollicular cells (also known as C-cells) of the thyroid , and in many other animals in the ultimobranchial body .	<i>Calcitonin</i>
1351: who hit the first home run at riverfront stadium? On June 30, 1970, the Reds hosted the Atlanta Braves in their grand opening, with Hank Aaron <b>hitting</b> the <b>first ever home run</b> at Riverfront.	<i>Riverfront Stadium</i>
1352: Who Makes Neken Tires <b>Neken Tire</b> is a <b>tire</b> manufacturer, headquartered in Yangsan , South Gyeongsang Province , and Seoul , both in South Korea .	<i>Neken Tire</i>
1353: who did richard nixon refer to as the silent majority The term was popularized (though not first used) by U.S. President <b>Richard Nixon</b> in a November 3, 1969, speech in which he said, "And so tonight – to you, the great <b>silent majority</b> of my fellow Americans—I ask for your support."	<i>Silent majority</i>
1355: what are the lateral and median apertures of the brain? It is an opening in each <b>lateral</b> extremity of the <b>lateral</b> recess of the fourth ventricle of the human <b>brain</b> , which also has a single <b>median aperture</b> .	<i>Lateral aperture</i>
1358: what part of the pig do pork chops come from A <b>pork chop</b> is a <b>chop</b> of <b>pork</b> (a meat <b>chop</b> ) cut perpendicularly to the spine of the <b>pig</b> and usually containing a rib or part of a vertebra, served as an individual portion.	<i>Pork chop</i>
1362: what are grits made from Modern <b>grits</b> are commonly <b>made</b> of alkali-treated corn known as hominy .	<i>Grits</i>
1367: what state is new england in New England is a region in the northeastern corner of the United States consisting of the six states of Maine , New Hampshire , Vermont , Massachusetts , Rhode Island , and Connecticut .	<i>New England</i>
1378: what are stanzas in poetry A <b>stanza</b> consists of a grouping of two or more lines , set off by a space , that usually has a set pattern of meter and rhyme. The <b>stanza in poetry</b> is analogous with the paragraph that is seen in prose , related thoughts are grouped into units.	<i>Stanza</i>
1382: who shot franz ferdinand Assassination of Archduke Franz Ferdinand of Austria	<i>Assassination of Archduke Franz Ferdinand of Austria</i>
On 28 June 1914, Archduke <b>Franz Ferdinand</b> of Austria , heir presumptive to the Austro-Hungarian throne, and his wife, Sophie, Duchess of Hohenberg , were <b>shot</b> dead in Sarajevo , by Gavrilo Princip , one of a group of six Bosnian Serb assassins coordinated by Danilo Ilic .	
1386: what are the three ossicles The <b>ossicles</b> (also called auditory <b>ossicles</b> ) are the <b>three</b> smallest bones in the human body, the malleus , the incus and the stapes .	<i>Ossicles</i>
1391: what is an army specialist <b>Specialist</b> (abbreviated "SPC") is one of the four junior enlisted ranks in the U.S. <b>Army</b> , just above Private First Class and equivalent in pay grade to <b>Corporal</b> .	<i>Specialist (rank)</i>
1399: what to make with linen Many products are <b>made</b> of <b>linen</b> : aprons, bags, towels (swimmers, bath, beach, <b>body</b> and wash towels), napkins, <b>bed linens</b> , <b>linen</b> tablecloths, runners, chair covers, and men's & women's wear.	<i>Linen</i>
1404: what south dakota county is wakonda in <b>Wakonda</b> is a town in <b>Clay County , South Dakota</b> , United States .	<i>Wakonda, South Dakota</i>
1408: what are american people of japanese descent called <b>Japanese Americans</b> have historically been among the three largest Asian <b>American</b> communities, but in recent decades, it has become the sixth largest group at roughly 1,304,286, including those of mixed-race or mixed-ethnicity.	<i>Japanese American</i>
1420: what is a bus adapter in a computer? In <b>computer</b> hardware , a host controller, host <b>adapter</b> , or host <b>bus adapter</b> (HBA) connects a host system (the <b>computer</b> ) to other network and storage devices.	<i>Host adapter</i>
1421: where does cashmere come from	<i>Cashmere wool</i>

1423: Who invented egg rolls	Egg roll
2_Egg roll	Varieties of egg rolls are found in mainland China , many Chinese-speaking regions of Asia, and Chinese immigrant communities around the world.
1424: what westerners include in pain and suffering	Pain and suffering
Pain and suffering	is the legal term for the physical and emotional stress caused from an injury (see also pain and suffering ).
1429: who won the 2010 world cup	2010 FIFA World Cup
In the final , Spain , the European champions , defeated third-time finalists the Netherlands 1-0 after extra time , with Andrés Iniesta 's goal in the 116th minute giving Spain their first world title , becoming the eighth nation to win the tournament, and the first European nation to win the tournament outside its home continent.	
1431: who wrote the song in the mood	In the Mood
"In the Mood" is a big band era #1 hit recorded by American bandleader Glenn Miller .	
1434: what is active learning strategies	Active learning
Active learning	is an umbrella term that refers to several models of instruction that focus the responsibility of learning on learners .
1442: who sang I want to dance with somebody	I Wanna Dance With Somebody (Who Loves Me)
"I Wanna Dance with Somebody (Who Loves Me)" is the first single from Whitney Houston 's second studio album Whitney .	
1456: who sang the banana boat song	Day-O (The Banana Boat Song)
"Day-O (The Banana Boat Song)" is a traditional Jamaican mento folk song , the best-known version of which was sung by Harry Belafonte and an alternate version interspersed with another Jamaican folksong , Hill and Gully Rider , by Dame Shirley Bassey .	
#1459: what year did keeping up with the kardashians begin	Keeping Up with the Kardashians
Keeping Up with the Kardashians (often referred to simply as The Kardashians) is an American reality television series that premiered on October 14, 2007, on E ! .	
1464: what is a store confectioner	Confectionery store
A confectionery store (more commonly referred to as a sweet shop in the United Kingdom, a candy store in the North America, or a lolly shop in Australia) sells confectionery and is usually targeted to children.	
#1465: when did qing dynasty begin	Qing Dynasty
The Qing Dynasty , also Empire of the Great Qing or Great Qing , was the last imperial dynasty of China , ruling from 1644 to 1912 with a brief, abortive restoration in 1917.	
1489: who wrote serenity prayer	Serenity Prayer
The Serenity Prayer is the common name for an originally untitled prayer by the American theologian Reinhold Niebuhr (1892–1971).	
1490: what are the Declaration and Resolves of the First Continental Congress aboutDeclaration and Resolves of the First Continental Congress	
The Declaration and Resolves of the First Continental Congress (also known as the Declaration of Colonial Rights, or the Declaration of Rights), was a statement adopted by the First Continental Congress on October 14, 1774, in response to the Intolerable Acts passed by the British Parliament .	
1499: who created massey ferguson	Massey Ferguson
The company was formed by a merger between Massey Harris and the Ferguson Company farm machinery manufacturer in 1953, creating the company Massey Harris Ferguson .	
1506: where do the mohawks live	Mohawk people
Their traditional homeland stretched southward of the Mohawk River , eastward to the Green Mountains of Vermont , westward to the border with the Oneida Nation 's traditional homeland territory , and northward to the St Lawrence River .	
#1511: when did proof die	Proof (rapper)
In 2006, Proof was shot and killed during an altercation at the CCC nightclub in Detroit.	
DeShaun Dupree Holton (October 2, 1973 – April 11, 2006), better known by his stage name Proof , was an American rapper and actor from Detroit , Michigan .	
1524: what are social security taxes	Social Security (United States)
Tax deposits are formally entrusted to the Federal Old-Age and Survivors Insurance Trust Fund , the Federal Disability Insurance Trust Fund , the Federal Hospital Insurance Trust Fund , or the Federal Supplementary Medical Insurance Trust Fund which comprise the Social Security Trust Fund .	
1527: what town is laurel hollow my in	Laurel Hollow, New York
Laurel Hollow is a village in the Town of Oyster Bay in Nassau County, New York in the United States .	
1530: who won fifa world cup 2010	2010 FIFA World Cup
In the final , Spain , the European champions , defeated third-time finalists the Netherlands 1-0 after extra time , with Andrés Iniesta 's goal in the 116th minute giving Spain their first world title , becoming the eighth nation to win the tournament, and the first European nation to win the tournament outside its home continent.	
#1531: when did egg McMuffin get invented	McMuffin
The Egg McMuffin is the signature breakfast sandwich; it was invented by the late McDonald's franchisee Herb Peterson in the late 1960s and was introduced nationwide in 1972.	
1544: what are the names of destiny's child	Destiny's Child
Formed in 1990 in Houston , Texas, Destiny's Child members began their musical endeavors as Girl's Tyme comprising, among others, Knowles, Rowland, LaTavia Roberson and LeToya Luckett .	
Destiny's Child was launched into mainstream recognition following the release of their best-selling second album, The Writing's on the Wall , which contained the number-one singles " Bills , Bills , Bills " and " Say My Name ".	
Destiny's Child was an American R&B girl group whose final, and perhaps most recognizable, line-up comprised Beyoncé Knowles , Kelly Rowland and Michelle Williams .	
Destiny's Child has sold more than 50 million records worldwide to date.	
After years of limited success, they were signed to Columbia Records as Destiny's Child .	
1551: What are procedure codes in coding?	Procedure codes
Procedure codes	are numbers or alphanumeric codes used to identify specific health interventions taken by medical professionals.
1562: what is a duvet cover used for	Duvet
A duvet ( or ; from the French duvet "down"), also known as a doona in Australian English or a continental quilt (or simply quilt) in British English , but this usage is no longer common, is a type of bedding , a soft flat bag filled with down , feathers , wool , silk or a synthetic alternative , and protected with a removable cover , analogous to a pillow and pillow case.	
1563: who did mr bojangles	Mr. Bojangles (song)
"Mr. Bojangles" is the title of a song originally written and recorded by American country music artist Jeff Walker for his 1968 album of the same title.	
1574: what is an agents job role in film	Talent agent
A talent agent , or booking agent , is a person who finds jobs for actors , authors , film directors , musicians , models , producers , professional athletes , writers and other people in various entertainment businesses .	
1576: WHAT IS A LAW ENFORCEMENT MURDER BOOK	Murder book
In law enforcement parlance, the term murder book refers to the case file of a murder investigation.	
#1581: when album love always released	Love Always
Love Always is the debut album of American R&B duo K-Ci & JoJo , released on June 17, 1997, by MCA Records .	
1583: who wrote west side story	West Side Story
West Side Story is an American musical with a book by Arthur Laurents , music by Leonard Bernstein , lyrics by Stephen Sondheim , and conception and choreography by Jerome Robbins .	
1600: what state does interstate 70 travel through	Interstate 70
Interstate 70 (I-70) is an Interstate Highway in the United States that runs from Interstate 15 near Cove Fort , Utah , to a Park and Ride near Baltimore , Maryland .	
1604: where do mangos come from	Mango
The mango is native to South Asia , from where it has been distributed worldwide to become one of the most cultivated fruits in the tropics .	
1607: what video format will play in a DVD player	DVD-Video
Discs using the DVD-Video specification require a DVD drive and an H.262/MPEG-2 Part 2 decoder (e.g., a DVD player , or a computer DVD drive with a software DVD player).	
1610: what part of the government governs the US post office?	United States Postal Service
The United States Postal Service (USPS), also known as the Post Office and U.S. Mail, is an independent agency of the United States federal government responsible for providing postal service in the United States .	
1614: what motor does a 2001 monte carlo	Chevrolet Monte Carlo
The Monte Carlo SS was revived from 2000 to 2007 and initially powered by a 3.8-liter V6 (supercharged in 2004–2005), later to be replaced by a 5.3-liter V8 for 2006–2007 .	
#1619: what year was President kennedy president?	John F. Kennedy
John Fitzgerald "Jack" Kennedy (May 29, 1917 – November 22, 1963), often referred to by his initials JFK, was the 35th President of the United States , serving from 1961 until his death in 1963 .	
1620: who owns popeyes chicken	Popeyes Louisiana Kitchen
Often referred to as Popeyes and sometimes as Popeyes Chicken & Biscuits or Popeyes Chicken & Seafood , It was acquired by Sandy Springs, Georgia -based AFC Enterprises , originally America's Favorite Chicken Company, in 1993.	
1624: who said give me liberty or give me death	Give me Liberty, or give me Death!
"Give me liberty, or give me death!" is a quotation attributed to Patrick Henry from a speech he made to the Virginia Convention .	
1630: what are the most known sports in america	Sports in the United States
Four of the nation's five most popular team sports were developed in North America: American football , baseball , basketball and ice hockey , whereas soccer was developed in England.	
1639: what is aden disease	Acute disseminated encephalomyelitis
Acute disseminated encephalomyelitis (ADEM) is an immune mediated disease of the brain .	
#1642: what year did john adams become president?	John Adams
John Adams (October 30, 1735 (O.S. October 19, 1735) – July 4, 1826) was the second president of the United States (1797–1801), having earlier served as the first vice president of the United States .	
1648: who created bingo	Bing
Bing (known previously as Live Search , Windows Live Search , and MSN Search ) is a web search engine (advertised as a "decision engine ") from Microsoft .	
1653: what are anti inflammatories	Anti-inflammatory
Anti-inflammatory refers to the property of a substance or treatment that reduces inflammation .	
Anti-inflammatory drugs make up about half of analgesics , remedying pain by reducing inflammation as opposed to opioids , which affect the central nervous system .	
1656: what sang momma told me not to come	Mama Told Me Not to Come
"Mama Told Me (Not to Come)" is a song by Randy Newman written for Eric Burdon 's first solo album in 1966.	
1660: what are dialysis tubes	Dialysis tubing
Dialysis tubing is a type of semi- or partially permeable membrane tubing made from regenerated cellulose or cellophane .	
1662: who designed the statue of liberty	Statue of Liberty
The Statue of Liberty (Liberty Enlightening the World; French: La Liberté éclairant le monde) is a neoclassical sculpture on Liberty Island in New York Harbor , designed by Frédéric Bartholdi and dedicated on October 28, 1886.	
1667: what is am 2201	AM-2201
The toxicity of AM-2201 is still a matter of debate and there may be long term side effects .	
Repeated recreational use of AM-2201 in the United States has led to it being specifically listed in a proposed 2011 amendment to the Controlled Substances Act , aiming to add a number of synthetic drugs into Schedule 1 .	
AM-2201 (1-(5-fluoropropyl)-3-(1-naphthoyl)indole) is a research chemical that acts as a potent but nonselective full agonist for the cannabinoid receptor .	
#1677: when did george washington died?	George Washington
George Washington (–1799) was the first President of the United States (1789–1797), the commander-in-chief of the Continental Army during the American Revolutionary War , and one of the Founding Fathers of the United States .	
1683: what are four thirds cameras	Four Thirds system
The Four Thirds system is a standard created by Olympus and Kodak for digital single-lens reflex camera (DSLR) design and development . [ref name="DPReview.com 2001-02-13"]	
1691: who wrote nature boy	Nature Boy
"Nature Boy" is a song by Eden Ahbez , published in 1947.	
Nat King Cole 's 1948 recording of the song was a major hit and "Nature Boy" has since become a pop and jazz standard , with dozens of major artists interpreting the song .	
#1698: when did playstation 3 first launch	PlayStation 3 launch
The launch of the PlayStation 3 into the Japanese market on 10 November 2006 marked the second major seventh generation entertainment system to be released.	
1708: who wrote the song cocaine	Cocaine (song)
"Cocaine" is a song written and recorded by JJ Cale in 1976, but also known as a cover version recorded by Eric Clapton .	
#1710: when did pearl harbor get bombed	Attack on Pearl Harbor
The attack on Pearl Harbor (called Hawaii Operation or Operation AI by the Japanese Imperial General Headquarters (Operation Z in planning) and the Battle of Pearl Harbor) was a surprise military strike conducted by the Imperial Japanese Navy against the United States naval base at Pearl Harbor , Hawaii, on the morning of December 7, 1941 (December 8 in Japan).	
#1714: when did marlins start	Miami Marlins
The Miami Marlins began play in the 1993 season as the Florida Marlins .	
Per agreement with the city and Miami-Dade County (which owns the park), the Marlins officially changed their name to the "Miami Marlins" on November 11, 2011.	
1716: where do dungensis crab live	Dungeness crab
The Dungeness crab , Metacarcinus magister (formerly Cancer magister), is a species of crab that inhabits eelgrass beds and water bottoms on the west coast of North America .	
1720: who plays mary poppins in the movie	Mary Poppins (film)
Mary Poppins is a 1964 musical film starring Julie Andrews , Dick Van Dyke , David Tomlinson , and Glynis Johns , produced by Walt Disney , and based on the Mary Poppins books series by P. L. Travers .	
1722: What are the different translations for the Bible	Bible translations
The full Bible has been translated into 518 languages , and 2798 languages have at least some portion of the Bible .	
#1723: when does black friday start for christmas	Black Friday (shopping)
Black Friday is the day following Thanksgiving Day in the United States , often regarded as the beginning of the Christmas shopping season .	
#1730: when did Emperor Justinian die	Justinian I
Justinian I ( ) (c. 482 – 14 November 565), commonly known as Justinian the Great , was Byzantine Emperor from 527 to 565 .	
1734: what movement was elizabeth cady stanton a part of	Elizabeth Cady Stanton
Elizabeth Cady Stanton (November 12, 1815 – October 26, 1902) was an American social activist , abolitionist, and leading figure of the early women's rights movement .	
1738: what radio station are the boston bruins on?	Boston Bruins Radio Network
The Boston Bruins Radio Network is a 25-station (17 A.M. , 8 F.M. plus 2 F.M. translators) CBS Radio network which carries live game broadcasts of the Boston Bruins .	
#1739: when did ms. drgs go into effect	Diagnosis-related group
The system is also referred to as the DRGs , and its intent was to identify the "products" that a hospital provides .	
DRGs may be further grouped into Major Diagnostic Categories (MDCs) .	
DRGs have been used in the US since 1982 to determine how much Medicare pays the hospital for each "product", since patients within each category are clinically similar and are expected to use the same level of hospital resources .	
DRGs are assigned by a "grouper" program based on ICD (International Classification of Diseases) diagnoses, procedures, age, sex, discharge status, and the presence of complications or comorbidities .	
DRGs are also standard practice for establishing reimbursements for other Medicare related reimbursements such as to home healthcare providers .	
1741: what are the characteristics of an ethnic group	Ethnic group
Ethnic identity is constantly reinforced through common characteristics which set the group apart from other groups .	
1746: what mountains are on appalachian trail	Appalachian Trail
The Appalachian National Scenic Trail , generally known as the Appalachian Trail or simply the A.T., is a marked hiking trail in the eastern United States extending between Springer Mountain in Georgia and Mount Katahdin in Maine .	
1748: what are add ons	Add-on
Add-ons for Firefox (addons.mozilla.org), the official Mozilla Foundation web site which serves as a repository of add-ons for Mozilla-based applications .	
Add-on (Mozilla) , a piece of software that enhances and customizes Mozilla-based applications .	
Add-on might mean:	
1749: what would be the deliverables	Deliverable
A deliverable also differs from a project document in that project document is typically part of a project deliverable , or a project deliverable may contain number of documents and physical things .	
It may be either an outcome to be achieved (as in "The corporation says that becoming profitable this year is a deliverable .") or an output to be provided (as in "The deliverable for the completed project consists of a special-purpose electronic device and its controlling software .").	
Deliverable is a term used in project management to describe a tangible or intangible object produced as a result of the project that is intended to be delivered to a customer (either internal or external) .	

A <b>deliverable</b> may be composed of multiple smaller <b>deliverables</b> .	
A <b>deliverable</b> differs from a project milestone in that a milestone is a measurement of progress toward an output whereas the <b>deliverable</b> is the result of the process.	
In technical projects, <b>deliverables</b> can further be classified as hardware , software , or design documents .	
In contracted efforts, <b>deliverable</b> may refer to an item specifically required by contract documents, such as an item on a Contract Data Requirements List or mentioned in the Statement Of Work .	
For a typical project, a milestone might be the completion of a product design while the <b>deliverable</b> might be the technical diagram of the product.	
A <b>deliverable</b> could be a report, a document, a server upgrade or any other building block of an overall project.	
1754: what school did oprah winfrey help	Oprah Winfrey Leadership Academy for Girls
The <b>Oprah Winfrey Leadership Academy for Girls</b> - South Africa is a female boarding <b>school</b> founded in January 2007 and located in Henley on Klip near Meyerton , south of Johannesburg , South Africa .	
# 1757: when did scott pilgrim come out	Scott Pilgrim vs. the World
Scott Pilgrim vs. the World is a 2010 American <b>comedy</b> film directed by Edgar Wright , based on the graphic novel series Scott Pilgrim by Bryan Lee O'Malley .	
1760: what is a google in math	Googol
In the binary numeral system , one would need 333 bits to represent a <b>googol</b> , i.e., $1 \text{ googol} = 2^{332.19}$ , or exactly .	
Other names for <b>googol</b> include ten duotrigintillion on the short scale , ten thousand sexdecillion on the long scale , or ten sexdecilliard on the Peletier long scale .	
It is also the namesake of the internet company Google , with the name "Google" being a misspelling of " <b>googol</b> " by the company's founders.	
A <b>googol</b> is the large number 10100; that is, the digit 1 followed by 100 zeroes :	
A <b>googol</b> is approximately 70!	
A <b>googol</b> has no particular significance in mathematics, but is useful when comparing with other very large quantities such as the number of subatomic particles in the visible universe or the number of hypothetically possible chess games.	
1 <b>Googol</b> (1E100) is a small fraction more, than the biggest number a usual hand-held calculator can display and handle, which is 9.999...	
1763: what structure is disulfide bonds	Disulfide bond
In chemistry , a <b>disulfide bond</b> (Br.E. disulphide <b>bond</b> ) is a covalent <b>bond</b> , usually derived by the coupling of two thiol groups .	
1768: who rules communism government	Communist state
It has a form of <b>government</b> characterized by single-party <b>rule</b> or dominant-party <b>rule</b> of a communist party (referred as Dictatorship of the Proletariat by its proponents) and a professed allegiance to a Leninist or Marxist-Leninist ideology as the guiding principle of the state.	
1772: who makes skittles?	Skittles (confectionery)
<b>Skittles</b> is a brand of fruit-flavoured sweets , currently produced and marketed by the Wm. Wrigley Jr. Company , a division of Mars, Inc.	
1778: what is a discipline of study	List of academic disciplines
An academic <b>discipline</b> , or field of <b>study</b> , is a branch of knowledge that is taught and researched at the college or university level.	
1783: who wrote stand by me	Stand by Me (song)
"Stand By Me" is a song originally performed by Ben E. King and <b>written</b> by King, Jerry Leiber , and Mike Stoller, inspired by the spiritual "Lord Stand By Me," plus two lines rooted in Psalms 46:2-3.	
1784: what are metaphors used for	Metaphor
This quote is a <b>metaphor</b> because the world is not literally a stage.	
One of the most prominent examples of a <b>metaphor</b> in English literature is the All the world's a stage monologue from As You Like It :	
<b>Metaphor</b> is a type of analogy and is closely related to other rhetorical figures of speech that achieve their effects via association, comparison or resemblance including allegory , hyperbole , and simile .	
In simpler terms, a <b>metaphor</b> compares two objects or things without using the words "like" or "as".	
A <b>metaphor</b> is a figure of speech that describes a subject by asserting that it is, on some point of comparison, the same as another otherwise unrelated object.	
1788: what are the names of the ll divos	Il Divo
Il Divo is an English multinational operatic pop vocal group created by music manager, executive, and reality TV star Simon Cowell .	
Il Divo is a group of four male singers: French pop singer Sébastien Izambard , Spanish baritone Carlos Marín , American tenor David Miller , and Swiss tenor Urs Bühler .	
1792: What is a 28 day cycle?	Menstrual cycle
This article focuses on the human menstrual <b>cycle</b> , a "monthly" <b>cycle</b> that can vary around an average of <b>-28 days</b> per <b>cycle</b> .	
1797: what is a llc company?	Limited liability company
An <b>LLC</b> is a legal form of <b>company</b> that <b>provides</b> limited liability to its owners in the vast majority of <b>United States jurisdictions</b> .	
1798: what are k cups	K-Cup
K-Cup portion packs are used with Keurig or other single <b>cup brewing systems</b> to brew a <b>cup</b> of coffee , tea , or hot chocolate .	
1800: what is a chronograph watch	Chronograph
A <b>chronograph</b> is a specific type of <b>watch</b> that is used as a stopwatch combined with a display <b>watch</b> .	
1801: what is a PCI port used for	Conventional PCI
<b>Conventional PCI</b> (PCI) is an initialism formed from Peripheral Component Interconnect, part of the <b>PCI</b> Local Bus <b>standard</b> and often shortened to just <b>PCI</b> is a local <b>computer</b> bus for attaching hardware devices in a <b>computer</b> .	
# 1813: when did the movie deep blue sea come out	Deep Blue Sea
Deep Blue Sea is a 1999 science fiction horror film that stars Saffron Burrows , Thomas Jane , LL Cool J , Michael Rapaport , Stellan Skarsgård and Samuel L. Jackson .	
1818: what are a and r reps	Artists and repertoire
<b>Artists</b> and <b>repertoire</b> (A&R) is the division of a record label or music publishing company that is responsible for talent scouting and overseeing the <b>artistic</b> development of recording <b>artists</b> and/or songwriters.	
1819: what measurement is a furlong	Furlong
A <b>furlong</b> is a <b>measure</b> of distance in imperial units and U.S. customary units equal to one-eighth of a mile , equivalent to 220 yards , 660 feet , 40 rods , or 10 chains .	
# 1825: when the wind blows james patterson	When the Wind Blows (James Patterson novel)
When the Wind Blows is a novel by <b>James Patterson</b> .	
1829: what are the catholic gifts of the holy spirit	Seven gifts of the Holy Spirit
The <b>seven gifts</b> of the <b>Holy Spirit</b> is an enumeration of <b>seven spiritual gifts</b> originating with patristic authors, later elaborated by five intellectual <b>virtues</b> and four other groups of ethical characteristics.	
1832: what are lobbying groups	Lobbying
<b>Lobbying</b> is done by many different types of people and organized <b>groups</b> , including individuals in the private sector , corporations , fellow legislators or government officials, or advocacy <b>groups</b> (interest groups).	
1840: what nationality is kris jenner	Kris Jenner
Kristen Mary "Kris" Jenner (née Houghton, previously Kardashian; born November 5, 1955) is an American socialite, author and television personality.	
# 1850: when did lucy stone died	Lucy Stone
Lucy Stone (August 13, 1818 – October 19, 1893) was a prominent American abolitionist and suffragist , and a vocal advocate and organizer promoting rights for women .	
1851: what are stocks and bonds?	Bond (finance)
<b>Bonds</b> and <b>stocks</b> are both securities , but the major <b>difference</b> between the two is that (capital) stockholders have an equity stake in the <b>company</b> (i.e. they are owners), whereas <b>bondholders</b> have a creditor stake in the <b>company</b> (i.e. they are lenders).	
1852: what are the boundaries of the pelvic outlet	Pelvic outlet
The lower circumference of the lesser pelvis is very irregular; the space enclosed by it is named the inferior aperture or <b>pelvic outlet</b> .	
# 1858: when did germans enter parts in ww2	Liberation of Paris
The Liberation of <b>Paris</b> (also known as the Battle for <b>Paris</b> ) took place during World War II from 19 <b>August</b> 1944 until the surrender of the occupying <b>German</b> garrison on <b>25 August</b> .	
1859: who wrote the song feelin alright	Feelin' Alright
" <b>Feelin'</b> Alright?" also known as " <b>Feeling</b> Alright" is a <b>song</b> <b>written</b> by Dave Mason of the English rock band Traffic from their eponymous 1968 album, Traffic .	
1860: what is a neuro tract	Neural pathway
A neural <b>pathway</b> , neural <b>tract</b> , or neural face, connects one part of the nervous system with another and usually consists of bundles of elongated, myelin-insulated neurons , known collectively as white matter .	
# 1863: when did dr carter g woodson die	Carter G. Woodson
<b>Carter Godwin Woodson</b> (December 19, 1875April 3, 1950) was an African-American historian , author , journalist and the founder of the Association for the Study of African American Life and History .	
1864: what are loan origination	Loan origination
<b>Loan origination</b> is the <b>process</b> by which a borrower <b>applies</b> for a new <b>loan</b> , and a lender <b>processes</b> that <b>application</b> .	
# 1872: what year did aerosmith i dont want to miss a thing	I Don't Want to Miss a Thing
"I Don't Want to Miss a Thing" is a song performed by American rock <b>band</b> Aerosmith for the 1998 film Armageddon .	
1886: who makes triumph motorcycles	Triumph Motorcycles
<b>Triumph Motorcycles</b> Ltd , a current British <b>motorcycle</b> manufacturer	
1894: what produces primary xylem?	Xylem
The word <b>xylem</b> is derived from the Greek word ξύλον (xylon), meaning " <b>wood</b> "; the best-known <b>xylem</b> tissue is <b>wood</b> , though it is <b>found</b> throughout the plant.	
1897: what is a ti 82?	TI-82
The <b>TI-82</b> was designed in 1993 as a stripped down, more user friendly version of the TI-85 , and as a replacement for the TI-81 .	
The <b>TI-82</b> is a graphing calculator made by Texas Instruments .	
Like the TI-81 , the <b>TI-82</b> features a 96x64 pixel display, and the core feature set of the TI-81 with many new features.	
1899: what makes of the united states	United States
The <b>United States</b> of America (USA or U.S.A.), commonly called the <b>United States</b> (US or U.S.) or America, is a <b>federal republic</b> consisting of fifty <b>states</b> and a <b>federal district</b> .	
# 1919: when did thomas jefferson become president	Thomas Jefferson
<b>Thomas Jefferson</b> (April 13, 1743 (April 2, 1743 O.S.) – July 4, 1826) was an <b>American</b> Founding <b>Father</b> , the principal author of the Declaration of Independence (1776) and the third <b>President</b> of the United States (1801–1809).	
1920: who sang sun city	Sun City (song)
" <b>Sun City</b> " is a 1985 protest song written by Steven Van Zandt , produced by Van Zandt and Arthur Baker and recorded by Artists United Against <b>Apartheid</b> to convey opposition to the South African policy of <b>apartheid</b> .	
1931: who wrote the song hallelujah	Hallelujah (Leonard Cohen song)
" <b>Hallelujah</b> " is a song written by Canadian singer-songwriter Leonard Cohen , originally released on his album Various Positions (1984).	
1932: what are garnishments	Garnishment
A <b>garnishment</b> is a means of collecting a <b>monetary</b> judgment against a defendant by ordering a third party (the <b>garnishee</b> ) to pay <b>money</b> ; otherwise owed to the defendant , directly to the plaintiff.	
1944: what are land parcels	Parcel
A <b>land lot</b> , a piece of <b>land</b> ;	
1945: what state is Mn	Minnesota
Minnesota () is a <b>U.S. state</b> located in the Midwestern United States .	
1947: what naturally occurring isotopes does cobalt have	Isotopes of cobalt
<b>Naturally occurring cobalt</b> (Co) is composed of 1 stable <b>isotope</b> , 59Co. 28 radioisotopes have been characterized with the most stable being 60Co with a half-life of 5.2714 years, 57Co with a half-life of 271.79 days, 56Co with a half-life of 77.27 days, and 58Co with a half-life of 70.86 days.	
1948: Who Moved My Cheese Synopsis	Who Moved My Cheese?
It describes <b>change</b> in one's work and life, and four typical reactions to said <b>change</b> by two mice and two "littlepeople", during their hunt for <b>cheese</b> .	
1949: who replaced nikita khrushchev	Nikita Khrushchev
<b>Khrushchev's party colleagues</b> removed him from power in 1964, replacing him with <b>Leonid Brezhnev</b> as First Secretary and Alexei Kosygin as Premier.	
1953: What Are Mnemonic Devices	Mnemonic
A <b>mnemonic</b> (, with a silent "m"), or <b>mnemonic device</b> , is any <b>learning</b> technique that aids <b>information</b> retention.	
1962: What is a surveyor's wheel	Surveyor's wheel
A <b>surveyor's wheel</b> , also called a clickwheel, hodometer, waywiser, trundle <b>wheel</b> , measuring <b>wheel</b> , or perambulator is a device for measuring distance.	
# 1967: what year was girls just want to have fun release	Girls Just Want to Have Fun
" <b>Girls Just Want to Have Fun</b> " is a 1979 song originally written by Robert Hazard and made famous by singer Cyndi Lauper .	
1972: what are the players in quidditch?	Quidditch
Matches are <b>played</b> between two teams of seven <b>players</b> riding flying broomsticks, using four balls and six elevated ring-shaped goals, three on each side of the <b>Quidditch</b> pitch (field).	
1975: who wrote the song a little more country than that>	A Little More Country Than That
" <b>A Little More Country Than That</b> " is the title of a <b>song</b> written by Joey + Rory 's Rory Lee Feek , Wynne Varble , and Don Poythress, and recorded by American <b>country</b> artist Easton Corbin .	
1976: who makes blackberry	BlackBerry
The <b>BlackBerry</b> is a line of wireless handheld devices and services designed and marketed by Research In Motion <b>Limited</b> (RIM) operating as <b>BlackBerry</b> .	
1979: what are arizona's symbols	List of Arizona state symbols
The newest adopted <b>symbol</b> of <b>Arizona</b> is the Colt Single Action Army in 2011.	
The following is a list of <b>symbols</b> of the U.S. state of <b>Arizona</b> .	
The first <b>symbol</b> was the motto, which was made official in 1864 for the <b>Arizona</b> Territory .	
Fourteen of the state <b>symbols</b> are on display on the <b>Arizona</b> Capitol Museum .	
1983: what is adoration catholic church	Eucharistic adoration
Eucharistic <b>adoration</b> is a practice in the Roman <b>Catholic Church</b> , and in a few Anglican and Lutheran <b>churches</b> , in which the Blessed Sacrament is exposed and <b>adored</b> by the faithful.	
# 1984: when the body is systemic	Systemic
<b>Systemic</b> refers to something that is spread throughout, <b>system-wide</b> , affecting a group or <b>system</b> such as a <b>body</b> , economy, market or society as a whole.	
1985: what are slr cameras	Single-lens reflex camera
A single-lens reflex (SLR) <b>camera</b> is a <b>camera</b> that typically uses a mirror and prism system (hence "reflex", from the mirror's reflection) that permits the photographer to view through the lens and see exactly what will be captured, contrary to viewfinder <b>camera</b> where the image could be significantly different from what will be captured.	
1986: what is a vm server	Virtual machine
A virtual machine (VM) is a software implemented <b>abstraction</b> of the underlying hardware , which is presented to the application layer of the <b>system</b> .	
# 1998: when did the titanic sink	RMS Titanic
RMS <b>Titanic</b> was a British passenger liner that <b>sank</b> in the North Atlantic Ocean on 15 April 1912 after colliding with an <b>iceberg</b> during her <b>maiden voyage</b> from Southampton , UK to New York City , US.	
1999: where is the human thigh located?	Human leg
The <b>human leg</b> is the entire <b>lower extremity</b> or <b>limb</b> of the <b>human</b> body , including the foot , <b>thigh</b> and even the hip or gluteal region; however, the precise definition in <b>human</b> anatomy refers only to the section of the <b>lower limb</b> extending from the <b>knee</b> to the <b>ankle</b> .	
2001: what is considered a large car	Full-size car
A full-size <b>car</b> is a marketing term used in North America for an automobile <b>larger</b> than a mid-size <b>car</b> .	
2004: who is flo from progressive	Flo (Progressive Insurance)
In 2011, <b>Progressive</b> introduced an Australian counterpart to <b>Flo</b> , named Kitty, played by Australian actress Holly Austin. <b>Flo</b> is a fictional character who appears in commercials for <b>Progressive</b> Insurance .	
2009: what is petit le mans	Petit Le Mans
The <b>Petit Le Mans</b> (French for little Le Mans) is a sports car endurance race held annually at Road Atlanta in Braselton, Georgia , USA.	
The <b>Petit Le Mans</b> covers a maximum of (which is approximately 394 laps) or a maximum of 10 hours, whichever comes first; only once, in the rain-stopped 2009 race, has the leading team failed to complete .	
2015: what is considered to be a disasters	Disaster
Developing countries suffer the greatest costs when a <b>disaster</b> hits – more than 95 percent of all deaths caused by <b>disasters</b> occur in developing countries, and losses due to natural <b>disasters</b> are 20 times greater (as a percentage of GDP ) in developing countries than in industrialized countries.	
Ruins from the 1906 San Francisco earthquake , remembered as one of the worst natural <b>disasters</b> in United States history .	
In contemporary academia, <b>disasters</b> are seen as the consequence of inappropriately managed risk .	
Hazards that strike in areas with low vulnerability will never become <b>disasters</b> , as is the case in uninhabited regions.	

2017: what is firebird server	<i>Firebird (database server)</i>	
The database forked from Borland 's open source edition of InterBase in 2000, but since <b>Firebird</b> 1.5 the code has been largely rewritten.		
<b>Firebird</b> is an open source SQL relational database management system that runs on Linux , Windows , and a variety of Unix .		
2019: What is caused by the human immunodeficiency virus?	<i>HIV</i>	
<b>Human immunodeficiency virus</b> (HIV) is a lentivirus (slowly replicating retrovirus ) that <b>causes</b> acquired immunodeficiency syndrome (AIDS), a condition in <b>humans</b> in which progressive failure of the immune system <b>allows</b> life-threatening opportunistic infections and cancers to thrive.		
2026: what is the function of the vas defens? <i>Vas defens</i>		
<b>vas defens</b> (plural: <i>vasa defensita</i> ), also called <b>ductus defens</b> ( Latin : "carrying-away vessel"; plural: <i>ductus defens</i> ), is part of the male anatomy of many vertebrates ; they transport sperm from the epididymis in anticipation of ejaculation .		
2027: What is Ischemia or infarction?	<i>Ischemia</i>	
<b>Ischemia</b> is the toes with characteristic cyanosis .		
<b>Ischemia</b> is generally caused by problems with blood vessels , with resultant damage to or dysfunction of tissue . In medicine , <b>ischemia</b> , also spelled as <i>ischæmia</i> or <i>ischæma</i> ; (; from Greek language <i>ἰσχαιμή</i> , <i>ischaimia</i> ; <i>isch-</i> root denoting a restriction or thinning or to make or grow thin/lean, <i>blood</i> a) is a restriction in blood supply to tissues , causing a shortage of oxygen and glucose needed for cellular metabolism (to keep tissue alive).		
2035: where is humboldt ks	<i>Humboldt, Kansas</i>	
<b>Humboldt</b> is a city situated along the Neosho River in the southwest part of Allen County , located in southeast Kansas , in the Central United States .		
2036: what was nixon accused of	<i>Watergate scandal</i>	
The Watergate <b>scandal</b> was a political <b>scandal</b> that occurred in the United States in the 1970s as a result of the June 17, 1972 break-in at the Democratic National Committee headquarters at the Watergate office complex in Washington, D.C., and the <b>Nixon administration</b> 's attempted cover-up of its involvement.		
2043: who is victoria jackson from saturday night live	<i>Victoria Jackson</i>	
<b>Victoria Jackson</b> (born August 2, 1959) is an American comedian, actress, satirist, singer and internet blogger best known as a cast member of the NBC television sketch comedy series <b>Saturday Night Live</b> (SNL) from 1986 to 1992.		
2044: what was the post modernist era in literature?	<i>Postmodern literature</i>	
Postmodern <b>literature</b> is <b>literature</b> characterized by heavy reliance on techniques like <b>fragmentation</b> , paradox, and questionable narrators, and is often (though not exclusively) defined as a style or <b>trend</b> which emerged in the post-World War II era.		
2045: who is the group enigma	<i>Enigma (musical project)</i>	
The Romanian-born Cretu conceived the <b>Enigma</b> project while working in Germany, but based his recording studio A.R.T. Studios in Ibiza, Spain , from the early 1990s until May 2009, where he has recorded all of <b>Enigma</b> 's studio releases to date.		
2049: what is the federal death tax?	<i>Estate tax in the United States</i>	
The estate <b>tax</b> in the United States is a <b>tax</b> imposed on the <b>transfer</b> of the " <b>taxable estate</b> " of a deceased <b>person</b> , whether such <b>property</b> is <b>transferred</b> via a will , according to the <b>state</b> laws of intestacy or otherwise made as an incident of the <b>death</b> of the owner, such as a <b>transfer</b> of <b>property</b> from an intestate estate or trust, or the payment of <b>certain</b> life insurance benefits or financial account sums to beneficiaries.		
#2052: when was raphael born	<i>Raphael</i>	
Raffaello Sanzio da <b>Urbino</b> (April 6 or March 28, 1483 – April 6, 1520), better known simply as <b>Raphael</b> , was an Italian painter and architect of the High Renaissance .		
2059: where is the seed located in an artichoke	<i>Artichoke</i>	
The edible matter is <b>buds</b> that form within the <b>flower</b> heads before the <b>flowers</b> come into bloom.		
2065: WHAT IS PARESTHESIAS OF HANDS	<i>Paresthesia</i>	
The manifestation of <b>paresthesia</b> may be transient or chronic .		
<b>Paresthesia</b> ( or ), is a sensation of tickling, tingling, burning, pricking, or numbness of a person's skin with no apparent long-term physical effect.		
2072: who is paul avery to the zodiac killings	<i>Paul Avery</i>	
<b>Paul Avery</b> (April 2, 1934December 10, 2000) was an American police reporter, best known for his stories on the infamous serial killer known as the <b>Zodiac</b> , and later for his work on the Patricia Hearst kidnapping.		
2074: what is sump pump used for	<i>Sump pump</i>	
A <b>sump pump</b> is a <b>pump</b> used to remove water that has accumulated in a water collecting <b>sump</b> basin, commonly found in the basement of homes.		
#2075: when was bloody kansas	<i>Bleeding Kansas</i>	
Bleeding <b>Kansas</b> , <b>Bloody Kansas</b> or the Border War, was a series of violent political confrontations involving anti-slavery Free-Staters and pro-slavery "Border Ruffian" elements, that took place in the <b>Kansas</b> Territory and the neighboring towns of Missouri between 1854 and 1861.		
2076: What Is Benzene Used For	<i>Benzene</i>	
It is mainly <b>used</b> as a precursor to heavy <b>chemicals</b> , such as <b>ethylbenzene</b> and cumene , which are <b>produced</b> on a billion kilogram scale.		
2079: what is homebrew for wii?	<i>Wii homebrew</i>	
<b>Wii homebrew</b> refers to the reuse of Nintendo 's <b>Wii</b> game console to run software that has not been authorized by Nintendo .		
In more general terms, <b>Wii Homebrew</b> refers to the use of the <b>Wii</b> 's hardware, accessories and software for purposes outside those intended by the manufacturer.		
2081: who was bush's national security advisor	<i>Stephen Hadley</i>	
<b>Stephen John Hadley</b> (born February 13, 1947) was the 21st U.S. Assistant to the President for <b>National Security</b> Affairs (commonly referred to as <b>National Security Advisor</b> ), serving under President George W. <b>Bush</b> .		
2088: what is another name for the large intestine	<i>Large intestine</i>	
The <b>large intestine</b> (or bowel, colon) is the last part of the digestive system in vertebrate animals .		
The <b>large intestine</b> is about long, which is about one-fifth of the whole length of the intestinal canal .		
The <b>large intestine</b> consists of the cecum , rectum and anal canal .		
In Terminologia Anatomica the <b>large intestine</b> includes the cecum, colon, rectum, and anal canal.		
2090: where is the valley at in los angeles	<i>San Fernando Valley</i>	
The <b>San Fernando Valley</b> (locally known as "The <b>Valley</b> ") is an urbanized <b>valley</b> located in the <b>Los Angeles</b> metropolitan area of southern California , United States , defined by the mountains of the Transverse Ranges circling it.		
#2094: when is susan smith eligible for parole	<i>Susan Smith</i>	
According to the South Carolina Department of Corrections , <b>Smith</b> will be <b>eligible</b> for <b>parole</b> on November 4, 2024, after serving a minimum of thirty years.		
#2098: when was everybody hates chris made	<i>Everybody Hates Chris</i>	
Everybody <b>Hates Chris</b> is an American television period sitcom that depicts the teenage experiences of comedian <b>Chris Rock</b> (who is also the narrator ) while growing up in the Bedford-Stuyvesant neighborhood of Brooklyn, New York .		
2099: what is lung effusion	<i>Pleural effusion</i>	
Pleural <b>effusion</b> is excess fluid that accumulates between the two pleural layers , the fluid-filled <b>space</b> that surrounds the <b>lungs</b> .		
2104: who is shem in the bible	<i>Shem</i>	
<b>Shem</b> ( ; Sēm; Arabic : سام ; Sām; "renown; prosperity; name") was one of the sons of Noah in the Hebrew Bible as well as in Islamic literature . However, the New American Standard <b>Bible</b> gives, "Also to <b>Shem</b> , the father of all the children of Eber, and the older brother of Japheth, children were born."		
2105: what is go daddy.com?	<i>Go Daddy</i>	
<b>Go Daddy</b> or <b>Go Daddy Group Inc.</b> is a privately held <b>company</b> that is primarily an internet domain <b>registrar</b> and web hosting <b>company</b> .		
2106: where is andy whitfield from?	<i>Andy Whitfield</i>	
<b>Andy Whitfield</b> (died 11 September 2011) was a Welsh Australian actor and model.		
2107: where is the country andorra located	<i>Andorra</i>	
<b>Andorra</b> ( ; ), officially the Principality of <b>Andorra</b> (), also called the Principality of the Valleys of <b>Andorra</b> , () is a landlocked microstate in Southwestern Europe , located in the eastern Pyrenees mountains and bordered by <b>Spain</b> and France .		
2110: what is the highest point in oahu	<i>Oahu</i>	
The <b>highest point</b> is Mt. Ka'ala in the Waianae Range, rising to above sea level.		
2113: what is spelt flour	<i>Spelt</i>	
<b>Spelt</b> was an important staple in parts of Europe from the Bronze Age to medieval times ; it now survives as a relict crop in Central Europe and northern Spain and has found a new market as a health food.		
<b>Spelt</b> is sometimes considered a subspecies of the closely related species common wheat ( <i>T. aestivum</i> ), in which case its botanical name is considered to be <i>Triticum aestivum subsp. <i>spelta</i></i> .		
<b>Spelt</b> , also known as dinkel wheat, or hulled wheat, is an ancient species of wheat from the fifth millennium BC.		
2118: where is kennedywood in pittsburgh	<i>Kennywood</i>	
<b>Kennywood</b> is an amusement park located in West Mifflin , <b>Pennsylvania</b> , a suburb of <b>Pittsburgh</b> .		
2125: Who is the rap singer in right round with kesha in the background? *(dancer)*		<i>Right Round</i>
2128: what is rock of ages about		<i>Rock of Ages (musical)</i>
<b>Rock</b> of <b>Age</b> s is a <b>rock</b> / jukebox <b>musical</b> , with a book by Chris D'Arienzo, built around classic <b>rock</b> hits from the 1980s, especially from the famous glam metal bands of the decade.		
2134: what is the singer steve wonder full name		<i>Stevie wonder</i>
<b>Stevland Hardaway Morris</b> (born May 13, 1950 as <b>Stevland Hardaway Judkins</b> , known by his stage name <b>Steve Wonder</b> , is an American <b>singer</b> , songwriter, and multi-instrumentalist, a child prodigy who developed into one of the most creative <b>musical</b> figures of the late 20th century.		
2140: what is general chiu chicken		<i>General Tso's chicken</i>
<b>General Tso's chicken</b> (sometimes Governor <b>Tso's chicken</b> , <b>General Gau's chicken</b> , <b>General Tao's chicken</b> , <b>General Tsao's chicken</b> , <b>General Tong's chicken</b> , <b>General Tang's chicken</b> or simply <b>General's Chicken</b> ) is a sweet, slightly spicy, deep-fried <b>chicken</b> dish that is popularly served in North American Chinese restaurants.		
2142: where is david ortiz from		<i>David Ortiz</i>
<b>David Américo Ortiz</b> Arias (born November 18, 1975), nicknamed "Big Papi", is a Dominican-American professional baseball designated hitter with the Boston Red Sox of Major League Baseball (MLB).		
2143: what is an irregular heartbeat pvc		<i>Premature ventricular contraction</i>
A premature ventricular contraction (PVC), also known as a premature ventricular <b>complex</b> , ventricular premature contraction (or <b>complex</b> or <b>complexes</b> ) (VPC), ventricular premature beat (VPB), or ventricular extrasystole (VES), is a relatively common event where the <b>heartbeat</b> is initiated by Purkinje fibres in the ventricles rather than by the sinoatrial node , the normal <b>heartbeat</b> initiator.		
2150: what is flour made from		<i>Flour</i>
<b>Flour</b> is a powder which is <b>made</b> by grinding cereal grains , or other seeds or roots (like <b>Cassava</b> ).		
2156: what is cubic ft		<i>Cubic foot</i>
To calculate <b>cubic</b> feet multiply length X width X height.		
The term <b>cubic</b> foot is an Imperial and US customary (non- metric ) unit of volume , used in the United States and the United Kingdom.		
It is defined as the volume of a <b>cube</b> with sides of one foot (0.3048 m) in length .		
#2157: when was kirstie alley on cheers		<i>Kirstie Alley</i>
<b>Kirstie Louise Alley</b> (born January 12, 1951) is an American actress and comedian known for her role in the TV series <b>Cheers</b> , in which she played Rebecca Howe from 1987–1993, winning an Emmy Award and a Golden Globe Award as the Outstanding Lead Actress in a Comedy Series in 1991.		
2161: what is another name for cpu		<i>Central processing unit</i>
A central processing unit (CPU), also referred to as a central processor unit, is the <b>hardware</b> within a computer that carries out the instructions of a computer program by performing the basic arithmetical, logical, and input/output operations of the system.		
2169: what was the city of Mithridates		<i>Mithridates VI of Pontus</i>
<b>Mithridates</b> VI or Mithridates VI (), from Old Persian Mithradatha, "gift of Mithra "; 134–63 BC, also known as Mithridates the Great (Megas) and Eupator Dionysius, was king of Pontus and Armenia Minor in northern Anatolia (now Turkey ) from about 120–63 BC.		
<b>Mithridates</b> is remembered as one of the Roman Republic 's most formidable and successful enemies, who engaged three of the prominent generals from the late Roman Republic in the Mithridatic Wars : Lucius Cornelius Sulla , Lucullus and Pompey .		
2174: what is extreme right wing		<i>Far-right politics</i>
The far <b>right</b> is commonly associated with persons or groups who hold <b>extreme</b> nationalist , xenophobic , racist , religious fundamentalist or reactionary views.		
The far-right (also known as the <b>extreme right</b> ) refers to the highest degree of rightism in right-wing politics .		
2177: what is high emotional intelligence?		<i>Emotional intelligence</i>
<b>Emotional intelligence</b> (EI) is the ability to identify, assess, and control the <b>emotions</b> of oneself, of others, and of groups.		
2178: what is the role of heredity		<i>Heredity</i>
<b>Heredity</b> is the passing of traits to offspring from its parents or ancestor.		
#2186: what is prefix phone number		<i>Telephone prefix</i>
A telephone <b>prefix</b> is the first set of digits of a telephone <b>number</b> ; in the North American Numbering Plan countries (country code 1), it is the first three digits of a seven-digit <b>phone number</b> .		
2187: what is metformin used for		<i>Metformin</i>
<b>Metformin</b> ( BP , pronounced ; originally sold as <b>Glucophage</b> ) is an <b>oral</b> antidiabetic drug in the biguanide class.		
2190: what is mount rushmore?		<i>Mount Rushmore</i>
<b>Mount Rushmore</b> National Memorial is a sculpture carved into the granite face of <b>Mount Rushmore</b> near Keystone , South Dakota , in the United States .		
2195: who is keith whitely from		<i>Keith Whitley</i>
Jackie <b>Keith Whitley</b> (July 1, 1954Stambler, Irwin, and Grelun Landon (2000). - Country Music: The Encyclopedia. - New York: St. Martin's Press . p.533. - ISBN 978-0-312-26487-1 . -Carlin, Richard (2003). - Country Music: A Biographical Dictionary. - New York: Routledge . p.427. - ISBN 978-0-415-93802-0 . -Larkin, Colin (1995). - The Guinness Encyclopedia of Popular Music . - New York: pocket Books . p.395. - ISBN 978-0-85112-662-3 . -Stanton, Scott (2003). - The Tombstone Tourist. Musicians . - New York: Pocket Books . p.395. - ISBN 978-0-7434-6330-0 . -Hicks, Jack . - "Singer <b>Keith Whitley</b> 's Memory Alive Through Songs, Love in Home Town". - The Kentucky Post . - September 25, 1991 . -"Country Music Star <b>Keith Whitley</b> Dead at 33". - Lexington Herald-Leader . - May 10, 1989 . —Hurst, Jack . - "Whitley's Last Days". - Chicago Tribune . - May 14, 1989.—"Alcohol Kills Country Singer <b>Keith Whitley</b> ". - United Press International . - (c/o The San Francisco Chronicle ) . - May 10, 1989. —May 9, 1989), known professionally as <b>Keith Whitley</b> , was an American country music singer.		
2200: who is basketball star antoine walker		<i>Antoine Walker</i>
<b>Antoine Devon Walker</b> (born August 12, 1976) is an American former professional <b>basketball</b> player.		
2202: where are the internal and external iliac arteries		<i>External iliac artery</i>
The <b>external iliac arteries</b> are two major <b>arteries</b> which bifurcate off the common <b>iliac arteries</b> anterior to the sacroiliac joint of the pelvis.		
2207: where is the rhine river located on a map		<i>Rhine</i>
The <b>Rhine</b> ( ; ) is a European <b>river</b> that runs from the Swiss canton of Grisons in the southeastern Swiss Alps through <b>Germany</b> and eventually flows into the North Sea coast in the <b>Netherlands</b> and is the twelfth longest <b>river</b> in Europe , at about , with an average discharge of more than .		
2208: what is sanskrit shri		<i>Sri</i>
Sri ( Devanagari : श्री, IAST ; Śrī), also transliterated as Sree or Shri or Shree is a word of <b>Sanskrit</b> origin, used in the Indian subcontinent as polite form of address equivalent to the English "Mr." in written and spoken language, or as a title of veneration for deities (usually translated as "Holy").		
2214: WHere is a famous alluvial plain		<i>Mississippi Alluvial Plain</i>
The Mississippi River <b>Alluvial Plain</b> is an <b>alluvial plain</b> created by the Mississippi River on which lies parts of seven U.S. states , from southern Louisiana to southern Illinois .		
2215: Where is south beach in Miami		<i>South Beach</i>
<b>South Beach</b> , also nicknamed <b>SoBe</b> , is a neighborhood in the city of <b>Miami Beach</b> , Florida , United States , located due east of <b>Miami</b> city proper between Biscayne Bay and the Atlantic Ocean .		
2233: what is bones job		<i>Bones (TV series)</i>
The show is based on forensic anthropology and forensic archaeology , with each <b>episode</b> focusing on an FBI case file concerning the mystery behind human remains brought by FBI Special Agent Seeley Booth ( David Boreanaz ) to the forensic anthropologist Dr. Temperance "Bones" Brennan ( Emily Deschanel ).		
#2235: when was malcolm x assassinated		<i>Malcolm X</i>
<b>Malcolm X</b> (, May 19, 1925February 21, 1965), born <b>Malcolm Little</b> and also known as El-Hajj Malik El-Shabazz (), was an African-American Muslim minister and human rights activist.		
2239: what is penn state stadium		<i>Beaver Stadium</i>
Beaver <b>Stadium</b> is an outdoor college football <b>stadium</b> in University Park , Pennsylvania , United States , on the campus of The Pennsylvania State University .		
2240: what is a tequila sunrise?		<i>Tequila Sunrise (cocktail)</i>
The <b>Tequila Sunrise</b> is a cocktail made in two different ways, the original ( <b>tequila</b> , crème de cassis , lime juice and soda water ) and the more popular concoction ( <b>tequila</b> , orange juice , and grenadine syrup ).		

2244: where are facial sinuses	Paranasal sinuses
Paranasal sinuses are a group of four paired air-filled spaces that surround the nasal cavity ( maxillary sinuses ), above the eyes ( frontal sinuses ), between the eyes ( ethmoid sinuses ), and behind the ethmoids ( sphenoid sinuses ).	
2245: what is lean manufacturing and who developed	Lean manufacturing
Lean manufacturing, lean enterprise, or lean production, often simply, "Lean," is a production practice that considers the expenditure of resources for any goal other than the creation of value for the end customer to be wasteful, and thus a target for elimination.	
2253: what is the capital city of california.	Sacramento, California
Sacramento is the <b>capital city</b> of the U.S. state of <b>California</b> and the seat of government of Sacramento County .	
2255: What was "Freedom Summer"?	Freedom Summer
<b>Freedom Summer</b> (also known as the Mississippi Summer Project) was a campaign in the United States launched in June 1964 to attempt to register as many African American voters as possible in Mississippi , which had historically excluded most blacks from voting	
2264: what was the Name the first electronic handheld calculator	Sumlock ANITA calculator
The ANITA Mark VII and ANITA Mark VIII calculators were launched simultaneously in late 1961 as the world's first all-electronic desktop calculators .	
#2265: when was washington elected president	George Washington
George Washington (- , 1799) was the <b>first President of the United States</b> (1789–1797), the commander-in-chief of the Continental Army during the American Revolutionary War , and one of the Founding Fathers of the <b>United States</b> .	
2279: what is customary at shiva?	Shiva (Judaism)
This state lasts for seven days , during which family members traditionally gather in one home (preferably the home of the deceased) and receive visitors.	
2281: who is the guy in the wheelchair who is smart	Stephen Hawking in popular culture
Professor Stephen Hawking , known for being a theoretical physicist , has appeared in many works of popular culture .	
2283: what is in a hot toddy	Hot toddy
A hot toddy , also hot totty; and hot tottie, is typically a mixed drink made of liquor and water with sugar and spices and served hot.	
2284: what is honey bee propolis	Propolis
Propolis is a resinous mixture that <b>honey bees</b> collect from tree buds, sap flows, or other botanical sources.	
#2285: when was steven tyler born	Steven Tyler
Steven Tyler (born Steven Victor Tallarico; March 26, 1948) is an American singer , songwriter, and multi-instrumentalist, best known as the frontman and lead singer of the Boston -based rock band Aerosmith , in which he also plays the harmonica, and occasional piano and percussion.	
2286: what is name of national anthem song of switzerland	Swiss Psalm
The Swiss Psalm ( , , ) is the <b>national anthem</b> of Switzerland .	
2287: what is blood urea	Blood urea nitrogen
Normal human adult <b>blood</b> should contain between 7 to 21 mg of <b>urea</b> nitrogen per 100 ml (7–21 mg/ dl.) of <b>blood</b> .	
2291: what is corpus christi holiday	Corpus Christi (feast)
The Feast of <b>Corpus Christi</b> (Latin for Body of Christ) , also known as Corpus Domini , is a Latin Rite liturgical solemnity celebrating the tradition and belief in the body and blood of Jesus Christ and his Real Presence in the Eucharist .	
2296: what is the difference between multistage and cluster sampling	Cluster sampling
In this technique, the total population is divided into these groups (or <b>clusters</b> ) and a <b>simple random sample</b> of the groups is selected.	
<b>Cluster sampling</b> is a <b>sampling</b> technique used when "natural" but relatively homogeneous groupings are evident in a statistical population .	
2297: where is jamestown north carolina	Jamestown, North Carolina
<b>Jamestown</b> is a town in Guilford County , <b>North Carolina</b> , United States , and is a suburb of the nearby cities of Greensboro and High Point .	
2302: what is stepwise linear regression	Stepwise regression
In statistics , <b>stepwise regression</b> includes <b>regression</b> models in which the choice of predictive variables is carried out by an automatic procedure.	
#2314: when was andy griffith born	Andy Griffith
Andy Samuel Griffith (June 1, 1926 – July 3, 2012) was an American actor, television producer, Grammy Award -winning Southern-gospel singer, and writer.	
2315: what is the definition of a hung jury	Hung jury
A hung jury or deadlocked jury is a <b>jury</b> that cannot, by the required voting threshold, agree upon a verdict after an extended period of deliberation and is unable to change its votes.	
2321: what is the capacity of the cowboy stadium	Cowboys Stadium
The stadium seats 85,000, making it the third <b>largest stadium</b> in the NFL by seating <b>capacity</b> .	
2326: what is the color puce	Puce
The colors in the boxes at right are two of the various shades and varieties of <b>puce</b> .	
<b>Puce</b> (often misspelled as "puze", "peuze" or "peuce") is defined in the United States as a brownish-purple <b>color</b> .	
2333: where was jfk buried	State funeral of John F. Kennedy
After the Requiem Mass at St. Matthew's Cathedral , the late president was <b>buried</b> at Arlington National Cemetery in Virginia.	
2337: what is renaissance english	English Renaissance
The English Renaissance was a cultural and <b>artistic</b> movement in <b>England</b> dating from the late 15th and early 16th centuries to the early 17th century.	
2338: Who is the husband of Betty Ford	Betty Ford
Elizabeth Ann Bloomer Warren "Betty" Ford (April 8, 1918 – July 8, 2011), was First Lady of the United States from 1974 to 1977 during the presidency of her <b>husband</b> Gerald Ford .	
2344: what was Coco Chanel's real first name?	Coco Chanel
Gabrielle "Coco" Bonheur Chanel (August 19, 1883 – January 10, 1971) was a French fashion designer and founder of the <b>Chanel</b> brand.	
2345: what is grist mill stone	Gristmill
A <b>gristmill</b> (also: <b>grist mill</b> , corn mill or flour mill) grinds grain into flour .	
2352: where is big pokey	Big Pokey
Milton Powell (born December 4, 1977), better known by his stage name Big Pokey , is a rap artist from Houston, Texas and is one of the original members of the Screwed Up Click .	
Big Pokey joined up with DJ Screw in the early 1990s and started releasing songs on DJ Screw's many mixtapes .	
2353: what is baklava recipe	Baklava
Baklava ( , or; also Baklawa) is a rich, sweet pastry made of layers of phyllo pastry filled with chopped nuts and sweetened with syrup or honey .	
2354: what is the rule of the 9s	Total body surface area
In adults, the <b>rule of nines</b> is used to determine the total percentage of area burned for each major section of the body.	
2357: what is the sign for degrees	Degree symbol
The degree symbol (°) is a typographical symbol that is used, among other things, to represent <b>degrees</b> of arc (e.g. in geographic coordinate systems ), hours (in the medical field), or <b>degrees</b> of temperature .	
2361: who is carlos pena on big time rush	Big Time Rush
<b>Big Time Rush</b> (BTR) is an American television series created by Scott Fellows about the Hollywood misadventures of four hockey players from Minnesota —Kendall, James, <b>Carlos</b> , and Logan, after they are selected to form a boy band .	
2370: WHO IS HENRY SAMPSON JR.	Henry Sampson (inventor)
Henry T. Thomas Sampson, Jr. (born in Jackson, Mississippi in 1934) is an African-American inventor.	
2383: what is the fundamental theorem of calculus used for	Fundamental theorem of calculus
The <b>fundamental theorem of calculus</b> is a <b>theorem</b> that links the concept of the derivative of a <b>function</b> with the concept of the <b>integral</b> .	
2385: what is nicki minaj real name	Nicki Minaj
Onika Tanya Maraj (born December 8, 1982), known by her stage <b>name Nicki Minaj</b> (), is a Trinidadian-born American rapper, <b>singer</b> , songwriter and <b>television</b> personality.	
2390: what is the disease osteonecrosis of the jaw?	Osteonecrosis of the jaw
Osteonecrosis of the jaw (ONJ) is a <b>severe bone disease</b> that affects the maxilla and the mandible .	
2395: who is suicide tna	Suicide (character)
<b>Suicide</b> is a fictional character from <b>TNA Impact!</b>	
2397: where is hickory located nc	Hickory, North Carolina
<b>Hickory</b> is the principal city in the Hickory-Lenoir-Morganton MSA , in which the population at the 2010 Census was 365,497.	
Hickory is a city in Catawba County , with parts also in Burke County and Caldwell County .	
2398: what is Polyester in packaging PET	Polyethylene terephthalate

2537: what is la tour de france	<i>Tour de France</i>	The <b>Tour de France</b> () is an annual multiple stage bicycle race primarily held in <b>France</b> , while also occasionally making passes through nearby countries.
2540: what is associates arts degree	<i>Associate degree</i>	An <b>associate degree</b> is an undergraduate academic <b>degree</b> awarded by community colleges , junior colleges , technical colleges, and bachelor's <b>degree</b> -granting colleges and universities upon completion of a course of study usually lasting two years.
2543: what is vat tax?	<i>Value added tax</i>	A <b>VAT</b> is like a sales <b>tax</b> in that ultimately only the end <b>consumer</b> is <b>taxed</b> .
2544: what is high sticking in hockey	<i>High-Sticking</i>	-sticking is the name of two infractions in the sport of ice <b>hockey</b> that may occur when a player intentionally or inadvertently plays with his or her <b>stick</b> above the height of the shoulders or above the cross bar of a <b>hockey</b> goal.
#2553: when was the trojan war	<i>Trojan War</i>	The ancient <b>Greeks</b> thought that the Trojan <b>War</b> was a <b>historical event</b> that had taken <b>place</b> in the 13th or 12th <b>century BC</b> , and believed that Troy was <b>located</b> in modern-day Turkey near the Dardanelles .
2556: where is the palatine canal	<i>Greater palatine canal</i>	The greater <b>palatine canal</b> (or pterygopalatine <b>canal</b> ) is a passage in the skull that transmits the greater <b>palatine artery</b> , vein, and nerve between the pterygopalatine fossa and the oral cavity .
2557: what is the lowest temperature ever recorded in antarctica	<i>Climate of Antarctica</i>	<b>Antarctica</b> has the <b>lowest</b> naturally occurring <b>temperature</b> ever recorded on the ground on Earth: $-89.2^{\circ}\text{C}$ ( $-128.6^{\circ}\text{F}$ ) at Vostok Station .
2560: what is oregon institute of technology like	<i>Oregon Institute of Technology</i>	<b>Oregon Institute of Technology</b> , also known as <b>Oregon Tech</b> or OIT, is one of seven Universities in the <b>Oregon</b> University System , and the only public <b>institute of technology</b> in the Pacific Northwest .
#2562: When Is Passover Over	<i>Passover</i>	In Judaism , a <b>day</b> commences at dusk and lasts until the following dusk, thus the first <b>day</b> of <b>Passover</b> only begins after dusk of the 14th of Nisan and ends at dusk of the 15th <b>day</b> of the month of Nisan .
2580: WHERE IS ROUGH AND READY, CA	<i>Rough and Ready, California</i>	<b>Rough and Ready</b> is a census-designated place in Nevada County, California , United States .
2584: where is rashard lewis from???	<i>Rashard Lewis</i>	<b>Rashard Quwon Lewis</b> (born August 8, 1979 in Pineville, Louisiana ) is an American professional basketball player who currently plays for the Miami Heat of the NBA .
2590: what is doxycycline hyclate used for	<i>Doxycycline</i>	<b>Doxycycline</b> is a member of the tetracycline antibiotics group, and is commonly <b>used</b> to treat a variety of infections .
2592: what is sodium hypochlorite solution	<i>Sodium hypochlorite</i>	<b>Sodium hypochlorite solution</b> , commonly known as bleach or liquid bleach, is frequently used as a disinfectant or a bleaching agent.
#2597: when was jacques cousteau born	<i>Jacques Cousteau</i>	Jacques-Yves <b>Cousteau</b> (, commonly known in English as <b>Jacques Cousteau</b> ; 11 June 1910 – 25 June 1997) was a French naval officer, explorer , <b>conservationist</b> , filmmaker, innovator, scientist, photographer, author and researcher who studied the sea and all forms of life in water .
2607: where was the fugitive slave law made	<i>Fugitive Slave Act of 1850</i>	The <b>Fugitive Slave Law</b> or <b>Fugitive Slave Act</b> was passed by the United States Congress on September 18, 1850, as part of the Compromise of 1850 between Southern <b>slave-holding</b> interests and Northern Free-Soilers .
2612: what is eit earned income credit	<i>Earned Income Tax Credit</i>	The United States <b>federal earned income tax credit</b> or <b>earned income credit</b> (EITC or EIC) is a refundable <b>credit</b> for low- and medium-income individuals and couples, primarily for those who have qualifying children .
2620: where is valley village ca	<i>Valley Village, Los Angeles</i>	<b>Valley Village</b> is a district in the San Fernando <b>Valley</b> region of Los Angeles, California .
2630: what is sign of cancer	<i>Cancer (astrology)</i>	<b>Cancer</b> (♋) is an astrological <b>sign</b> , which is associated with the constellation <b>Cancer</b> .
2643: what is a gsm cell phone	<i>GSM</i>	<b>GSM</b> (Global System for Mobile Communications, originally '), is a <b>standard</b> set developed by the European Telecommunications <b>Standards</b> Institute (ETSI) to describe protocols for second generation ( 2G ) digital cellular networks used by mobile <b>phones</b> .
2644: what was the cash and carry lend lease	<i>Lend-Lease</i>	<b>Lend Lease</b> act was an act where the United States had supported its allies.
2646: what is el morro puerto rico	<i>Castillo San Felipe del Morro</i>	Castillo San <b>Felipe del Morro</b> also known as Fort San <b>Felipe del Morro</b> or <b>Morro</b> Castle, is a 16th-century citadel located in San Juan, <b>Puerto Rico</b> .
2651: where are poison dart frog seen	<i>Poison dart frog</i>	<b>Poison dart frog</b> (also Dart-poison <b>frog</b> , poison <b>frog</b> or formerly poison arrow <b>frog</b> ) is the common name of a group of <b>frogs</b> in the family Dendrobatidae which are native to Central and South America .
2659: where were the Winter Olympics in 2006	<i>2006 Winter Olympics</i>	The <b>2006 Winter Olympics</b> , officially known as the XX <b>Olympic Winter Games</b> , was a <b>winter</b> multi-sport event which was celebrated in Turin , Italy from February 10, 2006 , through February 26, 2006 .
2671: WHAT IS NON BINDING?	<i>Non-binding arbitration</i>	Subsequent to a non-binding arbitration, the parties remain free to pursue their claims either through the courts, or by way of a <b>binding</b> arbitration, although in practice a settlement is the most common outcome .
2675: who was mr big on sex and the city	<i>Mr. Big (Sex and the City)</i>	Non-binding arbitration is a type of arbitration in which the arbitrator makes a determination of the rights of the parties to the dispute, but this determination is not <b>binding</b> upon them, and no enforceable arbitration award is issued .
2676: who was mr big on sex and the city	<i>John James "Mr. Big"</i>	Preston is a recurring fictional <b>character</b> in the HBO series <b>Sex</b> and the <b>City</b> , portrayed by Chris Noth .
2682: what is the actresses name that played in walk that line?	<i>Walk the Line</i>	The <b>film</b> was nominated for five Academy Awards including Best <b>Actor</b> (Joaquin Phoenix), Best <b>Actress</b> (Reese Witherspoon, which she won), and Best Costume Design (Arianne Phillips).
2686: where is university of nelson mandela metropolitan located	<i>Nelson Mandela Metropolitan University</i>	<b>Nelson Mandela Metropolitan University</b> (NMMU) is a South African tertiary education institution with its main administration in the coastal city of Port Elizabeth .
#2687: when was How the west was won filmed?	<i>How the West Was Won (film)</i>	How the <b>West Was Won</b> is a 1962 American epic - <b>Western film</b> .
#2703: when was the great fire in chicago	<i>Great Chicago Fire</i>	The <b>Great Chicago Fire</b> was a conflagration that burned from Sunday, <b>October</b> 8, to early Tuesday, <b>October</b> 10, 1871, killing hundreds and destroying about in <b>Chicago</b> , Illinois .
2717: what is pci Interface	<i>Conventional PCI</i>	<b>Conventional PCI</b> (PCI) is an initialism formed from Peripheral Component Interconnect, part of the <b>PCI Local Bus standard</b> and often shortened to just <b>PCI</b> is a local <b>computer</b> bus for attaching hardware devices in a <b>computer</b> .
2718: what is jagger bombs	<i>Jägerbomb</i>	The <b>Jägerbomb</b> () is a bomb shot cocktail that was originally mixed by dropping a shot of Jägermeister into a glass of beer and in recent years evolved by the Bagheri brothers (UK) with Red Bull or other energy drinks . A long drink mixed with Jägermeister and Red Bull is called " <b>JägerBull</b> " as it is adopted from Jägermeister and RedBull.
2724: what is the name of mountains along california	<i>California Coast Ranges</i>	The other three coastal <b>California mountain</b> ranges are the Transverse Ranges , Peninsular Ranges and the Klamath <b>Mountains</b> .
2733: who is the author of tree grows in brooklyn	<i>A Tree Grows in Brooklyn (novel)</i>	A <b>Tree Grows in Brooklyn</b> is a 1943 <b>novel</b> written by Betty Smith .
2734: where is al jazeera based	<i>Al Jazeera</i>	<b>Al Jazeera</b> (, literally "The Island", abbreviating "The [Arabian] Peninsula"), also known as Aljazeera and JSC <b>Jazeera Satellite Channel</b> ), is a broadcaster owned by the privately held <b>Al Jazeera Media Network</b> and headquartered in Doha , Qatar .
2756: what is the location of coldwater ms	<i>Coldwater, Mississippi</i>	<b>Coldwater</b> is a small town in Tate County , Mississippi .
2757: what is the function of the hard palate	<i>Hard palate</i>	The <b>hard palate</b> is a thin horizontal bony plate of the skull , located in the roof of the mouth. Also on the anterior portion of the roof of the <b>hard palate</b> are the irregular ridges in the mucous membrane that help facilitate the movement of food backwards towards the pharynx.
2761: where is Chayanne from?	<i>Chayanne</i>	Elmer Figueroa Arce (born June 28, 1968), best known under the stage name <b>Chayanne</b> , is a Puerto Rican Latin pop singer and actor.
2768: what is eggning made of	<i>Egg nog</i>	As a solo artist, <b>Chayanne</b> has released 21 solo albums and sold over 15 million albums worldwide.
2773: what is gravy made of	<i>Gravy</i>	The <b>gravy</b> may be further colored and flavored with <b>gravy</b> salt (a simple mix of salt and caramel food colouring) or <b>gravy</b> browning ( <b>gravy</b> salt dissolved in water) or ready-made cubes and powders can be used as a substitute for natural meat or <b>vegetable</b> extracts.
2786: what is cta used for	<i>Computed tomography angiography</i>	Computed tomography angiography ( <b>CTA</b> ) is a computed tomography technique <b>used</b> to visualize arterial and venous vessels throughout the body.
2795: what is vitamin b12 used for	<i>Vitamin B12</i>	<b>Vitamin B12</b> , <b>vitamin B12</b> or <b>vitamin B-12</b> , also <b>called</b> cobalamin, is a water-soluble <b>vitamin</b> with a key role in the normal functioning of the brain and nervous system , and for the formation of blood.
2797: What is the North American Free Trade Agreement?	<i>North American Free Trade Agreement</i>	The <b>North American Free Trade Agreement</b> (NAFTA) is an <b>agreement</b> signed by Canada , Mexico , and the United States , creating a trilateral trade <b>block</b> in <b>North America</b> .
2798: what is tmz stand for	<i>TMZ (website)</i>	The name <b>TMZ</b> stands for the historic " studio zone " or 30-mile zone radius from the intersection of West Beverly Boulevard and North La Cienega Boulevard in Los Angeles.
2800: what is impingement of the shoulder	<i>Impingement syndrome</i>	<b>Shoulder impingement</b> syndrome, also called painful arc syndrome, supraspinatus syndrome, swimmer's <b>shoulder</b> , and thrower's <b>shoulder</b> , is a clinical syndrome which occurs when the tendons of the rotator cuff muscles become irritated and inflamed as they pass through the subacromial space, the passage beneath the acromion .
2818: who was the first one to invent medicine	<i>History of medicine</i>	The ancient Egyptians had a system of <b>medicine</b> that was very advanced for its time and influenced later medical <b>traditions</b> .
2822: Who was Daniel J Daly?	<i>Daniel Daly</i>	Sergeant Major <b>Daniel Joseph "Dan" Daly</b> (November 11, 1873 – April 27, 1937) was a United States Marine and one of only nineteen men (including seven Marines) to have received the Medal of Honor twice .
2829: what is the plot of the shawshank redemption?	<i>The Shawshank Redemption</i>	Adapted from the Stephen <b>King</b> novella Rita Hayworth and <b>Shawshank Redemption</b> , the film tells the story of Andy Dufresne, a banker who spends nearly two decades in <b>Shawshank</b> State Prison for the murder of his wife and her lover <b>despite</b> his claims of innocence.
2830: what is an information technology manager	<i>Information technology management</i>	Managing this responsibility within a company entails many of the basic management functions, like budgeting , staffing, and organizing and controlling, along with other aspects that are unique to <b>technology</b> , like change management, software design , network planning, tech support etc.
2830: what is an information technology manager	<i>Information technology management</i>	IT management is the discipline whereby all of the <b>technology</b> resources of a firm are managed in accordance with its needs and priorities.
2832: what is the scientific name of a cardinal bird	<i>Northern Cardinal</i>	The Northern <b>Cardinal</b> ( <i>Cardinalis cardinalis</i> ) is a North American <b>bird</b> in the genus <i>Cardinalis</i> ; it is also known colloquially as the redbird or common <b>cardinal</b> .
2838: where are kenworth trucks built	<i>Kenworth</i>	<b>Kenworth</b> is an American manufacturer of medium and heavy-duty Class 8 <b>trucks</b> based in Kirkland, <b>Washington</b> , United States , a suburb of Seattle, <b>Washington</b> .
#2844: when is the feast of St. Rita	<i>Rita of Cascia</i>	The Roman Catholic <b>Church</b> , under the pontificate of Pope Leo XIII officially <b>canonized</b> <b>Rita</b> on May 24, 1900, while her <b>feast day</b> is celebrated every May 22.
2846: where is cole from	<i>J. Cole</i>	Jermaine Lamarr <b>Cole</b> (born January 28, 1985), better known by his stage name <b>J. Cole</b> , is an American hip-hop recording artist and record producer from Fayetteville, North Carolina .
2856: what is the gdp for greenland 2010?	<i>Economy of Greenland</i>	GDP per capita is similar to the average European economies but the economy is critically dependent upon substantial support from the Danish government, which supplies about half the revenues of the home rule government who in turn employ about 8,000 <b>Greenlanders</b> out of a labor force of 40,156 (Jan. 2012).
2857: What is hydrogen in	<i>Hydrogen</i>	Most of the <b>hydrogen</b> on Earth is in molecules such as water and organic compounds because <b>hydrogen</b> readily forms covalent compounds with most non-metallic elements.
In 1766–81, Henry Cavendish was the first to recognize that <b>hydrogen</b> gas was a discrete substance, and that it produces water when burned, a property which later gave it its name: in Greek, <b>hydrogen</b> means "water-former".		
2858: where is UWA on world list of universities	<i>University of Western Australia</i>	One of <b>Australia</b> 's best and most prestigious <b>universities</b> , <b>UWA</b> is highly ranked <b>internationally</b> in various <b>publications</b> . The 2011 QS <b>World University Rankings</b> placed <b>UWA</b> at 73rd <b>internationally</b> .
#2878: when was purple haze by jimi hendrix made?	<i>Purple haze</i>	"Purple Haze" is a song written and recorded by <b>Jimi Hendrix</b> in 1967, released as the second single by The <b>Jimi Hendrix Experience</b> in both the United Kingdom and the United States .
2880: where was martin luther born	<i>Martin Luther</i>	<b>Martin Luther</b> (, 10 November 1483 – 18 February 1546) was a German monk , former Catholic priest , professor of theology and seminal figure of a reform movement in sixteenth century Christianity , subsequently known as the <b>Protestant Reformation</b> .
2891: what is the mortality rate of sepsis	<i>Septic shock</i>	The <b>mortality</b> rate from <b>septic</b> shock is approximately 25–50%.
2893: where is kos from?	<i>Kos</i>	The alias "K-os" , spelled with a lower case "k" , was intended to be less aggressive than the pseudonyms of other rappers whose names were all upper case, such as KRS-One .
K-os usually performs with a live band, something that is uncommon in the hip hop genre.		
K-os received his first musical exposure with the single "Musical Essence", released in 1993.		
Kevin Brereton (born February 20, 1972), better known by his stage name <b>k-os</b> (, "chaos"), is a Canadian rapper , singer , songwriter and record producer .		
2904: what is the sign called?	<i>At sign</i>	It is an acronym for "Knowledge of Self" , although in a later interview he said that it originally stood for "Kevin's Original Sound." <b>k-os</b> music incorporates a wide variety of music genres, including rap , funk , rock , and reggae .
2905: who is the book the catcher in the rye by?	<i>The Catcher in the Rye</i>	A musician as well as a producer, <b>k-os</b> has written and produced nearly every part of all four of his albums.
The at <b>sign</b> is also commonly called the at symbol, ampersat, apetail or commercial at in <b>English</b> —and less commonly a wide range of other terms.		
2906: who is the book the catcher in the rye by?	<i>The Catcher in the Rye</i>	The <b>Catcher in the Rye</b> is a 1951 novel by J. D. <b>Salinger</b> .
2910: what is prince williams last name	<i>Prince William, Duke of Cambridge</i>	Prince <b>William</b> , Duke of <b>Cambridge</b> (William Arthur Philip Louis; born 21 June 1982), is the elder son of Charles, Prince of Wales , and Diana, Princess of Wales , and third-eldest grandchild of Queen Elizabeth II and Prince Philip, Duke of Edinburgh .
2914: what is IBRIX	<i>IBRIX Fusion</i>	The software was produced, sold, and supported by <b>IBRIX</b> Incorporated of Billerica, Massachusetts . Subsequent to the acquisition, the software components of <b>IBRIX</b> have been combined with ProLiant servers to form the X9000 series of storage systems.
<b>IBRIX</b> Fusion is a scalable parallel file system combined with integrated logical volume manager , availability features and a management interface.		
HP announced on July 17, 2009 that it had reached a definitive agreement to acquire <b>IBRIX</b> .		
2915: what is bilirubin total	<i>Bilirubin</i>	<b>Bilirubin</b> is excreted in bile and urine , and elevated levels may indicate certain diseases.
<b>Bilirubin</b> (formerly referred to as hematoidin) is the yellow breakdown product of normal heme catabolism .		
2918: what is quasi judicial agency	<i>Quasi-judicial body</i>	

294: who is the founder of twitter	Twitter	
Twitter was created in March 2006 by Jack Dorsey and by July, the social networking site was launched.		
Twitter is an online social networking service and microblogging service that enables its users to send and read text-based messages of up to 140 characters , known as "tweets".		
Twitter Inc is based in San Francisco , with additional servers and offices in New York City , Boston , and San Antonio . Since its launch, Twitter has become one of the ten most visited websites on the Internet, and has been described as "the SMS of the Internet."	Bj's Wholesale Club	
2932: who is bj's wholesale club	Bj's Wholesale Club	
Bj's Wholesale Club, Inc., commonly referred to simply as <b>Bj's</b> , is a membership-only warehouse club chain operating on the United States East Coast , as well as in the state of Ohio .		
2936: what is the name of the wizard of oz	Wizard of Oz (character)	
The Wizard of Oz , known during his reign as The Great and Powerful Oz , is the epithet of Oscar Zoroaster Phadrig Isaac Norman Henkel Emmmanuel Ambroise Diggs, a fictional character in the Land of Oz , created by American author L. Frank Baum .		
2940: what is homebrew for wii	Wii homebrew	
Wii homebrew refers to the reuse of Nintendo 's Wii game console to run software that has not been authorized by Nintendo		
In more general terms, Wii Homebrew refers to the use of the Wii 's hardware, accessories and software for purposes outside those intended by the manufacturer.		
2952: what is korean money called	South Korean won	
The won () sign : ₩; code : KRW is the currency of South Korea .		
2955: what is three phase electrical	Three-phase electric power	
Three-phase electric power is a common method of alternating-current electric power generation , transmission , and distribution .		
2957: what is the significance of Good Friday?	Good Friday	
It is also known as <b>Holy Friday</b> , Great <b>Friday</b> , Black <b>Friday</b> , or Easter <b>Friday</b> , though the latter properly refers to the <b>Friday in Easter</b> week .		
2959: where is the great basin located on a us map	Great Basin	
It is noted for both its arid conditions and its <b>Basin</b> and <b>range</b> topography that varies from the North American low point at Badwater <b>Basin</b> to the highest point of the contiguous United States , less than away at the summit of Mount Whitney .		
2966: what was the date of pearl harbor	Attack on Pearl Harbor	
The attack on Pearl Harbor (called Hawaii Operation or Operation AI by the Japanese Imperial General Headquarters (Operation Z in planning) and the Battle of Pearl Harbor ) was a surprise military strike conducted by the Imperial Japanese Navy against the United States naval base at Pearl Harbor , Hawaii, on the morning of December 7, 1941 (December 8 in Japan).		
2969: what is disney's magic kingdom	Magic Kingdom	
Magic Kingdom Park , also known as <b>Magic Kingdom</b> , is the first of four theme parks built at the Walt Disney World Resort in Bay Lake , Florida .		
2972: what is the kanji for language	Kanji	
Kanji ( ) are the adopted logographic Chinese characters ( hanzi ) that are used in the modern Japanese writing system along with hiragana (ひらがな, ), katakana (カタカナ, ), Hindu-Arabic numerals , and the occasional use of the Latin alphabet .		
2973: what is in fruitcake	Fruit cake	
Fruit cake (or fruitcake) is a cake made with chopped candied fruit and/or dried fruit , nuts , and spices , and (optionally) soaked in spirits .		
2978: who was the congressman who was caught with an escort in ny	Eliot Spitzer prostitution scandal	
		On March 10, 2008, The New York Times reported that New York Governor Eliot Spitzer had patronized an elite escort service run by Emperors Club VIP .
		2982: what is lockton affinity
		Lockton Companies
		Lockton Affinity: Lockton Affinity, an affiliate of Lockton Companies, meets the insurance needs of affinity groups, franchises, professional organizations, and associations of all sizes.
		2991: what is the name of chris cornell's band?
		Chris Cornell
		Chris Cornell (born Christopher John Boyle; July 20, 1964) is an American rock <b>musician</b> best known as the lead vocalist and rhythm <b>guitarist</b> for Soundgarden and as the former lead vocalist for Audioslave .
		2994: what is preciosia crystal?
		Preciosa (corporation)
		Preciosa is the luxury brand name for the range of precision-cut lead <b>crystal</b> glass and related products produced by Preciosa a.s. of Jablonec nad Nisou , Czech Republic .
		2998: where is the chupacabra found
		Chupacabra
		The Chupacabra (, from chupar "to suck" and cabra "goat", literally "goat sucker") is a legendary cryptid rumored to inhabit parts of the Americas .
		2999: What is the purpose of North American Free Trade Agreement
		North American Free Trade Agreement
		The North American Free Trade Agreement (NAFTA) is an agreement signed by Canada , Mexico , and the United States , creating a trilateral trade bloc in North America .
		3004: who are all of the jonas brothers
		Jonas Brothers
		Formed in 2005, they gained popularity from the Disney Channel children's television network and consists of three brothers from Wyckoff, New Jersey : Paul Kevin Jonas II , Joseph Adam Jonas and Nicholas Jerry Jonas .
		3012: where is the brisket from
		Brisket
		The brisket muscles include the superficial and deep pectorals.
		The beef brisket is one of the nine beef prime cuts.
		Brisket is a cut of meat from the breast or lower chest of beef or veal .
		American cuts of beef including the brisket
		According to the Random House Dictionary of the English Language , Second Edition, the term derives from the Middle English <b>brusket</b> which comes from the earlier Old Norse , meaning cartilage .
		3016: what is human chorionic
		Human chorionic gonadotropin
		In molecular biology , human chorionic gonadotropin (hCG) is a hormone produced by the fertilized egg after conception .
		3021: what is busiest airport in us
		List of the busiest airports in the United States
		Hartsfield-Jackson Atlanta International Airport is the busiest single airport in the United States .
		3027: where is osaka japan
		Osaka
		is a city in the Kansai region of Japan 's main island of Honshu , a designated city under the Local Autonomy Law , the capital city of Osaka Prefecture and also the largest part of the Keihanshin metropolis, which comprises three major cities of Japan Kyoto , Osaka and Kobe .
		3030: what is surname and given name
		Given name
		A given name , in Western contexts often referred to as a <b>first name</b> , is a personal name that specifies and differentiates between members of a group of individuals, especially in a family , all of whose members usually share the same family name (surname).
		3037: What is an economic feature?
		Economics
		Other broad distinctions include those between positive economics (describing "what is") and normative economics (advocating "what ought to be"); between economic theory and applied economics ; between rational and behavioral economics ; and between mainstream economics (more "orthodox" and dealing with the "rationality-individualism-equilibrium nexus") and heterodox economics (more "radical" and dealing with the "institutions-history-social structure nexus").
		3043: what is section eight housing
		Section 8 (housing)
		Section 8 of the Housing Act of 1937 (), often simply known as Section 8, as repeatedly amended, authorizes the payment of rental housing assistance to private landlords on behalf of approximately 3.1 million low-income households .

*Example 4.15.4.* In Table S12 below, we list the questions where our top-scoring answers have all been officially labeled as incorrect by the WikiQA team, including 167 “quantitative” questions, out of 662 total. It should be noted that some questions in this table are simply too tricky to be handled by literal interpretations. For example, our top-scoring answer to Q1838 differs from the official one (shown in blue):

1838: who sang cool jerk	Cool Jerk
In the feature film Home Alone 2: Lost in New York , Uncle Frank ( Gerry Bamman ) sings "Cool Jerk" in the shower.	
"Cool Jerk" is a popular song written by Donald Storball and originally performed by The Capitols .	
Since our algorithm fails to reinterpret "who sang" as "who originally performed", our answer is outshone by the official correct answer. Doubtlessly, this particular entry in Table S12 illustrates the limitations of our automated question comprehension.	
<i>Example 4.15.5.</i> It is still arguable whether dozens of entries in Table S12 represent miscarriage of justice by the WikiQA team. In the cases selected below, we believe that our top-scoring answers are not any worse than the official ones (painted in blue):	
#129: how many stripes on the flag	Flag of the United States
The national flag of the United States of America , often simply referred to as the American flag , consists of thirteen equal horizontal stripes of red (top and bottom) alternating with white , with a blue rectangle in the canton (referred to specifically as the "union") bearing fifty small, white, five-pointed stars arranged in nine offset horizontal rows of six stars (top and bottom) alternating with rows of five stars .	
The 50 stars on the flag represent the 50 states of the United States of America and the 13 stripes represent the thirteen British colonies that declared independence from the Kingdom of Great Britain and became the first states in the Union .	
214: what does salesforce.com do?	Salesforce.com
Salesforce.com Inc. is a global enterprise software company headquartered in San Francisco , California .	
Though best known for its customer relationship management (CRM) product, Salesforce has also expanded into the "social enterprise arena" through acquisitions.	
#557: how many spanish speaking people are there in the us	Spanish language in the United States
There are 45 million Hispanophones who speak Spanish as a first or second language , as well as six million Spanish language students, composing the largest national Spanish-speaking community outside of Mexico and making Spanish the Romance language and the Indo-European language with largest number of native speakers in the world .	
According to the 2010 American Community Survey conducted by the U.S. Census Bureau , Spanish is the primary language spoken at home by almost 37 million people aged five or older, a figure more than double that of 1990 .	
1831: where do you find iodine	Iodine
Because of this function, radioisotopes of iodine are concentrated in the thyroid gland along with nonradioactive iodine .	
Iodine is found on Earth mainly as the highly water-soluble iodide ion, I <sup>-</sup> , which concentrates it in oceans and brine pools .	
2463: who is E from entourage	Eric Murphy
Murphy is a fictional character on the comedy-drama television series Entourage .	
Eric "E"	
2532: what is the official language of america?	Languages of the United States
The situation is quite varied at the state and territorial levels, with some states mirroring the federal policy of adopting no official language in a de jure capacity, others adopting English alone, others officially adopting English as well as local languages , and still others adopting a policy of de facto bilingualism.	
The most commonly used language is English .	
2841: what is the type of democracy in which all citizens have the right to make government decisions	Democracy
One form of democracy is direct democracy , in which all eligible citizens have direct and active participation in the decision making of the government .	
Democracy is a form of government in which all eligible citizens have an equal say in the decisions that affect their lives.	

In particular, the penultimate official answer listed above (accompanying Q2532) is logically irrelevant,<sup>48</sup> if not factually misleading; the last official answer listed above (accompanying Q2841) does not even address the “type of democracy” in the question (*i.e.* direct, representative *etc.*) at all. Despite our disagreement over certain judgements made by the WikiQA team, we are not going to alter their official labels for correct/incorrect answers, so as to ensure a fair comparison in [1, Fig. 7b] and Fig. S2.

**Table S12. List of WikiQA misses**

4: how a water pump works	Pump	The <b>Big Ten Conference</b> , formerly Western <b>Conference</b> and <b>Big Nine Conference</b> , is the oldest <b>Division I</b> college athletic conference in the United States .
4: how a water pump works	Pump	Sound mass
<b>Pumps</b> can be classified into three major groups according to the method they use to move the fluid: direct lift, displacement, and gravity <b>pumps</b> .		In contrast to more traditional musical textures , <b>sound mass composition</b> “minimizes the importance of individual pitches in preference for texture , timbre , and dynamics as primary shapers of gesture and impact.”
#11: how big is bmc software in houston, tx	BMC Software	Spades
Headquartered in <b>Houston</b> , Texas , <b>BMC</b> develops, markets and sells <b>software</b> used for multiple functions, including IT service management, data center <b>automation</b> , <b>performance</b> management, virtualization lifecycle management and cloud computing management.		The continent of <b>Australia</b> lies on a continental shelf overlain by shallow seas which divide it into several landmasses — the Arafura Sea and Torres Strait between <b>mainland Australia</b> and New Guinea, and Bass Strait between <b>mainland Australia</b> and Tasmania.
#17: how much are the harry potter movies worth	Harry Potter	Cardiovascular disease
Also due to the <b>success</b> of the books and films, <b>Harry Potter</b> has been used for a theme <b>park</b> . The <b>Wizarding World of Harry Potter</b> in <b>Universal Parks &amp; Resorts</b> ' Islands of Adventure .		Cardiovascular disease refers to any <b>disease</b> that affects the cardiovascular system , principally <b>cardiac disease</b> , vascular <b>diseases</b> of the brain and kidney , and peripheral <b>arterial disease</b> .
#20: how old was sue lyon when she made lolita	Lolita (1962 film)	History of immigration to the United States
The film stars <b>James Mason</b> as Humbert Humbert, <b>Sue Lyon</b> as Dolores Haze ( <b>Lolita</b> ) and Shelley Winters as Charlotte Haze with Peter Sellers as Clare Quilty.		During the nation's history , the growing <b>country</b> experienced successive waves of <b>immigration</b> which rose and fell over time, particularly from Europe, with the cost of transoceanic transportation sometimes paid by travelers becoming indentured servants after their arrival in the New World.
21: how are cholera and typhus transmitted and prevented	Cholera	Atlanta
<b>Cholera</b> is an infection in the small intestine caused by the bacterium <b>Vibrio cholerae</b> .		<b>Atlanta</b> is the primary transportation hub of the Southeastern United States , via highway, railroad, and air , with Hartsfield-Jackson <b>Atlanta</b> International Airport being the world's <b>busiest</b> airport since 1998.
28: how are aircraft radial engines built	Radial engine	Commonwealth
The <b>radial configuration</b> was very commonly used in <b>aircraft engines</b> before turbine <b>engines</b> became predominant.		Most notably, the <b>Commonwealth of Nations</b> , an association primarily of former members of the British Empire , is often referred to as simply “the <b>Commonwealth</b> ”.
#30: how deep can be drill for deep underwater	Deepwater drilling	Supreme Court of the United States
There are basically two kinds of mobile <b>deepwater drilling</b> rigs: semi-submersible <b>drilling</b> rigs and <b>drillships</b> .		It has ultimate (and largely discretionary ) appellate <b>jurisdiction</b> over all <b>federal courts</b> and over <b>state court cases</b> involving issues of <b>federal law</b> , and original <b>jurisdiction</b> over a small range of <b>cases</b> .
#32: how long was frank sinatra famous	Frank Sinatra	Microsoft SQL Server
<b>Sinatra</b> left Capitol to found his own record label, Reprise Records in 1961 (finding success with albums such as Ring-a-Ding-Ding! , <b>Sinatra</b> at the Sands and <b>Francis Albert Sinatra &amp; Antonio Carlos Jobim</b> ), toured internationally, was a founding member of the Rat Pack and fraternized with celebrities and statesmen, including John F. <b>Kennedy</b> .		There are at least a dozen different <b>editions</b> of Microsoft SQL Server aimed at different audiences and for different workloads (ranging from small applications that store and retrieve <b>data</b> on the same computer, to millions of users and computers that access huge amounts of <b>data</b> from the Internet at the same time).
#42: how much is jk rowling worth	J. K. Rowling	Oklahoma City bombing
Joanne "Jo" <b>Rowling</b> (), born 31 July 1965), pen name J. <b>K. Rowling</b> , is a British novelist, best known as the author of the Harry Potter fantasy series .		The <b>Oklahoma City bombing</b> was a domestic terrorist <b>bomb</b> attack on the Alfred P. Murrah <b>Federal Building</b> in downtown <b>Oklahoma City</b> on April 19, 1995.
#43: how big is Auburndale Florida	Auburndale, Florida	ARIA Music Awards
<b>Auburndale</b> is a city in Polk County, <b>Florida</b> , United States .		The <b>ARIA Recording Industry Association Music Awards</b> (commonly known as <b>ARIA Music Awards</b> or <b>ARIA Awards</b> ) is an annual series of awards nights celebrating the <b>Australian music industry</b> , put on by the <b>Australian Recording Industry Association</b> (ARIA).
#46: how old is the singer bob seger	Bob Seger	Benedict arnold
As a locally successful Detroit-area artist, he performed and recorded as <b>Bob Seger</b> and the Last Heard and <b>Bob Seger</b> System throughout the 1960s.		In the winter of 1782, <b>Arnold</b> moved to London with his second wife, <b>Margaret "Peggy" Shippem Arnold</b> .
#48: how long was i love lucy on the air	I Love Lucy	Port (computer networking)
<b>I Love Lucy</b> is an American <b>television sitcom</b> starring Lucille Ball , Desi Arnaz , Vivian Vance , and William <b>Frawley</b> .		In the client-server model of application architecture, <b>ports</b> are used to provide a multiplexing service on each server-side <b>port</b> number that <b>network</b> clients connect to for service initiation, after which communication can be reestablished on other connection-specific <b>port</b> numbers.
#59: HOW MUCH IS CENTAVOS IN MEXICO	Mexican peso	Comparison of the health care systems in Canada and the United States
The name was originally used in reference to <b>pesos</b> oro (gold weights) or <b>pesos</b> plata (silver weights).		Total <b>government spending</b> per capita in the U.S. on <b>health care</b> was 23% higher than <b>Canadian government spending</b> , and U.S. <b>government expenditure</b> on <b>health care</b> was just under 83% of total <b>Canadian spending (public and private)</b> though these statistics don't take into account population differences.
66: how did armando christian perez become famous	Pitbull (entertainer)	List of Xbox 360 games
In 2004, he released his debut album titled M.I.A.M.I. (short for Money Is A Major Issue) under T.V.T Records .		For a list of downloadable <b>Xbox Live Arcade games</b> , see the List of <b>Xbox Live Arcade games</b> .
#75: how old were the twin towers when destroyed	World Trade Center	History of wine
At the time of their completion, the original 1 World Trade Center (the North <b>Tower</b> ) and 2 World Trade Center (the South Tower), known collectively as the “ <b>Twin Towers</b> ”, were the <b>tallest buildings</b> in the world.		Following the decline of Rome and its industrial-scale <b>wine</b> production for <b>export</b> , the Christian Church in medieval Europe also became a firm supporter of <b>wine</b> , necessary for celebration of the Catholic Mass. Whereas <b>wine</b> was forbidden in medieval Islamic cultures, its use in Christian liturgy was widely tolerated.
#89: how many presidents of the us	List of Presidents of the United States	List of districts of West Bengal
Upon the death, resignation, or removal from <b>office</b> of an incumbent <b>President</b> , the <b>Vice President</b> assumes the <b>office</b> .		Geographically, <b>West Bengal</b> is divided into a variety of <b>regions</b> —Darjeeling Himalayan hill <b>region</b> , Terai and Dooars <b>region</b> , North <b>Bengal plains</b> , Rarhi <b>region</b> , Western plateau and high lands , coastal <b>plains</b> , Sunderbans and the Ganges Delta .
#100: what happens to the light independent reactions of photosynthesis?	Light-independent reactions	List of cities and towns in New Hampshire
This <b>process happens</b> when <b>light</b> is available <b>independent</b> of the kind of <b>photosynthesis</b> (C3 carbon fixation , C4 carbon fixation , and Crassulacean Acid Metabolism ); CAM plants store malic acid in their vacuoles every night and release it by day in order to make this <b>process</b> work.		However, <b>towns</b> currently are able to change their form of government by simple voter approval of a <b>new</b> municipal charter, with several of the more populous <b>towns</b> having already done so.
107: what countries are under the buddhism religion	Buddhism by country	Salesforce.com
Government policies in these <b>countries</b> may encourage the under-reporting or non-reporting of <b>religious</b> adherence, resulting in official totals that may drastically underestimate the number of <b>religious</b> practitioners in these <b>countries</b> .		Salesforce.com Inc. is a global enterprise software company <b>headquartered</b> in San Francisco , California .
108: how did wild bill's father die	Wild Bill Hickok	Earth
James Butler Hickok (May 27, 1837 – August 2, 1876), better known as <b>Wild Bill</b> Hickok, was a folk hero of the American Old West .		The Earth's <b>axis of rotation</b> is tilted 23.4° away from the perpendicular of its orbital <b>plane</b> , producing <b>seasonal</b> variations on the <b>planet</b> 's surface with a <b>period</b> of one tropical year (365.24 solar days).
#109: how many land rovers have landed on mars	Mars rover	Cellular respiration
It currently manages the <b>Mars Exploration Rover</b> mission's active Opportunity <b>rover</b> and inactive Spirit , and, as part of the Mars Science Laboratory mission, the Curiosity <b>rover</b> .		While the overall <b>reaction</b> is a combustion <b>reaction</b> , no single <b>reaction</b> that comprises it is a combustion <b>reaction</b> .
112: what bird family is the owl	Owl	Fire extinguisher
<b>Owls</b> hunt mostly small mammals , insects , and other <b>birds</b> , although a few species specialize in <b>hunting</b> fish .		A <b>fire extinguisher</b> , flame <b>extinguisher</b> , or simply an <b>extinguisher</b> , is an active <b>fire</b> protection device used to <b>extinguish</b> or control small <b>fires</b> , often in emergency situations.
113: how did harmon killebrew get strong	Harmon Killebrew	Central america
<b>Harmon Clayton Killebrew</b> ( June 29 , 1936 – May 17, 2011), nicknamed “ <b>Killer</b> ” and “Hammerin' <b>Harmon</b> ” , was an American professional baseball first baseman , third baseman , and left fielder .		<b>Central America</b> () is the central geographic region of the <b>Americas</b> .
116: What did the augurs use to interpret the will of the gods?	Augur	Steam engine
An <b>augur</b> holding a litmus , the curved wand often <b>used</b> as a symbol of <b>augury</b> on Roman coins		In general usage, the term <b>steam engine</b> can refer to either the integrated <b>steam</b> plants (including boilers etc.) such as railway <b>steam locomotives</b> and portable <b>engines</b> , or may refer to the piston or turbine machinery alone, as in the beam <b>engine</b> and stationary <b>steam engine</b> .
#117: How much did Waterboy grossed	The Waterboy	Steam engine
The <b>Waterboy</b> is a 1998 American sports/comedy film directed by Frank Coraci (who played Robert ‘‘Roberto’’ Boucher, Sr.), starring Adam <b>Sandler</b> , Kathy Bates , Fairuza Balk , Henry Winkler , Jerry Reed , Larry Gilliard, Jr. , Blake Clark , Peter Dante and Jonathan Loughran , and produced by Robert Simonds and Jack Giarraputo		Since the 1990s, developments in the downtown core include the <b>University of Washington Tacoma</b> ; <b>Tacoma Link</b> , the first modern electric light rail service in the <b>state</b> ; the <b>state's highest</b> density of <b>art</b> and <b>history museums</b> ; and a restored urban waterfront, the Thea Foss Waterway .
118: what county is Farmington Hills, MI in?	Farmington Hills, Michigan	Chula Vista, California
Although the two cities have separate services and addresses, <b>Farmington</b> and <b>Farmington Hills</b> are often thought of as the same community.		Located in the city is one of America's few year-round United States Olympic Training <b>centers</b> and popular tourist destinations include Cricket Wireless Amphitheater, the <b>Chula Vista</b> marina, and the <b>Chula Vista Nature Center</b> .
121: what does a groundhog look for on groundhog day	Groundhog Day	Isaac Newton
<b>Groundhog Day</b> , already a widely recognized and popular tradition, received widespread attention as a result of the 1993 film <b>Groundhog Day</b> , which was set in Punxsutawney, PA (though filmed primarily in Woodstock, IL) and portrayed Punxsutawney Phil .		In addition to his work on the calculus, as a <b>mathematician</b> <b>Newton</b> contributed to the study of power series , generalised the binomial theorem to non-integer exponents, and developed <b>Newton's</b> method for approximating the roots of a function .
#129: how many stripes on the flag	Flag of the United States	Casualties of the Iraq War
The national <b>flag</b> of the United States of America , often simply referred to as the American <b>flag</b> , consists of <b>thirteen</b> equal horizontal <b>stripes</b> of <b>red</b> (top and bottom) <b>alternating</b> with <b>white</b> , with a <b>blue</b> rectangle in the canton (referred to specifically as the “union”) bearing fifty small , <b>white</b> , <b>five-pointed stars</b> arranged in nine offset horizontal rows of <b>six stars</b> (top and bottom) <b>alternating</b> with rows of <b>five stars</b> .		<b>Casualties</b> of the conflict in <b>Iraq</b> since 2003 (beginning with the 2003 invasion of <b>Iraq</b> , and continuing with the ensuing occupation of <b>Iraq</b> , as well as the activities of the various <b>armed</b> groups operating in the country) have come in many forms , and the accuracy of the information available on different types of <b>Iraq War casualties</b> varies greatly.
132: how is whooping cough distinguished from similar diseases	Pertussis	
In some countries, this <b>disease</b> is called the 100 days' <b>cough</b> or <b>cough</b> of 100 days.		
136: what county is galveston in texas	Galveston County, Texas	
League City is the largest city in <b>Galveston County</b> in terms of population; URL_http between 2000 and 2005 it surpassed <b>Galveston</b> as the <b>county's</b> largest city.		
137: what cities are in the bahamas	List of cities in the Bahamas	
This is a list of <b>cities</b> in the <b>Bahamas</b> .		
#139: how many schools are in the big ten	Big Ten Conference	

<sup>48</sup> Think about this: Punjabi is the most commonly used language in Pakistan, whose official languages are English and Urdu.

#255: how many professional hockey teams in canada	National Hockey League
The National Hockey League (NHL) is an "unincorporated not-for-profit association" which operates a major professional ice hockey league of 30 franchised member clubs, of which seven are currently located in Canada and 23 in the United States.	
#267: how many games did brett favre start in a row	Brett Favre
Favre became the Packers' starting quarterback in the fourth game of the 1992 season , stepping in for injured quarterback Don Majkowski , and started every game through the 2007 season .	
273: what does a vote to table a motion mean?	Table (parliamentary procedure)
In the rest of the English-speaking world such as the United Kingdom , to table means to move to place [the topic] upon the table (or to move to place on the table): a proposal to begin consideration (or reconsideration) of a proposal.	
276: What does the class mean for SDHC cards?	Secure Digital Host
Host devices that comply with newer versions of the specification provide backward compatibility and accept older SD cards , but this article explains several factors that can prevent the use of a newer SD card:	
#295: how many apple store are there in total?	Apple Store
The Apple Store is a chain of retail stores owned and operated by Apple Inc. , dealing in computers and consumer electronics.	
300: how is jerky made	Jerky
Some makers still use just salt and sun-dry fresh sliced meat to make jerky.	
#305: what happened on the moon during the period of Late Heavy Bombardment?	Late Heavy Bombardment
The Late Heavy Bombardment (commonly referred to as the lunar cataclysm, or LHB) is a hypothetical event around 4.1 to 3.8 billion years ago ( Ga ).	
307: How did Edgar Allan Poe die?	Edgar Allan Poe
Edgar Allan Poe (born Edgar Poe; January 19, 1809 – October 7, 1849) was an American author, poet, editor and literary critic, considered part of the American Romantic Movement .	
#311: how many seasons heroes	Heroes (TV series)
The second season of Heroes attracted an average of 13.1 million viewers in the U.S., and marked NBC's sole series among the top 20 ranked programs in total viewership for the 2007–2008 season.	
#318: how old r Dylan and Cole Sprouse	Dylan and Cole Sprouse
They are twins and are collectively referred to as Dylan and Cole Sprouse or the Sprouse brothers, usually abbreviated as Sprouse Bros.	
#321: how many days does the chinese new year last	Chinese New Year
Chinese New Year	
324: what does s.h.i.e.l.d stand for	S.H.I.E.L.D.
S.H.I.E.L.D. is a fictional espionage and law-enforcement agency in the Marvel Comics Universe .	
#326: how many consoles has xbox 360 sold	Xbox 360
Several major features of the Xbox 360 are its integrated Xbox Live service that allows players to compete online ; download arcade games , game demos, trailers, TV shows, music and movies; and its Windows Media Center multimedia capabilities.	
328: WHAT COUNTRY IS MEXICO IN	Mexico
Mexico ranks sixth in the world and first in the Americas by number of UNESCO World Heritage Sites with 31 , and in 2007 was the tenth most visited country in the world with 21.4 million international arrivals per year.	
#330: how is slugging percentage calculated	Slugging percentage
The next year he slugged .846, and these records went unbroken until 2001 , when Barry Bonds achieved 411 bases in 476 at-bats, bringing his slugging percentage to .863, unmatched since.	
335: what does the federal reserve do	Federal Reserve System
The Federal Reserve System's structure is composed of the presidentially appointed Board of Governors (or Federal Reserve Board), the Federal Open Market Committee (FOMC), twelve regional Federal Reserve Banks located in major cities throughout the nation, numerous privately owned U.S. member banks and various advisory councils.	
#340: how many books are included in the protestant Bible?	Books of the Bible
The first part of Christian Bibles is the Old Testament , which contains, at minimum, the twenty-four books of the Hebrew Bible divided into thirty-nine books and ordered differently from the Hebrew Bible.	
#341: how many stars on the first american flag	Betsy Ross flag
Although the Betsy Ross story is accepted by most Americans , some flag historians and revisionists do not accept the Betsy Ross design as the first American flag.	
#345: what became of rich on price is right	Rich Fields
Richard Wayne "Rich" Fields is an American broadcaster, spokesman, announcer and meteorologist , best known for being the announcer of the American version of The Price Is Right from 2004–2010.	
#352: what happened to stevie ray vaughan	Stevie Ray Vaughan
As the younger brother of Jimmie Vaughan , Vaughan started playing the guitar at age seven and formed several bands that occasionally performed in local nightclubs.	
358: what caused ww	Causes of World War I
Most historians and popular commentators include causes from more than one category of explanation to provide a rounded account of the causes of the war.	
#362: how many seasons were there of the wire	The Wire
Despite only receiving modest ratings and never winning major television awards, The Wire has been described by many critics as one of the greatest TV dramas of all time.	
364: what kind of legal remedy is it to ask someone to fulfill their promise	Contract
Contract law varies greatly from one jurisdiction to another, including differences in common law compared to civil law , the impact of received law , particularly from England in common law countries, and of law codified in regional legislation.	
365: how is ASPNET different from .NET	ASP.NET
ASP.NET is built on the Common Language Runtime (CLR), allowing programmers to write ASP.NET code using any supported .NET language .	
#370: how many people live in memphis tennessee	Memphis, Tennessee
A resident of Memphis is referred to as a Memphian , and the Memphis region is known, particularly to media outlets, as "Memphis & The Mid-South".	
#372: how many muscles in the human body	List of muscles of the human body
The muscles of the human body can be categorized into a number of groups which include muscles relating to the head and neck, muscles of the torso or trunk, muscles of the upper limbs, and muscles of the lower limbs.	
373: how did seminole war end	Seminole Wars
The Seminole Wars , also known as the Florida Wars , were three conflicts in Florida between the Seminole — the collective name given to the amalgamation of various groups of native Americans and Black people who settled in Florida in the early 18th century — and the United States Army .	
#385: how many nature oceans are on earth	Ocean
This animation uses Earth science data from a variety of sensors on NASA Earth observing satellites to measure physical oceanography parameters such as ocean currents , ocean winds , sea surface height and sea surface temperature.	
388: what county is jacksonville florida in	Jacksonville, Florida
Jacksonville is the principal city in the Jacksonville metropolitan area , with a population of 1,345,596 in 2010.	
391: what does a liquid oxygen plant look like	Liquid oxygen
Liquid oxygen — abbreviated LOX , LOX or Lox in the aerospace , submarine and gas industries — is one of the physical forms of elemental oxygen .	
402: What does the term "mens rea" mean	Mens rea
Therefore, mens rea refers to the mental element of the offence that accompanies the actus reus.	
In Australia , for example, the elements of the federal offences are now designated as "fault elements" or "mental elements" (mens rea) and "physical elements" or "external elements" (actus reus).	
406: what do UAs detect	Drug test
Major uses of drug testing are to detect the presence of performance enhancing steroids in sport or for drugs prohibited by laws, such as cannabis , cocaine and heroin.	
414: What does Human sperm consist of?	Semen
The process that results in the discharge of semen is called ejaculation .	
416: how did the penguins acquire sidney crosby	Sidney Crosby
The Penguins returned to the Finals against Detroit the following year and won in seven games; Crosby became the youngest captain in NHL history to win the Stanley Cup .	
#420: how much total wealth in USA	Wealth in the United States
In addition, wealth is unevenly distributed, with the wealthiest 25% of US households owning 87% of the wealth in the United States , which was \$54.2 trillion in 2009.	
429: what does a plus-minus sign mean	Plus-minus sign
The sign is normally pronounced "plus or minus".	
433: how was the moon formed	Moon
Red and orange tinted Moon , as seen from Earth during a lunar eclipse , where the Earth comes between the Moon and Sun	
435: what did ronald reagan do as president	Ronald Reagan
Ronald Wilson Reagan (; February 6, 1911 – June 5, 2004) was the 40th President of the United States (1981–1989).	
439: what does judgment as a matter of law mean	Judgment as a matter of law
JMOL motions may also be made after the verdict is returned, where they are called "renewed" motions for judgment as a matter of law (RMOL), but the motion is still commonly known by its former name, judgment notwithstanding verdict , or n.o.v. (from the English judgment and the Latin non obstatre veredicto).	
#441: how many presidents have been assassinated	List of United States presidential assassination attempts and plots
Assassination attempts and plots on Presidents of the United States have been numerous: more than 20 attempts to kill sitting and former presidents , as well as the Presidents-elect , are known.	
443: what does oklahoma produce	Oklahoma
Oklahoma City and Tulsa serve as Oklahoma's primary economic anchors, with nearly two thirds of Oklahomans living within their metropolitan statistical areas .	
#444: how many redwall books are there	Redwall
It is the title of the first book of the series, published in 1986, as well as the name of the Abbey featured in the book and the name of an animated TV series based on three of the novels ( Redwall , Mattimeo , and Martin the Warrior ), which first aired in 1999.	
446: what date did the american civil war start	American Civil War
The American Civil War (ACW), also known as the War between the States or simply the Civil War (see naming) , was a civil war fought from 1861 to 1865 between the United States (the "Union" or the "North") and several Southern slave states that declared their secession and formed the Confederate States of America (the "Confederacy" or the "South").	
452: how jamesons irish whiskey is made	Jameson Irish Whiskey
With annual sales of over 31 million bottles, Jameson is by far the best selling Irish whiskey in the world, as it has been sold internationally since the early 19th century when John Jameson along with his son (also named John) was producing more than a million gallons annually.	
#464: How many consecutive games did Ken Jennings win?	Ken Jennings
Jennings is noted for holding the record for the longest winning streak on the U.S. syndicated game show Jeopardy! .	
468: what county is jennings, la	Jennings, Louisiana
Jennings is the principal city of the Jennings Micropolitan Statistical Area , which includes all of Jefferson Davis Parish.	
469: how does nanotechnology affect health	Nanotechnology
A more generalized description of nanotechnology was subsequently established by the National Nanotechnology Initiative , which defines nanotechnology as the manipulation of matter with at least one dimension sized from 1 to 100 nanometers .	
476: what does it take to start a lodge in freemason	Masonic Lodge
By exception the three surviving lodges that formed the world's first known Grand Lodge in London (today called the United Grand Lodge of England ) have the unique privilege to operate as time immemorial i.e. without such warrant; only one other lodge operates without a warrant - this is the Grand Stewards' Lodge in London, although it is not also entitled to the "time immemorial" title.	
480: how does black pepper grow	Black pepper
Peppercorns, and the ground pepper derived from them, may be described simply as pepper , or more precisely as black pepper (cooked and dried unripe fruit), green pepper (dried unripe fruit) and white pepper (dried ripe seeds).	
482: how south african leaders are elected	President of South Africa
The President of the Republic of South Africa is the head of state and head of government under South Africa's Constitution .	
485: how was the phone invented	Invention of the telephone
This article covers the early years from 1844 to 1898, from conception of the idea of an electric voice-transmission device, to failed attempts to use "make-and-break" current, to successful experiments with electromagnetic telephones by Alexander Graham Bell and Thomas Watson , and finally to commercially successful telephones in the late 19th century.	
491: what forms seasons	Season
Hot regions have two or three seasons; the rainy (or wet, or monsoon ) season and the dry season , and in some tropical areas, a cool or mild season .	
492: how did mohammed gandhi die	Mahatma Gandhi
Mohandas Karamchand Gandhi ( ; 2 October 1869 – 30 January 1948), commonly known as Mahatma Gandhi, was the preeminent leader of Indian nationalism in British-ruled India .	
#495: how many users do twitter have	Twitter
Twitter is an online social networking service and microblogging service that enables its users to send and read text-based messages of up to 140 characters , known as "tweets".	
513: how did the vietnam war end	Vietnam War
They viewed the conflict as a colonial war , fought initially against France, backed by the U.S., and later against South Vietnam , which it regarded as a U.S. puppet state .	
#516: how much more time does chemo give to people with renal cancer	Renal cell carcinoma
Renal cell carcinoma (RCC, also known as hypernephroma) is a kidney cancer that originates in the lining of the proximal convoluted tubule , the very small tubes in the kidney that transport GF (glomerular filtrate) from the glomerulus to the descending limb of the nephron.	
521: how does a solid state drive work	Solid-state drive
Hybrid drives or solid state hybrid drives (SSHD) combine the features of SSDs and HDDs in the same unit , containing a large hard disk drive and an SSD cache to improve performance of frequently accessed data.	
#529: how long did the roman empire last	Roman Empire
The Roman Empire () was the post-Republican period of the ancient Roman civilization , characterised by an autocratic form of government and large territorial holdings around the Mediterranean in Europe , Africa , and Asia .	
#533: how much does U.S. pay on health care per person	Comparison of the health care systems in Canada and the United States
Total government spending per capita in the U.S. on health care was 23% higher than Canadian government spending , and U.S. government expenditure on health care was just under 83% of total Canadian spending (public and private) though these statistics don't take into account population differences.	
537: how is rfid tagged power	Radio-frequency identification
Radio-frequency identification (RFID) is the wireless non-contact use of radio-frequency electromagnetic fields to transfer data, for the purposes of automatically identifying and tracking tags attached to objects.	
#556: how long have kanab ambersnail been endangered?	Kanab ambersnail
The Kanab ambersnail , scientific name Oxylama haydeni kanabensis or Oxylama kanabensis , is a critically endangered subspecies or species of small, air-breathing land snail , a terrestrial pulmonate gastropod mollusc in the family Succineidae , the amber snails.	
Now considered a Critically Endangered species on the IUCN Red List of Threatened Species due to a series of factors (including influence), the Kanab ambersnail has been reintroduced to three springs above the historic high water level along the Colorado River .	
#557: how many spanish speaking people are there in the us	Spanish language in the United States
There are 45 million Hispanophones who speak Spanish as a first or second language , as well as six million Spanish language students, composing the largest national Spanish-speaking community outside of Mexico and making Spanish the Romance language and the Indo-European language with largest number of native speakers in the world.	
558: what does (sic) mean?	Sic
Sic (noise artist) , styled as [sic], stage name of Jennifer Morris, a Canadian noise artist	
560: how does sedimentary rock form	Sedimentary rock
The sedimentary rock cover of the continents of the Earth's crust is extensive, but the total contribution of sedimentary rocks is estimated to be only 8% of the total volume of the crust.	
562: what did st patrick do	Saint Patrick
The text, however, distinguishes between "Old Patrick" (thought to mean Palladius ) and "Patrick , archapostle of the Scots , who died in 492.	
563: how do you know if something is the golden ratio	Golden ratio
Many 20th century artists and architects have proportioned their works to approximate the golden ratio—especially in the form of the golden rectangle , in which the ratio of the longer side to the shorter is the golden ratio—believing this proportion to be aesthetically pleasing (see Applications and observations below).	
570: How is a computer used?	Computer

Simple <b>computers</b> are small enough to fit into mobile devices , and mobile <b>computers</b> can be powered by small batteries . Personal <b>computers</b> in their various forms are icons of the Information Age and are what most people think of as "computers ."	
#581: how often does ham station need to ID? <i>Station identification</i>	Leap year
<b>Station</b> identification (ident or channel <b>ID</b> ) is the <b>practice</b> of radio or television <b>stations</b> or networks identifying themselves on air, typically by means of a call sign or brand name (sometimes known, particularly in the United States, as a "sounder" or "stinger", more generally as a <b>station</b> or network <b>ID</b> ).	The term <b>leap year</b> gets its name from the fact that while a fixed <b>date</b> in the Gregorian calendar normally advances one <b>day</b> of the <b>week</b> from one <b>year</b> to the next, in a <b>leap year</b> the <b>day</b> of the <b>week</b> will advance two <b>days</b> (from March onwards) due to the <b>year's extra day</b> inserted at the end of <b>February</b> (thus "leaping over" one of the <b>days</b> in the <b>week</b> ).
597: what county is bethlehem pa in <i>Bethlehem, Pennsylvania</i>	There are three general sections of the city: North <b>Bethlehem</b> , South <b>Bethlehem</b> and West <b>Bethlehem</b> .
598: what countries legalize marijuana <i>Legality of cannabis</i>	While <b>federal law</b> in the <b>United States</b> bans all sale and <b>possession</b> of <b>cannabis</b> , enforcement varies widely at the <b>state</b> level and some <b>states</b> have established medicinal <b>marijuana</b> programs in that contradict <b>federal law</b> ; two <b>states</b> (Colorado and <b>Washington</b> ) have repealed their <b>laws</b> prohibiting the recreational <b>use</b> of <b>cannabis</b> and replaced them with a regulatory regime, also contrary to <b>federal</b> statute.
#603: how many albums has dmx sold to this date <i>DMX (rapper)</i>	In 1999, <b>DMX</b> released his best-selling <b>album</b> ...And Then There Was <b>X</b> , which featured the hit <b>singel</b> " Party Up (Up in Here) "
610: what do mucous membranes secrete <i>Mucous membrane</i>	The term <b>mucous membrane</b> refers to where they are found in the body and not every <b>mucous membrane</b> secretes mucus.
611: what company is cricket wireless by <i>Cricket Wireless</i>	<b>Cricket Communications, Inc.</b> , ( d.b.a. <b>Cricket Wireless</b> ) founded in 1999, provides <b>wireless</b> services to over 7 million customers in the United States.
#615: how many numbers on a credit card <i>Bank card number</i>	<b>Payment card numbers</b> are found on <b>payment cards</b> , such as <b>credit cards</b> and <b>debit cards</b> , as well as stored-value <b>cards</b> , <b>gift cards</b> and other similar <b>cards</b> .
617: how is public policy created <i>Public policy</i>	The U.S. professional association of <b>public policy</b> practitioners, researchers, scholars, and students is the Association for <b>Public Policy</b> Analysis and Management .
	As an academic discipline , <b>public policy</b> is studied by professors and students at <b>public policy</b> schools of major universities throughout the country.
#619: how many wives did henry the 8th have <i>Henry VIII of England</i>	<b>Henry's</b> struggles with Rome led to the separation of the <b>Church of England</b> from papal authority, the Dissolution of the Monasteries , and his own establishment as the Supreme Head of the <b>Church of England</b> .
620: what country are bongo drums from? <i>Bongo drum</i>	The <b>drums</b> are of different size: the larger <b>drum</b> is called in Spanish the <b>hembra</b> ( female ) and the smaller the <b>macho</b> ( male ).
#621: how many books in bible <i>Books of the Bible</i>	The first part of Christian <b>Bibles</b> is the Old Testament , which contains, at minimum, the twenty-four <b>books</b> of the Hebrew <b>Bible</b> divided into thirty-nine <b>books</b> and <b>ordered</b> differently from the Hebrew <b>Bible</b> .
639: what does freedom of speech cover <i>Freedom of speech</i>	The <b>right</b> to <b>freedom</b> of expression is recognized as a <b>human right</b> under <b>Article 19</b> of the Universal Declaration of <b>Human Rights</b> and recognized in international <b>human rights</b> law in the International Covenant on Civil and Political <b>Rights</b> (ICCPR).
647: what county is san jose in? <i>San Jose, California</i>	<b>San Jose</b> was founded on November 29, 1777, as El Pueblo de <b>San José</b> de Guadalupe, the <b>first</b> civilian town in the Spanish colony of Nueva <b>California</b> , which later became Alta <b>California</b> .
651: what cars have smart key systems <i>Smart key</i>	A <b>smart key</b> is an electronic <b>access</b> and authorization <b>system</b> which is available as an option or standard in <b>several cars</b> .
676: how did John F. Kennedy die? <i>John F. Kennedy</i>	<b>John Fitzgerald "Jack" Kennedy</b> (May 29, 1917 – November 22, 1963), often referred to by his initials JFK, was the 35th President of the United States , serving from 1961 until his <b>death</b> in 1963.
#677: how many members are in the house of representatives <i>United States House of Representatives</i>	The Speaker of the <b>United States House of Representatives</b> , who <b>presides</b> over the chamber, is elected by the <b>members</b> of the <b>House</b> , and is therefore traditionally the <b>leader</b> of the <b>House Democratic Caucus</b> or the <b>House Republican Conference</b> , whichever party has more voting <b>members</b> .
678: what can silk be used for <i>Silk</i>	Many <b>silks</b> are mainly <b>produced</b> by the <b>larvae of insects</b> undergoing complete metamorphosis , but some adult <b>insects</b> such as webspinners <b>produce silk</b> , and some <b>insects</b> such as raspberry crickets <b>produce silk</b> throughout their lives.
682: what does uncle sam represent to the american people <i>Uncle Sam</i>	The first use of <b>Uncle Sam</b> in literature was in the 1816 allegorical book "The Adventures of Uncle Sam in Search After His Lost Honor" by Frederick Augustus Fiddday, Esq. An <b>Uncle Sam</b> is mentioned as early as 1775, in the original "Yankee Doodle" lyrics of the Revolutionary War.
684: what do cyberstalkers do <i>Cyberstalking</i>	<b>Cyberstalking</b> is a criminal offense that comes into play under state anti-stalking <b>laws</b> , slander <b>laws</b> , and harassment <b>laws</b> .
685: how do insulin syringes work <i>Syringe</i>	The word " <b>syringe</b> " is derived from the Greek <b>syringē</b> syrinx = "tube" via back-formation of a new singular from its Greek-type plural " <b>syringes</b> " ( <b>syngenes</b> syringes).
692: how did david carradine die <i>David Carradine</i>	<b>David Carradine</b> (born John Arthur Carradine; December 8, 1936 – June 3, 2009) was an American actor and martial artist, best known for his leading role as a warrior monk , Kwai Chang Caine , in the 1970s television series Kung Fu .
700: what affects the money supply <i>Money supply</i>	Second, if the velocity of <b>money</b> , i.e., the ratio between <b>nominal GDP</b> and <b>money supply</b> , changes, an <b>increase</b> in the <b>money supply</b> could have either no <b>effect</b> , an <b>exaggerated effect</b> , or an <b>unpredictable effect</b> on the growth of <b>nominal GDP</b> .
703: how is today special? <i>Today's Special</i>	It was <b>set</b> in a department store , based on the flagship location of the now defunct Simpson's in <b>Toronto</b> .
707: what does auld lang syne mean <i>Auld Lang Syne</i>	Consequently, "For <b>auld lang syne</b> " , as it appears in the first line of the chorus, might be loosely translated as "for (the sake of) old times."
715: how was the president involved in the gulf war <i>Gulf War</i>	The <b>war</b> is also known under other <b>names</b> , such as the Persian <b>Gulf War</b> , First <b>Gulf War</b> , <b>Gulf War I</b> , or the First <b>Iraq War</b> , before the term " <b>Iraq War</b> " became identified instead with the 2003 <b>Iraq War</b> (also referred to in the U.S. as " <b>Operation Iraqi Freedom</b> ").
#716: how many times has a player hit for the cycle <i>Hitting for the cycle</i>	One <b>NPB</b> player has also <b>hit</b> for the <b>cycle</b> in an <b>NPB All-Star game</b> .
723: what does it mean if i'm flat footed? <i>Flat feet</i>	Three studies (see citations below in military section) of military recruits have shown no evidence of later increased <b>injury</b> , or <b>foot</b> problems, due to <b>flat feet</b> , in a population of people who reach military service age without prior <b>foot</b> problems.
726: what age group is generation x <i>Generation X</i>	<b>Generation X</b> , commonly abbreviated to Gen <b>X</b> , is the <b>generation</b> born after the Western post-World War II baby boom .
729: how did James Dean die? <i>James Dean</i>	<b>James Byron Dean</b> (February 8, 1931 – September 30, 1955) was an American actor.
734: what glows in the dark <i>Glow-in-the-dark</i>	Glow in the <b>Dark</b> Tour , a 2008 concert tour by Kanye West
737: what did sparta do around 650 bc <i>Sparta</i>	<b>Sparta</b> ( Doric <b>Greek</b> : , Attic <b>Greek</b> : ), or Lacedaemon, was a <b>prominent</b> city-state in <b>ancient Greece</b> , situated on the banks of the Eurotas River in Laconia, in south-eastern Peloponnese .
738: what fantasy american football means <i>Fantasy football (American)</i>	<b>Fantasy football</b> has vastly increased in <b>popularity</b> , particularly because <b>Fantasy football</b> providers such as ESPN, Yahoo, CBS, and the NFL itself are able to keep track of statistics entirely online, eliminating the need to check box scores and newspapers regularly to keep track of players.
#741: how many muscles in the body <i>List of muscles of the human body</i>	The <b>muscles</b> of the human <b>body</b> can be categorized into a number of groups which include <b>muscles</b> relating to the head and neck, <b>muscles</b> of the torso or trunk, <b>muscles</b> of the upper limbs, and <b>muscles</b> of the lower limbs.
747: what does the green mean on the mexican flag <i>Flag of Mexico</i>	While the <b>meaning</b> of the <b>colors</b> has changed over time, these three <b>colors</b> were adopted by <b>Mexico</b> following <b>independence</b> from Spain during the country's War of <b>Independence</b> , and subsequent First <b>Mexican Empire</b> .
#749: how often do elk have sex <i>Elk</i>	Male <b>elk</b> have large antlers which are shed each year.
#757: how many players on a side for a football game <i>American football</i>	American football is the most popular sport in the <b>United States</b> today, and the <b>National Football League</b> (NFL) is its most popular league.
760: what does the temporal lobe part of the brain do <i>Temporal lobe</i>	The <b>temporal lobe</b> is a region of the cerebral cortex that is located beneath the lateral fissure on both cerebral hemispheres of the mammalian <b>brain</b> .
#761: how many lungs does a human have <i>Human lung</i>	The right <b>lung</b> consists of three <b>lobes</b> while the left <b>lung</b> is slightly smaller consisting of only two <b>lobes</b> (the left <b>lung</b> has a "cardiac notch" allowing space for the heart within the chest).
768: what is the population of center tx <i>Center, Texas</i>	The Rio Theater in <b>Center</b>
	It is named for its location near the <b>center</b> of Shelby County, not for its location in Texas as it is located near the Louisiana border.
	First Christian Church at 124 Cora Street in <b>Center</b> is one of the oldest congregations in the community.
	First Baptist Church at 117 Cora Street in <b>Center</b> is located next to the downtown section.
	Chamber of Commerce Building in <b>Center</b> .
	<b>Center</b> is a city in Shelby County , Texas , United States .
770: what causes photo red eye <i>Red-eye effect</i>	The red-eye effect in photography is the common appearance of red pupils in color photographs of eyes .
#771: what creates a cloud <i>Cloud</i>	These include strato- for low <b>clouds</b> with limited convection that form mostly in uneven layers, cumulo- for complex highly-convective storm <b>clouds</b> , nimbo- for thick layered <b>clouds</b> of some complexity that can produce moderate to heavy precipitation, alto- for middle <b>clouds</b> , and cirro- for high <b>clouds</b> ; the latter two of which may be of simple or moderately complex structure.
#772: how many students go to santa barbara <i>University of California, Santa Barbara</i>	The University of California, Santa Barbara (commonly referred to as UC Santa Barbara or UCSB) is a public research university and one of the ten general campuses of the University of California system.
774: what culture is mariah carey <i>Mariah Carey</i>	Under the guidance of Columbia Records executive Tommy Mottola , <b>Carey</b> released her self-titled debut studio album <b>Mariah Carey</b> in 1990; it went multi-platinum and spawned four consecutive number one singles on the U.S. Billboard Hot 100 chart.
779: how was color introduced in film? <i>Color motion picture film</i>	Color motion picture film refers both to unexposed <b>color</b> photographic film in a format suitable for use in a <b>motion picture camera</b> , and to finished <b>motion picture film</b> , ready for use in a <b>projector</b> , which bears images in <b>color</b> .
781: what artist has song with ashanti? <i>Ashanti (entertainer)</i>	<b>Ashanti</b> is most famous for her eponymous debut album , which <b>featured</b> the hit song " Foolish " , and sold over 503,000 copies in its first week of release throughout the U.S. in April 2002.
783: what came first army or air force <i>United States Army Air Forces</i>	Although other <b>nations</b> already had separate <b>air forces independent</b> of the <b>army</b> or <b>navy</b> (such as the British Royal Air Force and the German Luftwaffe ), the AAF remained a part of the <b>United States Army</b> until the <b>United States Air Force</b> came into being in September 1947.
784: what day is 2011 super bowl? <i>Super Bowl XLV</i>	Unlike other matchups, this game featured two title-abundant franchises: coming into the game, the <b>Packers</b> held the most NFL championships with 12 (9 league championships prior to the Super Bowl era and 3 Super Bowl championships), while the Steelers held the most Super Bowl championships with 6.
785: what does the FOIA apply to <i>Freedom of Information Act (United States)</i>	The Federal Government's Freedom of Information Act should not be confused with the different and varying Freedom of Information Acts passed by the individual states . Many of those state acts may be similar but not identical to the federal act.
#786: how many countries are member of the eu? <i>European Union</i>	With a combined <b>population</b> of over 500 million inhabitants, or 7.3% of the world <b>population</b> , the EU in 2012, generated a nominal gross domestic product (GDP) of 16.584 trillion US dollars, representing approximately 20% of the global GDP when measured in terms of purchasing power parity , and represents the <b>largest</b> nominal GDP and GDP PPP in the world.
#790: how many vehicles are registered in the us <i>Passenger vehicles in the United States</i>	This number, along with the average age of <b>vehicles</b> , has increased steadily since 1960, indicating a growing number of <b>vehicles</b> per capita.
792: what cards do you need in poker to get a royal flush <i>List of poker hands</i>	The ranking of a particular hand is increased by including multiple <b>cards</b> of the same <b>card</b> rank, by all five <b>cards</b> being from the same suit , or by all five <b>cards</b> being of consecutive rank.
793: what does xylem transport <i>Xylem</i>	The word <b>xylem</b> is derived from the Greek word ξύλον (xylon), meaning "wood"; the best-known <b>xylem</b> tissue is <b>wood</b> , though it is <b>found</b> throughout the plant.
#800: how many myps has kobe bryant won <i>Kobe Bryant</i>	In 2006, <b>Bryant</b> scored a career-high 81 points against the Toronto Raptors , the second most points scored in a single <b>game</b> in NBA history, second only to Wilt Chamberlain's 100-point <b>game</b> in 1962.
804: what do jehovah witnesses believe <i>Jehovah's Witnesses</i>	<b>Jehovah's Witnesses</b> is a millenarian restorationist Christian denomination with nontrinitarian beliefs distinct from mainstream Christianity. Sources for descriptors: • Millenarian: • Restorationist: • Christian: • Denomination: The organization reports worldwide membership of over 7.78 million adherents involved in evangelism , convention attendance of over 12 million, and annual Memorial attendance of over 19 million.
#806: how old old is xp operating system <i>Windows XP</i>	According to web analytics data generated by Net Applications , Windows XP was the most widely used operating system until August 2012, when Windows 7 overtook it.
#811: how many pawns in chess <i>Pawn (chess)</i>	It is also common to refer to a rook <b>pawn</b> , meaning any <b>pawn</b> on the a- or h-file, a knight <b>pawn</b> (on the b- or g-file), a bishop <b>pawn</b> (on the c- or f-file), a queen <b>pawn</b> (on the d-file), a king <b>pawn</b> (on the e-file), and a central <b>pawn</b> (on either the d- or e-file).
814: how is hydrogen produced <i>Hydrogen production</i>	Hydrogen production is the family of industrial methods for generating <b>hydrogen</b> .
833: what country is madrid spain in <i>Madrid</i>	As the capital city of <b>Spain</b> , seat of government , and residence of the Spanish monarch , <b>Madrid</b> is also the political, economic and cultural centre of <b>Spain</b> .
835: what causes rogue waves <i>Rogue wave</i>	Rogue waves (also known as freak waves , monster waves , killer waves , extreme waves , and abnormal waves) are relatively large and spontaneous ocean surface waves that occur far out at sea, and are a threat even to large ships and ocean liners .
#837: what happened to montgomery cliff <i>Montgomery Clift</i>	Edward <b>Montgomery Clift</b> (October 17, 1920July 23, 1966) was an American film and stage actor The New York Times' obituary noted his portrayal of "moody, sensitive young men".
#850: what happened to george o'malley on grey's anatomy? <i>George O'Malley</i>	Introduced as a surgical intern at the fictional Seattle Grace Hospital, <b>O'Malley</b> worked his way up to resident level, while his relationships with his colleagues Meredith <b>Grey</b> ( Ellen Pompeo ), Cristina Yang ( Sandra Oh ), Izzie <b>Stevens</b> ( Katherine Heigl ) and Alex Karev ( Justin Chambers ) formed a focal point of the series.
#852: how many people die from myasthenia gravis per year <i>Myasthenia Gravis</i>	<b>Myasthenia gravis</b> (from Greek μυϊκός "muscle", "weakness", and "serious"; abbreviated MG) is an autoimmune neuromuscular disease leading to fluctuating muscle weakness and fatigability .
#854: how many rooms in Borgata hotel <i>The Borgata</i>	The <b>Borgata Hotel Casino</b> is a luxury <b>hotel</b> , casino , and <b>spa</b> in Atlantic City, New Jersey , United States .
#859: what country has the most muslims in the world <i>Islam by country</i>	

In the Middle East, the non-Arab <b>countries</b> of Turkey and Iran are the largest <b>Muslim-majority countries</b> ; in Africa, Egypt and Nigeria have the most populous <b>Muslim</b> communities.	<i>Notary public</i>
861: what country is dubai in	Dubai
The Sheikhdom of <b>Dubai</b> was formally established in 1833 by Sheikh Maktoum bin Butti Al-Maktoum when he persuaded around 800 members of his tribe of the Bani Yas , living in what was then the Second Saudi State and now part of Saudi Arabia , to follow him to the <b>Dubai Creek</b> by the Abu Falasa clan of the Bani Yas.	
#863: how much of earth is covered ocean water	Ocean
The Mars <b>ocean</b> hypothesis <b>suggests</b> that nearly a third of the surface of Mars was once <b>covered</b> by <b>water</b> , though the <b>water</b> on Mars is no longer <b>oceanic</b> , and a runaway greenhouse effect may have boiled away the global <b>ocean</b> of Venus.	
#870: how much did yankee stadium cost	Yankee Stadium
The <b>first</b> game at the new <b>Yankee Stadium</b> was a pre-season exhibition game against the Chicago Cubs <b>played</b> on <b>April</b> 3, 2009, which the <b>Yankees</b> <b>won</b> 7-4.	
#879: what happened to "The Glades" tv series	<i>The Glades (TV series)</i>
The <b>Glades</b> was renewed by A&E for a third <b>season</b> on October 18, 2011, which aired from June 3 to August 12, 2012.	
#887: how many albums has eminem sold in his career	Eminem
His next two Records <b>Marshall Mathers LP</b> , and <b>The Eminem Show</b> , also <b>won</b> Best Rap <b>Album Grammy Awards</b> , making <b>Eminem</b> the first <b>artist</b> to <b>win</b> Best Rap <b>Album</b> for three consecutive LPs.	
891: what do pigs eat	Pig
Pigs include the domestic <b>pig</b> , its ancestor the wild boar , and several other wild relatives.	
#900: what happened during the Starving Time in Jamestown?	Starving Time
The <b>Starving Time</b> at <b>Jamestown</b> in the Colony of Virginia was a period of <b>starvation</b> during the winter of 1609–1610 in which all but 60 of 500 colonists died.	
905: what does the president of the usa do	President of the United States
It also prohibits a <b>person</b> from being elected to the <b>presidency</b> more than once if that <b>person</b> previously had <b>served</b> as <b>president</b> , or acting <b>president</b> , for more than two years of another <b>person's</b> term as <b>president</b> .	
911: how kimberlite pipes form	Volcanic pipe
Volcanic <b>pipes</b> are geological structures <b>formed</b> by the violent, supersonic eruption of deep-origin volcanoes .	
#913: how many seasons of grey's anatomy are there	<i>Grey's Anatomy</i>
Dr. Preston Burke (Isaiah Washington) <b>departs</b> at the conclusion of the third <b>season</b> , and is <b>replaced</b> by Dr. Erica Hahn (Brooke Smith ), who leaves the show during the fifth <b>season</b> , and later Dr. Teddy Altman (Kim Raver ), who <b>departs</b> at the end of the eighth <b>season</b>	
914: how is schizophrenia diagnosed?	Schizophrenia
Despite the etymology of the term from the Greek roots <i>skhizein</i> ("to split") and <i>phren-</i> (phren- (φρήν, φρενός; "mind"), <b>schizophrenia</b> does not imply a "split personality", or "multiple personality disorder" (which is known these days as dissociative identity disorder)—a condition with which it is often confused in public perception.	
931: what kind of company is Microsoft?	Microsoft
As of 2013, <b>Microsoft</b> is market dominant in both the PC operating <b>system</b> and office suite markets (the latter with <b>Microsoft Office</b> ).	
938: how does weather happen	Weather
Studying how the <b>weather</b> <b>works</b> on other planets has been helpful in understanding how <b>weather</b> <b>works</b> on Earth.	
#940: how many countries have english as an official language	<i>List of countries where English is an official language</i>
According to the Constitution of India , "Hindi in the Devanagari script" is the <b>official language</b> of the union; and <b>English</b> the 'subsidiary <b>official language</b> '; however, <b>English</b> is mandated for the authoritative texts of all federal laws and Supreme Court decisions and (along with Hindi) is one of the two <b>languages</b> of the Indian Parliament .	
941: what countries are in cono sur	<i>Southern Cone</i>
<b>High life expectancy</b> , the <b>highest Human Development Index</b> of Latin America, <b>high Standard</b> of living , significant participation in the global markets and the emerging economy of its members make the <b>Southern Cone</b> the most prosperous macro-region in <b>South America</b> .	
949: what branch of the military is delta force	<i>Delta Force</i>
The Central Intelligence Agency's highly secretive Special Activities Division (SAD) and more <b>specifically</b> its elite Special Operations Group (SOG) often works with – and recruits – <b>operators</b> from <b>Delta Force</b> .	
957: what causes the seasons	<i>Season</i>
Hot regions have two or three <b>seasons</b> ; the <b>rainy</b> (or wet, or monsoon ) <b>season</b> and the <b>dry season</b> , and in some <b>tropical</b> areas, a cool or mild <b>season</b> .	
959: what does 3g network mean	3G
3G finds application in <b>wireless</b> voice <b>telephony</b> , mobile Internet access, fixed <b>wireless</b> Internet access, video calls and mobile TV .	
960: what does Mazel tov!	Mazel tov
"Mazel tov" or "mazel tov" (Hebrew / Yiddish : מַלְאָכִיכְךָ , בָּרוּךְ Hebrew: "mazal tov"; Yiddish: "mazel tov"; lit.	
961: How Do You Get Hepatitis C	Hepatitis C
Hepatitis C is an infectious disease affecting primarily the liver , caused by the <b>hepatitis C</b> virus (HCV).	
962: how does flexible spending account work	Flexible spending account
The most common <b>type</b> of <b>flexible spending account</b> , the medical expense FSA (also medical FSA or <b>health</b> FSA), is similar to a <b>health</b> savings <b>account</b> (HSA) or a <b>health</b> reimbursement <b>account</b> (HRA).	
#966: how many baseball teams usa	<i>Major league baseball</i>
Major <b>League Baseball</b> (MLB) is a professional <b>baseball league</b> , consisting of <b>teams</b> that play in the <b>National League</b> and the <b>American League</b> .	
#972: how many amendments in the US constitution	<i>List of amendments to the United States Constitution</i>
To date, no convention for proposing <b>amendments</b> has been called by the states, and only once—in 1933 for the ratification of the twenty-first <b>amendment</b> —has the convention method of ratification been employed.	
979: what day is the federal holiday for Martin Luther King Jr.	<i>Martin Luther King, Jr. Day</i>
Martin Luther King, Jr. Day is a United States <b>federal holiday</b> marking the birthday of Rev. Dr. Martin Luther King, Jr.	
985: what area code is 949	<i>Area code 949</i>
The <b>area code</b> in red is <b>Area Code 949</b> ; all others in blue are California <b>area codes</b> .	
#994: how many babies are in a typical raccoon litter	Raccoon
The <b>raccoon</b> (, Procyon lotor), sometimes spelled <b>racoons</b> , also known as the common <b>raccoon</b> , North <b>American</b> <b>raccoon</b> , northern <b>raccoon</b> and colloquially as coon , is a medium-sized mammal <b>native</b> to North <b>America</b> .	
998: What does Rapture meaning in a theological sense?	Rapture
Denominations such as Roman Catholics , Orthodox Christians , Lutheran Christians , and Reformed Christians believe in a <b>rapture</b> only in the <b>sense</b> of a general final resurrection , when Christ returns a single time.	
1002: what are corporation balance	<i>Balance sheet</i>
In financial accounting , a <b>balance</b> sheet or statement of financial position is a summary of the financial <b>balances</b> of a sole proprietorship , a business partnership , a <b>corporation</b> or other business organization , such as an LLC or an LLP .	
1003: what is .17 hmr caliber	.17 HMR
It descended from the .22 Magnum by necking down the .22 Magnum case to take a .17 <b>caliber</b> (4.5 mm) bullet, and it is more costly to shoot than traditional .22 <b>caliber</b> rimfire cartridges.	
1007: what made the civil war different from others	<i>American Civil War</i>
The American Civil War (ACW), also known as the <b>War</b> between the States or simply the <b>Civil War</b> (see naming ), was a <b>civil war</b> fought from 1861 to 1865 between the United States (the "Union" or the "North") and several <b>Southern slave</b> states that declared their secession and formed the Confederate States of America (the "Confederacy" or the "South").	
1008: where to buy potato bread made without wheat	<i>Potato bread</i>
Potato <b>bread</b> is a form of <b>bread</b> in which <b>potato</b> replaces a portion of the regular <b>wheat</b> flour .	
1014: who killed julius caesar	<i>Julius caesar</i>
Gaius Julius Caesar (, July 100 BC – March 44 BC) was a Roman general , statesman , Consul and notable author of Latin prose.	
1024: who wrote white christmas	<i>White Christmas (song)</i>
I just <b>wrote</b> the best <b>song</b> I've ever <b>written</b> — heck, I just <b>wrote</b> the best <b>song</b> that anybody's ever <b>written</b> !"	
1027: WHAT IS A FY QUARTER	Fiscal year
Many universities have a <b>fiscal year</b> which ends during the summer, both to align the <b>fiscal year</b> with the school <b>year</b> (and, in some cases involving public universities, with the state government's <b>fiscal year</b> ), and because the school is normally less <b>busy</b> during the summer months.	
1028: who pulmonary hypertension	<i>Pulmonary hypertension</i>
In medicine , <b>pulmonary hypertension</b> (PH) is an increase of <b>blood pressure</b> in the <b>pulmonary</b> artery , <b>pulmonary</b> vein , or <b>pulmonary</b> capillaries, together known as the <b>lung</b> <b>vasculature</b> , leading to shortness of breath , dizziness , fainting , and other symptoms, all of which are exacerbated by exertion.	
1033: what is a notary for	<i>Notary public</i>
With the exceptions of Louisiana , Puerto Rico , Quebec , whose private law is based on civil law , and British Columbia , whose notarial tradition stems from scrivener <b>notary</b> practice, a <b>notary</b> public in the rest of the <b>United States</b> and most of Canada has powers that are far more limited than those of civil-law or other <b>common-law notaries</b> , both of whom are qualified lawyers admitted to the bar: such <b>notaries</b> may be referred to as notaries-at-law or lawyer <b>notaries</b> .	
#1038: when did the cold war start	<i>Cold War</i>
The tensest times were during the Berlin Blockade (1948–1949), the <b>Korean War</b> (1950–1953), the Suez Crisis (1956), the Berlin Crisis of 1961 , the Cuban missile crisis (1962), the Vietnam War (1959–1975), the Yom Kippur War (1973), the Soviet war in Afghanistan (1979–1989), the Soviet downing of Korean Air Lines Flight 007 (1983), and the "Able Archer" NATO military exercises (1983).	
1039: what state is pine's peak in?	<i>Pikes Peak</i>
It gets its name from the Iowa incarnation of <b>Pikes Peak</b> , a particularly high point overlooking the gorge of the Upper Mississippi, and like <b>Pikes Peak</b> in Colorado , is named for Zebulon <b>Pike</b> .	
1049: what is a book index	<i>Index (publishing)</i>
In a traditional back-of-the-book <b>index</b> the headings will include names of people, places and events, and concepts selected by a person as being relevant and of interest to a possible reader of the <b>book</b> .	
1050: who sings the song never ending story	<i>The NeverEnding Story (song)</i>
"The NeverEnding Story" (titled "The NeverEnding Story (L'histoire sans fin)" in the French version) is the title <b>song</b> from the English version of the 1984 film <i>The NeverEnding Story</i> .	
1053: who created facebook	<i>Facebook</i>
Facebook is an online social networking service , whose <b>name</b> stems from the colloquial <b>name</b> for the book given to students at the start of the academic year by some <b>university</b> administrations in the United States to help <b>students</b> get to know each other.	
1057: what are risk for infections	<i>Risk of infection</i>
Although anyone can become <b>infected</b> by a pathogen, patients with this diagnosis are at an elevated <b>risk</b> and extra <b>infection</b> controls should be considered.	
1061: where do cruises dock in new york city	<i>New York Passenger Ship Terminal</i>
With an upsurge in <b>cruise</b> ship traffic and the terminal's ability to handle comfortably only three large ships at a time, two <b>new</b> terminals have opened in the harbor — the Cape Liberty <b>Cruise</b> Port opened in 2004 in Bayonne, New Jersey (used by Royal Caribbean <b>Cruise</b> Line , Celebrity <b>Cruises</b> and Azamara <b>Cruises</b> ), and the Brooklyn <b>Cruise</b> Terminal (used by the Queen Mary 2 and other ships of the Carnival Corporation <b>cruise</b> brands) opened in 2006 in Brooklyn, New York .	
1064: who played the lead roles in the movie leaving las vegas	<i>Leaving Las Vegas</i>
After limited <b>release</b> in the United States on October 27, 1995, <i>Leaving Las Vegas</i> made its nationwide <b>release</b> on February 9, 1996, receiving strong praise from critics and audiences.	
1065: what is a CMM machine	<i>Coordinate-measuring machine</i>
A <b>machine</b> which takes readings in six degrees of freedom and displays these readings in mathematical form is known as a <b>CMM</b> .	
#1067: what percentage of water in in the body	<i>Body water</i>
In physiology , <b>body water</b> is the <b>water</b> content of the human <b>body</b> .	
1068: what type of batteries are 357 (LR44)	<i>LR44 battery</i>
Silver-oxide <b>batteries</b> type SR44 may provide extra capacity compared to <b>LR44 types</b> but have slightly different voltage characteristics.	
1069: where did hurricane katrina begin	<i>Hurricane Katrina</i>
The <b>hurricane</b> <b>strengthened</b> to a Category 5 <b>hurricane</b> over the warm Gulf water, but weakened before making its second landfall as a Category 3 <b>hurricane</b> on the morning of Monday, August 29 in southeast Louisiana .	
1073: what is 1 mil guarnaries in united states dollars	<i>Mil (currency)</i>
In the <b>United States</b> , it is a notional unit equivalent to a <b>United States dollar</b> (a tenth of a cent).	
1075: where did the persian war take place	<i>Greco-Persian Wars</i>
The allied <b>Greeks</b> followed up their success by destroying the rest of the <b>Persian fleet</b> at the <b>Battle</b> of Mycale , before expelling <b>Persian garrisons</b> from Sestos (479 BC) and Byzantium (478 BC).	
1078: what is a day care for?	<i>Day care</i>
Some childminders <b>care</b> for children from several <b>families</b> at the same time, either in their own <b>home</b> (commonly known as "family day care" in Australia) or in a <b>specialized child care</b> facility.	
1079: who discovered neptune the planet	<i>Discovery of Neptune</i>
Unfortunately, Le Verrier's triumph also led to a tense international dispute over priority, as, shortly after the <b>discovery</b> , George Airy , at the time British Astronomer <b>Royal</b> , announced that Adams had also predicted the <b>discovery</b> of the <b>planet</b> .	
1081: what separates me from your album	<i>What Separates Me from You</i>
The <b>album</b> debuted on the US Billboard 200 at <b>number</b> 11 with 58,000 first week sales, becoming A Day to Remember's personal best, as Homesick peaked at <b>number</b> 21.	
#1084: when did thomson make the plum-pudding model	<i>Plum pudding model</i>
Still, <b>Thomson's model</b> (along with a similar Saturnian ring <b>model</b> for atomic electrons , also put forward in 1904 by Naogaoka after James Clark Maxwell's <b>model</b> of Saturn's rings ), were earlier harbingers of the later and more successful solar-system-like Bohr <b>model</b> of the atom.	
1085: what part of the pig is bacon	<i>Bacon</i>
Meat from other animals, such as beef , lamb , chicken , goat , or turkey , may also be cut, cured, or otherwise prepared to resemble <b>bacon</b> , and may even be <b>referred</b> to as "bacon".	
1086: where did erisa come from	<i>Employee Retirement Income Security Act</i>
ERISA is sometimes used to refer to the full body of <b>laws</b> regulating employee benefit plans, which are found mainly in the Internal Revenue Code and <b>ERISA</b> itself.	
1087: what percent of illegal immigrants are from mexico and europe	<i>Illegal immigration to the United States</i>
Illegal immigration to the United States is the act of foreign <b>nationals</b> entering the United States , without <b>government</b> permission and in violation of <b>United</b> States nationality law , or staying beyond the termination date of a visa, also in violation of the law.	
#1092: when us subprime mortgage market collapse	<i>Subprime mortgage crisis</i>
As <b>adjustable-rate mortgages</b> began to reset at <b>higher</b> <b>interest</b> <b>rates</b> ( <b>causing</b> <b>higher</b> <b>monthly</b> <b>payments</b> ), <b>mortgage delinquencies</b> soared.	
1098: what type of game is heavy rain	<i>Heavy Rain</i>
It won 2010's Game of the Year from CNN and <b>Gaming Union</b> , and Best PS3 Game of 2010 by GameSpy and IGN .	
1099: what are superannuation contributions?	<i>Superannuation in Australia</i>
For example, <b>employers</b> are required to pay a proportion of an <b>employee</b> 's salaries and wages (currently 9%) into a <b>superannuation</b> fund, but people are encouraged to put aside <b>additional</b> funds into <b>superannuation</b> .	
1102: what is a group of deer called	<i>Deer</i>
Species in the Cervidae family include white-tailed <b>deer</b> , mule <b>deer</b> such as black-tailed <b>deer</b> , <b>elk</b> , moose , red <b>deer</b> , reindeer (caribou) , fallow <b>deer</b> , roe <b>deer</b> and chital .	
1103: what state was john mccain a senator in during the 2008 election	<i>United States presidential election, 2008</i>
This <b>election</b> was also notable for being the first time in U.S. history that both <b>major party</b> candidates were sitting U.S. <b>Senators</b> , only the third time (after 1920 and 1960 ) that any sitting U.S. <b>Senator</b> was <b>elected</b> <b>president</b> , and only the second time that the <b>winning President</b> and <b>Vice President</b> (Obama and Biden) were both sitting U.S. <b>Senators</b> .	
#1104: when does the electoral college votes	<i>Electoral College (United States)</i>
Maine and Nebraska use the "congressional district method", selecting one <b>elector</b> within each <b>congressional</b> district by <b>popular vote</b> and selecting the remaining <b>two electors</b> by a statewide <b>popular vote</b> .	
#1105: when barack obama was born	<i>Barack Obama</i>
He began his <b>presidential campaign</b> in 2007 , and in 2008 , after a close <b>primary campaign</b> against Hillary Rodham Clinton , he won sufficient delegates in the Democratic Party <b>primaries</b> to receive the <b>presidential nomination</b> .	
#1128: what year did disney's animal kingdom lodge open	<i>Disney's Animal Kingdom Lodge</i>
Disney's <b>Animal Kingdom</b> Lodge is located in the <b>Animal Kingdom</b> Resort Area , adjacent to Disney's <b>Animal Kingdom</b> .	
1130: what are two languages in nigeria?	<i>Languages of Nigeria</i>
Nigeria's linguistic diversity is a microcosm of Africa as a whole, encompassing three major African <b>languages</b> families : Afroasiatic , Nilo-Saharan , and Niger-Congo .	
1133: what war led to Pearl Harbor	<i>Events leading to the attack on Pearl Harbor</i>
Rather than seize and fortify the islands, and wait for the inevitable US counterattack, Japan's military <b>leaders</b> instead decided on the pre-emptive <b>Pearl Harbor</b> attack, which they assumed would negate the American forces needed for the liberation and reconquest of the islands.	
A series of events <b>led</b> to the attack on <b>Pearl Harbor</b> .	
1141: where in the world are smallpox common	<i>Smallpox</i>
The disease was originally known in English as the "pox" or "red plague"; the term " <b>smallpox</b> " was first used in <b>Britain</b> in the 15th century to distinguish <b>variola</b> from the "great pox" ( syphilis ).	

#1142: what year did the beatles came out with the song i wanna hold your hand "I Want to Hold Your Hand" is a song by the English rock band the <b>Beatles</b> .	<i>I Want to Hold Your Hand</i>	
1143: what is a gasser car <b>The gasser</b> is the predecessor of the modern <b>Funny Car</b> .	<i>Gasser (car)</i>	
1145: what tensions preceded the berlin blockade During the multinational <b>occupation</b> of post-World War II <b>Germany</b> , the Soviet Union <b>blocked</b> the Western Allies' railway, road, and canal access to the sectors of <b>Berlin</b> under <b>Allied control</b> .	<i>Berlin Blockade</i>	<i>Role-playing game</i>
1153: who sang that song-a-change is going to come The <b>song</b> has gained in popularity and critical acclaim in the decades since its release, and is #12 on Rolling Stone's 500 Greatest <b>Songs</b> of All Time .	<i>A Change Is Gonna Come</i>	<i>Forward (association football)</i>
#1155: what percentage of the human body is water In physiology, <b>body water</b> is the <b>water</b> content of the <b>human body</b> .	<i>Body water</i>	
1158: what is a form of legal ownership Given a short-sighted <b>owner</b> , however, a <b>private</b> property system can make these <b>tragedies</b> worse—for example, a <b>private owner</b> of a piece of oil-rich property, depending on his worldview, might be more interested in short-term financial gain than incremental use with an eye toward other's concerns (e.g., those of future generations, the disenfranchised, etc.).	<i>Ownership</i>	
1162: who wrote what's my name rihanna An accompanying music <b>video</b> , directed by Philip Andelman, portrays a romantic encounter between <b>Rihanna</b> and Drake in a grocery store along with romantic <b>scenes</b> between the pair and <b>Rihanna walking</b> through Manhattan's <b>Lower East Side</b> .	<i>What's My Name? (Rihanna song)</i>	
#1163: what year did isaac newton die In addition to his work on the <b>calculus</b> , as a <b>mathematician Newton</b> contributed to the study of power series , generalised the <b>binomial theorem</b> to non-integer exponents, and developed <b>Newton's</b> method for approximating the roots of a function	<i>Isaac Newton</i>	
1164: what is a constant in math? What it means for a <b>constant</b> to arise "naturally", and what makes a <b>constant</b> "interesting", is ultimately a matter of taste, and some mathematical constants are notable more for historical reasons than for their intrinsic mathematical interest.	<i>Mathematical constant</i>	
1178: who owns land rover The ongoing commercial success of the original <b>Land Rover</b> series models, and latterly the Range <b>Rover</b> in the 1970s in the midst of BL's well documented business troubles prompted the establishment of a separate <b>Land Rover company</b> but still under the BL umbrella, remaining part of the subsequent <b>Rover Group</b> in 1988, under the <b>ownership</b> of British Aerospace after the remains of British <b>Leyland</b> were broken up and privatised.	<i>Land Rover</i>	
1180: what states have legalized prostitution As with other countries, <b>prostitution</b> in the United States can be divided into three broad categories: <b>street prostitution</b> , <b>brothel prostitution</b> , and <b>escort prostitution</b> .	<i>Prostitution in the United States</i>	
1181: what states are on the east coast The <b>East Coast of the United States</b> , also known as the <b>Eastern</b> Seaboard or the Atlantic Seaboard and commonly shortened to <b>East Coast</b> , refers to the easternmost <b>coast</b> of the United States along the Atlantic Ocean .	<i>East Coast of the United States</i>	
#1183: when did sertraline come on the market Sertraline is primarily <b>prescribed</b> for major <b>depressive</b> disorder in adult outpatients as well as obsessive-compulsive panic , and social anxiety disorders in both adults and children.	<i>Sertraline</i>	
#1186: when did secretariat win <b>Secretariat</b> 's grandsire, Nasrullah, is also the great-great-grandson of 1977 Triple Crown <b>winner</b> Seattle Slew .	<i>Secretariat (horse)</i>	
1187: what are batteries made up of Battery recycling of <b>automotive batteries</b> reduces the need for resources required for manufacture of new <b>batteries</b> , diverts toxic <b>lead</b> from landfills, and prevents risk of improper disposal.	<i>Automotive battery</i>	
1189: what is a full job time?? <b>Full-time jobs</b> are often considered careers .	<i>Full-time</i>	
1190: what are the three primary colors in the subtractive color model A <b>subtractive color model</b> explains the mixing of a limited set of dyes, inks , paint pigments or natural <b>colorants</b> to create a wider range of <b>colors</b> , each the result of partially or completely <b>subtracting</b> (that is, absorbing) some wavelengths of light and not others.	<i>Subtractive color</i>	
1191: what are layers of the ionosphere It is distinguished because it is <b>ionized</b> by solar radiation.	<i>Ionosphere</i>	
1201: who plays as big bird <b>Big Bird</b> is a protagonist of the children's television show Sesame Street .	<i>Big Bird</i>	
1205: who founded walmart <b>Walmart</b> remains a family-owned <b>business</b> , as the company is controlled by the Walton family , who own a 48 percent stake in <b>Walmart</b> .	<i>Walmart</i>	
1206: what are some legal uses of meth Both dextromethamphetamine and racemic methamphetamine are <b>Schedule II controlled substances</b> in the United States, and similarly the production, distribution, sale, and possession of methamphetamine is restricted or <b>illegal</b> in many jurisdictions .	<i>Methamphetamine</i>	
1208: what is 6 pin din connector Mini-DIN is similar to the larger, older <b>DIN connector</b> .	<i>Mini-DIN connector</i>	
1211: Where does the word baptism come from While <b>John</b> the Baptist's use of a deep river for his <b>baptism</b> suggests <b>immersion</b> , pictorial and archaeological evidence of <b>Christian baptism</b> from the 3rd century onward indicates that a normal form was to have the candidate stand in <b>water</b> while <b>water</b> was poured over the upper body.	<i>Baptism</i>	
1212: what school did Zach Thomas play for before making it to the NFL Zachary Michael <b>Thomas</b> (born September 1, 1973) is a former American college and professional football <b>player</b> who was a linebacker in the National Football League (NFL) for thirteen seasons.	<i>Zach Thomas</i>	
1213: what order is the moth Most of this <b>order are moths</b> ; there are thought to be approximately 160,000 species of <b>moth</b> (nearly ten times the number of species of butterfly), with thousands of species yet to be described.	<i>Moth</i>	
1215: what are the uses for gui The term <b>GUI</b> is restricted to the scope of two-dimensional display screens with display resolutions able to describe generic information in the tradition of the computer <b>science</b> research at the PARC (Palo Alto Research Center).	<i>Graphical user interface</i>	
1217: who won anim cycle 12 <b>America's Next Top Model</b> , Cycle 12 is the twelfth <b>cycle</b> of <b>America's</b> Next Top <b>Model</b> and the sixth <b>season</b> to air on the CW network.	<i>America's Next Top Model, Cycle 12</i>	
1219: who built the globe A modern reconstruction of the <b>Globe</b> , named "Shakespeare's <b>Globe</b> ", opened in 1997 approximately from the site of the original theatre.	<i>Globe Theatre</i>	
1221: what are warehouse spreadsheets used for In computing , a data <b>warehouse</b> or enterprise data <b>warehouse</b> (DW, DWH, or EDW) is a database <b>used</b> for <b>reporting</b> and <b>data analysis</b> .	<i>Data warehouse</i>	
1224: Who Started the Mormon Church Today <b>Mormons</b> are understood to be <b>members</b> of The <b>Church of Jesus Christ of Latter-day Saints</b> (LDS Church).	<i>Mormons</i>	
1226: what part of the earth's structure is believed to consist of tectonic plates <b>Plate tectonics</b> (from the Late Latin <b>tectonicus</b> , from the "pertaining to building") is a scientific theory that describes the large-scale <b>movements</b> of Earth's <b>lithosphere</b> .	<i>Plate tectonics</i>	
1232: who won the 2009 super bowl The club became unexpected <b>winner</b> during the regular <b>season</b> , compiling a 9–7 <b>record</b> , and the playoffs with the aid of head <b>coach</b> Ken Whisenhunt , who was the Steelers' offensive coordinator in <b>Super Bowl XL</b> , and the re-emergence of <b>quarterback</b> Kurt Warner , who was the <b>Super Bowl MVP</b> in <b>Super Bowl XXXIV</b> with his former <b>team</b> , St. Louis Rams.	<i>Super Bowl XLIII</i>	
1233: what is a popular people meter The <b>People Meter</b> is an electronic method of television measurement that moved from active and diary-based to passive and <b>meter-monitored</b> .	<i>People meter</i>	
Finally in 1986, Nielsen developed an electronic <b>meter</b> , <b>People Meter</b> , to solve the problem.		
1235: what are the characteristics of bryophytes <b>Bryophytes</b> produce enclosed <b>reproductive</b> structures ( <b>gametangia</b> and <b>sporangia</b> ), but they produce neither flowers nor seeds , <b>reproducing</b> via spores .	<i>Bryophyte</i>	
#1236: when to use semicolon The <b>semicolon</b> (;) is a punctuation mark with several <b>uses</b> . The modern <b>uses</b> of the <b>semicolon</b> relate either to the listing of items or to the linking of related clauses . Ben Jonson was the first notable English writer to <b>use</b> the <b>semicolon</b> systematically.	<i>Semicolon</i>	
#1238: when did xbox release	<i>Xbox</i>	
Although these two are free while <b>Xbox</b> Live required a subscription, as well as broadband-only connection which was not completely adopted yet, <b>Xbox</b> Live was a success due to better servers , features such as a buddy list, and milestone like Halo 2 <b>released</b> in November 2004, which is the best-selling <b>Xbox</b> video game and was by far the most popular online game for years.		
1242: what is a roll play games Role-playing <b>games</b> also <b>include</b> single-player offline role-playing video <b>games</b> in which players <b>control</b> a character or team who undertake quests, and may <b>include</b> capabilities that advance using statistical mechanics.		
1254: what is a forward in soccer Modern team formations usually <b>include</b> one to three <b>forwards</b> ; <b>two</b> is the most common.		
1255: WHO QUALIFIES AS A DOMESTIC PARTNER IN CA Initially, <b>domestic partnerships</b> enjoyed very few privileges—principally just hospital-visitation <b>rights</b> and the <b>right</b> to be claimed as a next of kin of the estate of a deceased <b>partner</b> .		
1256: where did mark jackson play pro basketball? <b>Mark Jackson</b> (born April 1, 1963) is an American <b>basketball coach</b> , and former player who is the current head <b>coach</b> of the Golden State Warriors of the NBA .		
1267: what are stink bombs made of <b>Stink bomb</b> is a device designed to create an unpleasant smell . A prank <b>stink bomb</b> .		
1270: who said "A picture is worth a thousand words?" The discussion of "One Picture Worth Thousand Words" versus "One Picture Worth Ten Thousand Words" Wan yen I hue and 10,000 miles <b>worth</b> 10,000 books is cited in Information graphics where the concept of many in different disciplines and cultures.		
1287: what teams won super bowl Super Bowl III in 1969 was the first such game that carried the "Super Bowl" moniker, the names "Super Bowl I" and "Super Bowl II" were retroactively applied to the first two games.		
#1288: what year was christianity introduced to sub-saharan africa However, the most recent December 18, 2012 <b>Pew Forum</b> research <b>estimates</b> that in 2010, 6,010 <b>million Christians</b> , 3,270 <b>million traditional African religion</b> followers, 610,000 Muslims and 50,000 unaffiliated (no known <b>religion</b> ) peoples lived in South Sudan.		
#1295: when did expos become nationals In , the <b>Expos</b> won a division <b>championship</b> , won their first-ever playoff series by defeating the Philadelphia Phillies , 3-2, and advanced to the <b>National League Championship Series</b> , where they would go on to lose to the Los Angeles Dodgers , 3-2, in their only postseason appearance during the strike-shortened season.		
1296: what is a contingent fee with an attorney A <b>contingent fee</b> (in the United States ) or <b>conditional fee</b> (in England and Wales ) is any <b>fee</b> for services provided where the <b>fee</b> is payable only if there is a favourable result.		
1297: what is a D.O. stand for medical doctor Many <b>D.O.</b> physicians attend the same graduate <b>medical</b> education programs as their <b>M.D.</b> counterparts, and then take <b>M.D.</b> specialty <b>board</b> exams, while other <b>D.O.</b> graduates enter osteopathic programs and take <b>D.O.</b> specialty <b>board</b> examinations.		
1307: who reports the consumer price index A <b>consumer price index (CPI)</b> measures changes in the <b>price</b> level of a market basket of <b>consumer</b> goods and services purchased by households.		
1309: what are rocker arms? The effective leverage of the <b>arm</b> (and thus the force it can exert on the valve stem) is determined by the <b>rocker arm</b> ratio, the ratio of the distance from the <b>rocker arm</b> 's center of rotation to the tip divided by the distance from the center of rotation to the point acted on by the camshaft or pushrod.		
1313: where does the return address go on mail Should the <b>return address</b> be of a different state or country, the <b>mail</b> may be routed through that location for ease of <b>return</b> .		
#1314: when did steven adler play for guns and roses During the 2000s, <b>Adler</b> was the drummer of the band <b>Adler's Appetite</b> , and since 2012, he has held the same position in the band <b>Adler</b> .		
#1327: what year was gulf war The <b>war</b> is also known under other <b>names</b> , such as the Persian <b>Gulf War</b> , First <b>Gulf War</b> , <b>Gulf War I</b> , or the First <b>Iraq War</b> , before the term "Iraq <b>War</b> " became identified instead with the 2003 <b>Iraq War</b> (also referred to in the U.S. as "Operation Iraqi Freedom").		
#1335: what year was elvis born In 1968, after seven <b>years</b> away from the stage, he returned to live performance in a celebrated comeback television <b>special</b> that led to an extended <b>Las Vegas</b> concert residency and a string of profitable tours.		
1339: where did columbus really land in 1492 Though <b>Columbus</b> was not the first European explorer to reach the Americas (having been preceded by the Norse expedition led by Leif Ericson in the 11th century), <b>Columbus's voyages</b> led to the first <b>lasting</b> European contact with the Americas, inaugurating a period of European exploration, conquest, and colonization that <b>lasted</b> for several <b>centuries</b> .		
1340: what is a brindle boxer Boxers were <b>first</b> exhibited in a dog show for St. Bernards in Munich in 1895, the <b>first Boxer</b> club being founded the next year.		
1342: who invented the internet Since the mid-1990s, the <b>Internet</b> has had a revolutionary impact on culture and <b>commerce</b> , including the rise of near-instant communication by electronic <b>mail</b> , instant messaging , Voice over <b>Internet Protocol</b> (VoIP) "phone calls", two-way interactive video calls , and the <b>Wide Web</b> with its discussion forums , blogs , social networking , and online shopping sites.		
#1345: when does the royal standard fly? If the ancient <b>Royal Standard</b> of Scotland is flying above Holyrood <b>Palace</b> or Balmoral <b>Castle</b> , instead of the <b>Royal Standard</b> of the United Kingdom used in Scotland, it also indicates that the Queen is not in residence.		
1354: what the atmosphere on mercury The existence of a <b>atmosphere</b> had contentious before 1974, although by that time a consensus had formed that <b>Mercury</b> , like the Moon , lacked any substantial <b>atmosphere</b> .		
#1356: what year did the last monkeys go into space Before humans <b>went</b> into <b>space</b> , several animals were <b>launched</b> into <b>space</b> , including numerous <b>monkeys</b> , so that scientists could investigate the biological effects of <b>space travel</b> .		
1364: what religion is westminster abbey <b>Westminster Abbey</b> is a collegiate <b>church</b> governed by the <b>Dean</b> and Chapter of <b>Westminster</b> , as established by Royal charter of Queen Elizabeth I in 1560, which <b>created</b> it as the Collegiate <b>Church</b> of St Peter <b>Westminster</b> and a Royal Peculiar under the personal jurisdiction of the Sovereign.		
#1365: when did coca cola first come out The most common of these is Diet <b>Coke</b> , with others including <b>Caffeine-Free Coca-Cola</b> , Diet <b>Coke Caffeine-Free, Coca-Cola Cherry</b> , <b>Coca-Cola Zero</b> , <b>Coca-Cola Vanilla</b> , and special versions with lemon, lime or coffee.		
1366: what role do ombudsman play in the swedish government? In many countries where the <b>ombudsman</b> 's remit extends beyond dealing with alleged maladministration to promoting and protecting <b>human rights</b> , the <b>ombudsman</b> is recognised as the national <b>human rights</b> institution.		
1372: who played batman in dark knight Considered one of the <b>best films</b> of the 2000s and one of the <b>best superhero films</b> ever made , the <b>film</b> received <b>highly</b> positive <b>reviews</b> and <b>set</b> numerous records during its <b>theatrical run</b> .		
1384: who won the super in xli <b>Super Bowl XLI</b> was an <b>American football</b> game between the <b>American Football Conference</b> (AFC) <b>champion</b> Indianapolis <b>Colts</b> and the <b>National Football Conference</b> (NFC) <b>champion</b> Chicago Bears to decide the <b>National Football League</b> (NFL) <b>champion</b> for the 2006 <b>season</b> .		
1389: what are the parts of plant stems? In most <b>plants</b> <b>stems</b> are located above the soil surface but some <b>plants</b> have underground <b>stems</b> .		
#1394: what year did martin luther king die <b>Martin Luther King, Jr.</b> was established as a U.S. federal <b>holiday</b> in 1986.		
1400: what spanish speaking countries have the most world cup titles <b>The World Cup</b> is the world's most widely viewed sporting event; an <b>estimated</b> 715.1 <b>million</b> people <b>watched</b> the final match of the 2006 <b>FIFA World Cup held</b> in <b>Germany</b> .		
1402: who composed the Singapore national anthem <b>Singaporeans</b> are especially encouraged to sing the <b>national anthem</b> on occasions of <b>national</b> celebration or <b>national</b> significance such as at the <b>National Day Parade</b> , at <b>National Day</b> observance ceremonies conducted by educational institutions and government departments, and at sporting events at which <b>Singapore</b> teams are participating.		

#1411: when does v start	V (2009 TV series)
V stars Morena Baccarin , Lourdes Benedicto , Morris Chestnut , Joel Gretsch , Logan Huffman , Charles Mesure , Elizabeth Mitchell , Laura Vandervoort and Scott Wolf , and was executive produced by Scott Rosenbaum , Yves Simoneau , Scott Peters , and Jace Hall .	Death of John Lennon
#1416: when did spongebob first air	SpongeBob SquarePants
He teamed up with several Nickelodeon veterans and Rocko crew members, including creative director Derek Drymon (Action League Now! , Hey Arnold! , and Rocko's Modern Life) writers and directors Sherm Cohen , and Dan Povenmire , writer Tim Hill , actor and writer Martin Olson , animation director Alan Smart (all from Rocko's Modern Life), and story editor Merrivether Williams (The Angry Beavers ), who worked on the series for its first few seasons and switched to SpongeBob SquarePants in July 1999.	New Orleans
1418: what are d.o. of medicine	Doctor of Osteopathic Medicine
Many D.O. physicians attend the same graduate medical education programs as their M.D. counterparts, and then take M.D. specialty board exams while other D.O. graduates enter osteopathic programs and take D.O. specialty board examinations.	
1425: what is a medallion guarantee	Medallion signature guarantee
In the United States and Canada , a medallion signature guarantee is a special signature guarantee for the transfer of securities .	
A medallion signature guarantee is not the same as an acknowledgment by a notary public , in the sense that a "signature guarantee" is a certification by the institution that the signature is authentic, and an acknowledgment is a certification by a notary public attesting that the signer signed a document voluntarily.	
1438: what is "thin film" technology	Thin film
Thin films are also used in dye-sensitized solar cells .	
1443: what zones are tropical	Tropics
The tropics are also referred to as the tropical zone and the torrid zone (see geographical zone ).	
1444: what type of business is walmart	Walmart
Walmart remains a family-owned business , as the company is controlled by the Walton family , who own a 48 percent stake in Walmart .	
1445: what part of the plant are avocados	Avocado
The avocado (Persea americana) is a tree native to Central Mexico, classified in the flowering plant family Lauraceae along with cinnamon , camphor and bay laurel .	
#1449: when does air bag deploy	Airbag
Modern vehicles may contain multiple airbag modules in various side and frontal locations of the passenger seating positions, and sensors may deploy one or more airbags in an impact zone at variable rates based on the type, angle and severity of impact; the airbag is designed to only inflate in moderate to severe frontal crashes .	
1452: what are some chinese inventions	List of Chinese inventions
For the purposes of this list , inventions are regarded as technological firsts developed in China , and as such does not include foreign technologies which the Chinese acquired through contact, such as the windmill from the Middle East or the telescope from Early modern Europe .	
#1457: when did world begin	World War I
It was predominantly called the World War or the Great War from its occurrence until the start of World War II in 1939, and the First World War or World War I thereafter.	
1461: what two empires fought to control afghanistan	Afghanistan
Afghanistan ( ; ), officially the Islamic Republic of Afghanistan , is a landlocked sovereign state forming part of South Asia , Central Asia , and to some extent Western Asia .	
1462: what triggered the civil war	American Civil War
The American Civil War (ACW), also known as the War between the States or simply the Civil War (see naming ), was a civil war fought from 1861 to 1865 between the United States (the "Union" or the "North") and several Southern slave states that declared their secession and formed the Confederate States of America (the "Confederacy" or the "South").	
1466: who owns youtube	Youtube
Most of the content on YouTube has been uploaded by individuals, although media corporations including CBS , the BBC , Vevo , Hulu , and other organizations offer some of their material via the site, as part of the YouTube partnership program.	
1468: what are tires made of	Tire
Metal tires are still used on locomotives and railcars , and solid rubber (or other polymer) tires are still used in various non-automotive applications, such as some casters , carts , lawnmowers , and wheelbarrows .	
1469: who killed general warren in bunker hill	Joseph Warren
Warren had been commissioned a Major General in the colony's militia shortly before the June 17, 1775 Battle of Bunker Hill .	
His death, immortalized in John Trumbull 's painting, The Death of General Warren at the Battle of Bunker's Hill , June 17, 1775, galvanized the rebel forces, and he has been memorialized in many place names in the United States.	
1472: where do women ejaculation exactly coming from	Female ejaculation
The exact source and nature of the fluid continue to be a topic of debate among medical professionals, which is also related to doubts over the existence of the G-Spot .	
1473: who won the women's world cup	FIFA Women's World Cup
The first Women's World Cup tournament, named the Women's World Championship, was held in 1991, sixty-one years after the men's first FIFA World Cup tournament in 1930.	
#1478: when did texas become a state	Texas
Houston is the largest city in Texas and the fourth-largest in the United States , while San Antonio is the second largest in the state and seventh largest in the United States .	
1484: who owns smirnoff'	Smirnoff
Smirnoff products include vodka , flavoured vodka , and malt beverages .	
1485: who made the matrix	The Matrix
The success of the film led to the release of two feature film sequels, both written and directed by the Wachowskis, The Matrix Reloaded and The Matrix Revolutions .	
1486: what are some six sigma tools used	Six sigma
A six sigma process is one in which 99.9996% of the products manufactured are statistically expected to be free of defects (3.4 defects per million), although, as discussed below, this defect level corresponds to only a 4.5 sigma level.	
1493: what is a bad beat in poker	Bad beat
There is no consensus among poker players as to what exactly constitutes a bad beat and often players will disagree about whether a particular hand was a bad beat .	
1495: what part of beef are rouladen cut from?	Rouladen
Rouladen (or Rinderroulade, singular: rouleade ) is a German meat roulade usually consisting of bacon, onions, mustard and pickles wrapped in thinly sliced beef which is then cooked.	
1497: where does cellular respiration occur	Cellular respiration
While the overall reaction is a combustion reaction , no single reaction that comprises it is a combustion reaction .	
1500: what is a lapping machine	Lapping
The other form of lapping involves a softer material such as pitch or a ceramic for the lap , which is "charged" with the abrasive.	
The first type of lapping (traditionally called grinding ) typically involves rubbing a brittle material such as glass against a surface such as iron or glass itself (also known as the "lap" or grinding tool) with an abrasive such as aluminum oxide , jeweller's rouge , optician's rouge , emery , silicon carbide , diamond , etc., in between them.	
1502: who plays ethan in my babysitter's a vampire	My Babysitter's a Vampire (TV series)
My Babysitter's a Vampire ( French : Ma gardienne est un vampire ) is a 2011 Canadian television series, based on the television film of the same name .	
1513: what is a base SI unit	SI base unit
The International System of Units (SI) defines seven units of measure as a basic set from which all other SI units are derived .	
#1516: when did ww1 end?	World War I
It was predominantly called the World War or the Great War from its occurrence until the start of World War II in 1939, and the First World War or World War I thereafter.	
1519: what state is area code 419	Area codes 419 and 567
The main area code , 419 , was created as one of the original area codes in October 1947; the overlay area code 567 was created on January 1, 2002.	
1521: what is 9/11 bombings	September 11 attacks
The fourth plane , United Airlines Flight 93 , was targeted at the United States Capitol in Washington, D.C. , but crashed into a field near Shanksville, Pennsylvania , after its passengers tried to overcome the hijackers .	
#1528: when did charles dickens live	Charles dickens
Born in Portsmouth , England , Dickens left school to work in a factory after his father was thrown into debtors' prison .	
#1529: when did the civil rights movement begin	Civil rights movement
Civil rights movements ranging from the global LGBT rights movement to the global Women's rights movement to various racial minority rights movements around the world continue .	
#1532: who shot john lennon	Death of John Lennon
John Lennon was an English musician who gained worldwide fame as one of the founders of The Beatles , for his subsequent solo career, and for his political activism and pacifism .	
1534: what state is new orleans in	New Orleans
The New Orleans metropolitan area (New Orleans-Metairie-Kenner Metropolitan Statistical Area) had a population of 1,167,764 in 2010 and was the 46th largest in the United States .	
1535: what where the most important factors that led to the defeat of the democrates in 1968?	United States presidential election, 1968
This was the last election in which New York had the most votes in the electoral college (43 votes).	
#1537: when did kurt cobain kill himself	Kurt Cobain
Kurt Donald Cobain (February 20, 1967 – April 5, 1994) was an American musician and artist , best known as the lead singer , guitarist and primary songwriter of the grunge band Nirvana .	
1538: who starred in the original true grit	True Grit (1969 film)
True Grit is a 1969 American western film written by Marguerite Roberts and directed by Henry Hathaway .	
1539: what is a millwright worker	Millwright
Modern millwrights work with steel and other materials in addition to wood and must often combine the skills of several skilled trades in order to successfully fabricate industrial machinery or to assemble machines from pre-fabricated parts .	
1540: who created the tourbillon movement?	Tourbillon
Stéphane Tourbillon Movement ( ).	
1549: what is an msi file	Windows Installer
The installation information, and often the files themselves, are packaged in installation packages , loosely relational databases structured as COM Structured Storage and commonly known as "MSI files" , from their default file extension .	
#1550: what year did mexico gain independence from spain	Mexican War of Independence
The movement, which became known as the Mexican War of Independence , was led by Mexican-born Spaniards , Mestizos and Amerindians who sought independence from Spain .	
1553: what is a Four Lokos	Four Loko
The name "Four" is derived from the original energy drink's four main ingredients: alcohol , caffeine , taurine , and guarana .	
1555: what is a fret on a guitar	Fret
The neck of a guitar showing the nut (in the background, coloured white) and first four metal frets .	
#1559: what year was smokey the bear invented	Smokey Bear
Smokey Bear (often called Smokey the Bear or Smoke) is a mascot of the United States Forest Service created to educate the public about the dangers of forest fires .	
1561: who said tv is a vast wasteland	Newton N. Minow
His speech referring to television as a "vast wasteland" is cited even as the speech has passed its 50th anniversary .	
1575: who discovered the 2 moons of mars,phobos and deimos	Moons of Mars
It is possible that Mars may have moons smaller than 50 - 100 meters and a dust ring between Phobos and Deimos may be present but none have been discovered .	
1580: where does angelia jolie currently work	Angela Davis
Angela Yvonne Davis (born January 26, 1944) is an American political activist , scholar, and author .	
1587: who killed robert kennedy	Assassination of Robert F. Kennedy
The assassination of Robert Francis "Bobby" Kennedy , a United States Senator and brother of assassinated President John Fitzgerald "Jack" Kennedy , took place shortly after midnight on June 5, 1968, in Los Angeles , California , during the campaign season for the United States Presidential election, 1968 .	
#1589: when Harry met Sally case	When Harry Met Sally...
When Harry Met Sally... is a 1989 American romantic comedy film written by Nora Ephron and directed by Rob Reiner .	
When Harry Met Sally... grossed a total of US\$92 .8 million in North America .	
1590: who owns hamburger helper	Hamburger Helper
The Hamburger Helper mascot is the "Helping Hand", an anthropomorphic animated , four fingered left-hand glove, which appears in the product's telecommunications and on the packages .	
1591: what makes a dwarf planet	Dwarf planet
The exclusion of dwarf planets from the roster of planets by the IAU has been both praised and criticized; it was said to be the "right decision" by Mike Brown , who discovered and other new dwarf planets , but has been rejected by Alan Stern , who had coined the term dwarf planet in 1990 .	
#1593: when did jack lalanne die	Jack LaLanne
On the occasion of LaLanne's death , Schwarzenegger credited LaLanne for being "an apostle for fitness" by inspiring "billions all over the world to live healthier lives." and, as governor of California , had earlier placed him on his Governor's Council on Physical Fitness .	
1598: who made the original care bears	Care Bears
The Care Bears appeared in their own TV specials called The Care Bears in the Land Without Feelings (1983) and The Care Bears Battle the Freeze Machine (1984) .	
#1601: when do solar eclipses happen?	Solar eclipse
Earth's orbit is called the elliptic plane as the Moon's orbit must cross this plane in order for an eclipse (both solar as well as lunar ) to occur .	
1606: what latitude is tropic of cancer	Tropic of cancer
The Tropic of Cancer , also referred to as the Northern tropic , is the circle of latitude on the Earth that marks the most northerly position which the Sun may appear directly overhead at its zenith .	
1611: what state is the capital in	List of capitals in the United States
In addition, each of the 50 U.S. states and the five principal territories of the United States maintains its own capital .	
1612: what are square diamonds called?	Princess cut
The square princess cut diamond is usually slightly cheaper than round brilliant cut diamonds of the same carat weight because it retains about 80% of the rough diamond , as opposed to the round brilliant which retains only about 50% of the rough rough .	
1615: where scottsdale?	Scottsdale, Arizona
Scottsdale is bordered to the west by Phoenix and Paradise Valley , to the north by Carefree , to the south by Tempe , and to the east by Fountain Hills and the Salt River Pima-Maricopa Indian Community .	
1616: who has brad pitt dated	Brad Pitt
In addition, Pitt owns a production company, Plan B Entertainment , whose productions include The Departed (2006), which won the Academy Award for Best Picture , and Moneyball , which garnered a Best Picture nomination .	
1617: what is a dogs classification	Dog
MTDNA evidence shows an evolutionary split between the modern dog's lineage and the modern wolf's lineage around 100,000 years ago but , the oldest fossil specimens genetically linked to the modern dog's lineage date to approximately 33,000-36,000 years ago .	
1622: what nationality is wendy williams	Wendy Williams
She hosts a syndicated television talk show , The Wendy Williams Show .	
1626: what president was theodore roosevelt	Theodore Roosevelt
Theodore Roosevelt was 42 years old when sworn in as President of the United States in 1901, making him the youngest president ever; he beat out the youngest elected president , John F. Kennedy , by only one year .	
1628: who starred in webster	Webster (TV series)
Webster is an American situation comedy that aired on ABC from September 16, 1983 until May 8, 1987, and in first-run syndication from September 21, 1987 until March 10, 1989 .	
1633: what is a monarch to a monarchy	Monarchy
The monarchs of Cambodia , Japan , Jordan , Malaysia and Morocco "reign , but do not rule" although there is considerable variation in the amount of authority they wield .	
1637: what is an "N.M?"	Newton metre
The symbolic form is N m or N.m .	
1652: who made hubble telescope	Hubble Space Telescope
Although not the first space telescope , Hubble is one of the largest and most versatile, and is well known as both a vital research tool and a public relations boon for astronomy .	
#1659: when did andrea doria sink	SS Andrea Doria
Named after the 16th-century Genoese admiral Andrea Doria , the ship had a gross register tonnage of 29,100 and a capacity of about 1,200 passengers and 500 crew .	
1673: who wrote second corinthians	Second Epistle to the Corinthians
The Second Epistle to the Corinthians , often referred to as Second Corinthians (and written as 2 Corinthians), is the eighth book of the New Testament of the Bible .	
1675: who first synthesized heroin	Heroin

Mexican cartels are also known to <b>produce</b> a third type of illicit <b>heroin</b> , commonly called black tar , which results from a simplified, <b>quicker</b> synthesis procedure and contains a high percentage of <b>morphine</b> derivatives other than <b>heroin</b> , such as 6-monoacetylmorphine (6-MAM).	
1678: what is a wiki platform	<i>Wiki</i>
Wikis are powered by <b>wiki</b> software .	
1686: who plays dumbledore in harry potter 6	<i>Albus Dumbledore</i>
Dumbledore is <b>portrayed</b> by Richard <b>Harris</b> in the <b>film</b> adaptions of <b>Harry Potter</b> and the Philosopher's <b>Stone</b> and <b>Harry Potter</b> and the <b>Chamber</b> of Secrets .	
1689: what type of land is savannah	<i>Savanna</i>
<b>Typical tropical savanna</b> in Northern Australia demonstrating the high tree density and regular spacing characteristic of many <b>savannas</b> .	
1695: who set the world record for women for high jump	<i>High jump</i>
Javier Sotomayor ( Cuba ) is the current men's <b>record</b> holder with a <b>jump</b> of <b>set</b> in 1993, the longest standing <b>record</b> in the history of men's <b>high jump</b> .	
1701: where does ground pepper come from	<i>Black pepper</i>
Peppercorns, and the <b>ground pepper derived</b> from them, may be described simply as <b>pepper</b> , or more precisely as black <b>pepper</b> (cooked and dried unripe fruit), green <b>pepper</b> (dried unripe fruit) and white <b>pepper</b> (dried ripe seeds).	
1707: what type of ecosystem does stingrays live in	<i>Stingray</i>
They are classified in the suborder Myliobatoidei of the order Myliobatiformes , and consist of eight families: Hexatrygonidae (sixgill stingray), Plesiotidae (deep water stingray), Urolophidae (stingrays), Urotrygonidae (round rays), Dasyatidae (whiptail stingrays), Potamotrygonidae (river stingrays), Gymnuridae (butterfly rays), and Myliobatidae (eagle rays).	
1712: what religion is primary in africa?	<i>Religion in Africa</i>
<b>Religious distribution in Africa</b>	
1713: who won season 2 of project runway	<i>Project Runway (season 2)</i>
<b>Project Runway Season 2</b> was the second <b>season</b> of Bravo 's successful <b>Project Runway</b> , a reality competition for fashion designers.	
1715: what is a mule in coins	<i>Mule (coin)</i>
The name derives from the <b>mule</b> , the hybrid offspring of a horse and a donkey, due to such a <b>coin</b> having two sides intended for different <b>coins</b> , much as a <b>mule</b> has parents of two different species.	
1721: what is an arc in a story plot	<i>Story arc</i>
Although <b>story arcs</b> have existed for decades, the term " <b>story arc</b> " was coined in 1988 in relation to the television series <i>Wiseguy</i> , and was quickly adapted for other uses.	
1727: who sang what a wonderful world	<i>What a Wonderful World</i>
"What a <b>Wonderful World</b> " is a song written by Bob Thiele (as "George Douglas") and George David Weiss .	
1728: what month is the president inaugurated	<i>United States presidential inauguration</i>
(Prior to the Twentieth Amendment , the <b>date</b> was <b>March 4</b> , the <b>day</b> of the year on which the Constitution of the United States first took effect in 1789; the last <b>inauguration</b> to take place on the older <b>date</b> was Franklin D. Roosevelt 's first one on <b>March 4</b> , 1933.)	
1751: what are club seats	<i>Club seating</i>
<b>Club level seating</b> is a special section of <b>seating</b> in modern sports stadiums .	
1766: what is a redshirt freshman football player	<i>Redshirt (college sports)</i>
For example, a coach may choose to <b>redshirt</b> a <b>player</b> who is then referred to as a <b>redshirt freshman</b> or simply a <b>redshirt</b> .	
#1776: what year did South Africa become a team in rugby	<i>South Africa national rugby union team</i>
The <b>South Africa national rugby union team</b> (known as the Springboks) represents <b>South Africa in rugby union</b> .	
1782: who did john f kennedy run against?	<i>John F. Kennedy</i>
After military <b>service</b> as <b>commander</b> of the Motor Torpedo Boats PT-109 and PT-59 during World War II in the South Pacific , <b>Kennedy</b> represented Massachusetts' 11th congressional district in the U.S. House of Representatives from 1947 to 1953 as a Democrat .	
1785: what season is dexter on	<i>Dexter (TV series)</i>
Set in Miami , the show's first <b>season</b> was largely based on the <b>novel</b> Darkly Dreaming <b>Dexter</b> , the first of the <b>Dexter</b> series <b>novels</b> by Jeff Lindsay .	
1786: What are context effects of memory?	<i>Context-dependent memory</i>
However, the research <b>literature</b> on context-dependent <b>memory</b> describes a number of different types of contextual information that may affect recall such as environmental context-dependent <b>memory</b> , state-dependent learning , cognitive context-dependent <b>memory</b> and mood-congruent <b>memory</b> .	
1789: who won the 1998 world cup	<i>1998 FIFA World Cup</i>
The <b>1998 FIFA World Cup</b> was the 16th <b>FIFA World Cup</b> , the <b>world</b> championship for men's national association football teams.	
1791: where did the vietnamese settle in america	<i>Vietnamese American</i>
A <b>Vietnamese</b> American () is an American of <b>Vietnamese descent</b> .	
#1793: when did hitler kill himself	<i>Death of Adolf Hitler</i>
Accounts differ as to the cause of <b>death</b> ; one that he <b>died</b> by poison only and another that he <b>died</b> by a self-inflicted gunshot, while biting down on a cyanide capsule.	
#1796: when did world war 2 end	<i>World War II</i>
With an invasion of the <b>Japanese</b> archipelago imminent, and the Soviet <b>Union</b> having declared <b>war</b> on <b>Japan</b> by invading <b>Manchuria</b> , <b>Japan</b> surrendered on 15 August 1945, <b>ending</b> the <b>war</b> in Asia and cementing the total <b>victory</b> of the Allies over the Axis.	
1805: who won the 1967 nba championship	<i>1967 NBA Finals</i>
The <b>1967 NBA World Championship Series</b> was the <b>championship</b> series of the 1966-67 National Basketball Association season , and was the conclusion of the <b>1967 NBA Playoffs</b> .	
1806: who make airbus	<i>Airbus</i>
Airbus began as a consortium of <b>aerospace manufacturers</b> , <b>Airbus</b> Industrie.	
1809: who invented the television	<i>History of television</i>
As <b>electronic camera</b> and <b>display tubes</b> were perfected, electromechanical <b>television</b> gave way to all-electronic systems in nearly all <b>applications</b> .	
1812: what are the side effects for lyme disease	<i>Lyme disease</i>
<b>Lyme disease</b> , <b>Lyme</b> borreliosis is an infectious <b>disease</b> caused by at least three species of bacteria belonging to the genus <i>Borrelia</i> .	
1814: who wrote the song for star wars	<i>Star Wars music</i>
Additionally, music for <b>Star Wars</b> : The <b>Clone Wars</b> was <b>written</b> by Kevin Kiner, and further music has been composed for <b>Star Wars</b> video <b>games</b> and <b>works</b> in other media.	
1822: where does ray lamontagne live	<i>Ray lamontagne</i>
Raymond "Ray" Charles Jack LaMontagne (; born June 18, 1973) is an American singer-songwriter .	
1827: who made facebook	<i>Facebook</i>
Critics , such as <b>Facebook</b> Detox , state that <b>Facebook</b> has turned into a national obsession in the United States, resulting in vast amounts of time lost and encouraging narcissism.	
1831: where do you find iodine	<i>Iodine</i>
Because of this function, radioisotopes of <b>iodine</b> are concentrated in the thyroid gland along with nonradioactive <b>iodine</b> .	
1833: what kind of people are on the show skins	<i>Skins (UK TV series)</i>
Other ventures to expand the brand have <b>included</b> a failed North <b>American</b> adaptation , which <b>aired</b> on MTV in 2011 but it was <b>canceled</b> after one season after <b>advertisers</b> abandoned the <b>series</b> in response to low <b>ratings</b> and the <b>significant</b> controversy which arose over its depiction of <b>teen</b> sexuality.	
1834: who played the drums in the band cream back in 1968	<i>Cream (band)</i>
<b>Cream</b> were inducted into the Rock and Roll <b>Hall of Fame</b> in 1993.	
1837: what president made decision to buy louisiana	<i>Louisiana Purchase</i>
The <b>Louisiana Purchase</b> ("Sale of Louisiana") was the acquisition by the United States of America in 1803 of France 's claim to the territory of <b>Louisiana</b> .	
1838: who sang cool jerk	<i>Cool Jerk</i>
In the feature film Home Alone 2: Lost in New York , Uncle Frank ( Gerry Bamman ) <b>sings</b> "Cool Jerk" in the shower.	
1839: what is a league in the sea	<i>League (unit)</i>
A <b>league</b> is a <b>unit</b> of length (or, rarely, area ).	
#1845: when will ie9 be released	<i>Internet Explorer 9</i>
The system requirements for <b>Internet Explorer 9</b> are <b>Windows 7</b> , <b>Windows Server 2008 R2</b> , <b>Windows</b> Vista Service Pack 2 or <b>Windows Server 2008</b> SP2 with the <b>Platform Update</b> .	
#1849: when did classification of races begin	<i>Race (human classification)</i>
While <b>biologists</b> sometimes <b>use</b> the <b>concept</b> of <b>race</b> to make <b>distinctions</b> among fuzzy sets of traits, others in the <b>scientific</b> community suggest that the idea of <b>race</b> often is used in a naive or simplistic way, i.e. that among humans, <b>race</b> has no taxonomic significance: all living humans belong to the same <b>species</b> , <b>Homo sapiens</b> and <b>subspecies</b> , <b>Homo sapiens sapiens</b> .	
1855: what states has the electric chair	<i>Electric chair</i>
Execution by <b>electrocution</b> usually performed using an <b>electric chair</b> , is an execution <b>method</b> originating in the United States in which the condemned person is strapped to a specially built wooden <b>chair</b> and <b>electrocuted</b> through <b>electrodes</b> placed on the body.	
1857: what are private labels	<i>Private label</i>
McBride plc is an example of a European based provider of <b>private label</b> household and <b>personal care</b> products.	
#1866: When did the New Deal start	<i>New Deal</i>
The <b>New Deal</b> produced a <b>political</b> realignment, making the Democratic <b>Party</b> the majority (as well as the <b>party</b> that held the White House for seven out of nine Presidential terms from 1933 to 1969), with its base in <b>liberal</b> ideas, the white <b>South</b> , traditional <b>Democrats</b> , big city machines, and the <b>newly</b> empowered labor unions and ethnic minorities.	
#1873: when did daylight savings time start	<i>Daylight saving time</i>
Daylight saving time (DST)—also summer <b>time</b> in British <b>English</b> — is the practice of advancing clocks during the lighter months so that evenings have more <b>daylight</b> and mornings have less.	
1874: What region of France is Montargis in?	<i>Montargis</i>
Montargis is a commune in the Loiret department in north-central <b>France</b> on the Loing river.	
#1877: what year lord of rings made?	<i>The Lord of the Rings</i>
The enduring <b>popularity</b> of <b>The Lord of the Rings</b> has led to numerous references in <b>popular culture</b> , the founding of many societies by fans of <b>Tolkien's works</b> , and the <b>publication</b> of many <b>books</b> about <b>Tolkien</b> and his <b>works</b> .	
1878: what time will the world end on may 21	<i>2011 end times prediction</i>
The <b>2011 end times</b> prediction made by American Christian radio <b>host</b> Harold Camping stated that the Rapture and Judgment Day would take <b>place</b> on May 21, 2011 , and that the <b>end of the world</b> would take <b>place</b> five months later on October 21, 2011.	
1885: what is file based system	<i>File system</i>
Some <b>file systems</b> are "virtual", in that the "files" supplied are <b>computed</b> on request (e.g. <b>proofs</b> ) or are merely a mapping into a <b>different file system</b> used as a backing store.	
1887: What political conflicts marked the presidency of William Howard Taft?	<i>William Howard Taft</i>
William Howard Taft (September 15, 1857 – March 8, 1930) was the 27th <b>President</b> of the United States (1909–1913) and later the tenth <b>Chief Justice</b> of the United States (1921–1930).	
1890: what temperature is a salt ice bath	<i>Cooling bath</i>
Cooling <b>baths</b> are generally one of two types: (a) a cold fluid (particularly liquid nitrogen , water , or even air) — but most commonly the term refers to (b) a mixture of 3 components: (1) a cooling agent (such as dry <b>ice</b> or water <b>ice</b> ) ; (b) a liquid 'carrier' (such as liquid water, ethylene glycol, acetone, etc.) , which transfers heat between the <b>bath</b> and the vessel; ; and (c) an additive to depress the melting-point of the solid/liquid system.	
#1893: when slavery abolished	<i>Slavery in the United States</i>
However, by 1804, all <b>states</b> north of the Mason and Dixon Line had either <b>abolished slavery</b> outright or <b>passed laws</b> for the gradual <b>abolition of slavery</b> .	
1905: what are circumpolar constellations	<i>Circumpolar constellation</i>
In the northern hemisphere, we will always be able to see stars and <b>constellations</b> in the northern <b>circumpolar</b> sky, while in the southern hemisphere, we will always be able to see stars and <b>constellations</b> in the southern <b>circumpolar</b> sky.	
1906: where do sesame seeds come from	<i>Sesame</i>
The world's <b>largest</b> exporter of <b>sesame seeds</b> was India, and Japan the <b>largest</b> importer.	
1912: who owned kansas before it became a state	<i>Kansas</i>
When officially opened to <b>settlement</b> by the U.S. <b>government</b> in 1854, abolitionist Free-Staters from New England and pro-slavery <b>settlers</b> from neighboring Missouri rushed to the territory to determine if Kansas would become a free <b>state</b> or a <b>slave state</b> .	
1915: where did the olmecs come from	<i>Olmec</i>
Among other "firsts", the <b>Olmec</b> appeared to practice ritual bloodletting and played the <b>Mesoamerican</b> ballgame , hallmark of nearly all subsequent <b>Mesamerican</b> societies.	
1918: what states allow same sex marriage	<i>Same-sex marriage in the United States</i>
The Defense of Marriage Act (DOMA), enacted in 1996, prevents the <b>federal</b> government from <b>recognizing same-sex marriages</b> and allows each <b>state</b> to refuse recognition of <b>same-sex marriages</b> performed in other <b>states</b> .	
1923: who won the most nba championships	<i>List of NBA players with most championships</i>
Saul won consecutive <b>championships</b> with the Rochester Royals and the Minneapolis Lakers in the 1950s, while Kerr <b>won</b> consecutive <b>championships</b> with the Bulls and the Spurs in the 1990s.	
Horry <b>won</b> seven <b>championships</b> with the Houston Rockets , the Los Angeles Lakers and the San Antonio Spurs , while Salley <b>won</b> four <b>championships</b> with the Detroit Pistons , the Bulls and the Lakers.	
#1925: when did the word fuck begin	<i>Fuck</i>
In modern usage, <b>fuck</b> and its derivatives (such as <b>fucker</b> and <b>fucking</b> ) can be used in the position of a noun , a verb , an adjective or an adverb .	
1926: what are the 7 continents	<i>Continent</i>
Depending on the convention and model , some <b>continents</b> may be consolidated or subdivided: for example, Eurasia is most often subdivided into Europe and Asia (red shades), while North and South America are sometimes recognized as one <b>American continent</b> (green shades).	
1930: what state is milwaukee in	<i>Milwaukee</i>
Known for its brewing traditions, major new additions to the city include the Milwaukee Riverwalk , the Delta Center (formerly "Frontier Airlines Center"), Miller Park , an internationally renowned addition to the Milwaukee Art Museum , Milwaukee Repertory Theater , and Pier Wisconsin , as well as major <b>renovations</b> to the U.S. Cellular Arena .	
#1935: what years was the 18th century	<i>18th century</i>
To historians who expand the <b>century</b> to include larger <b>historical</b> movements, the "long" 18th century may run from the Glorious Revolution of 1688 to the battle of Waterloo in 1815 or even later.	
#1937: what year was the eiffel tower made	<i>Eiffel Tower</i>
The Eiffel Tower ( ) is an iron lattice tower located on the Champ de Mars in Paris , named after the engineer Gustave Eiffel , whose company designed and built the tower.	
#1942: when monopoly came out	<i>History of the board game Monopoly</i>
Also in the 1970s, Professor Ralph Anspach , who had himself published a board game intended to illustrate the principles of both <b>monopolies</b> and trust busting , fought Parker Brothers and its then parent company, General Mills , over the trademarks of the <b>Monopoly</b> board game.	
1954: who passed no child left behind	<i>No Child Left Behind Act</i>
The No Child Left Behind Act of 2001 (NCLB) is a United States <b>Act</b> of Congress that is a <b>reauthorization</b> of the <b>Elementary</b> and Secondary Education <b>Act</b> , which included Title I , the government's flagship program for disadvantaged students.	
1956: what are the charges against Casey Anthony	<i>Death of Caylee Anthony</i>
Caylee Marie Anthony (August 9, 2005 – 2008) was a two-year-old American girl who <b>lived</b> in Orlando, Florida with her mother, Casey Marie Anthony , and her <b>maternal</b> grandparents, George and Cindy Anthony .	
1965: when did world war 1 start	<i>World War I</i>
It was predominantly called the <b>World War</b> or the <b>Great War</b> from its occurrence until the start of World War II in 1939, and the <b>First World War</b> or <b>World War I</b> thereafter.	
1966: what species is a spider	<i>Spider</i>
It now appears that the spiral orb <b>web</b> may be one of the earliest forms, and <b>spiders</b> that produce <b>tangled</b> cobwebs are more abundant and diverse than <b>orb-web spiders</b> .	
#1968: when can you use a defibrillator	<i>Defibrillation</i>
Some external units, known as <b>automated</b> external defibrillators (AEDs), <b>automate</b> the diagnosis of treatable rhythms, meaning that lay responders or bystanders are able to <b>use</b> them successfully with little, or in some cases no training at all.	
1969: who created the cato institute	<i>Cato Institute</i>
According to the 2011 Global Go To Think Tank Index, Cato is the 6th most influential US based think tank , ranking 3rd Economic Policy and 2nd in Social Policy.	
#1977: when does a demand curve shift?	<i>Demand curve</i>
The <b>demand curve</b> for all consumers together follows from the <b>demand curve</b> of every individual consumer: the individual <b>demands</b> at each price are added together.	
1992: Where Elephants Live	<i>Elephant</i>
Traditionally, two species are recognised, the <b>African elephant</b> ( <i>Loxodonta africana</i> ) and the <b>Asian elephant</b> ( <i>Elephas maximus</i> ), although some evidence suggests that <b>African bush elephants</b> and <b>African forest elephants</b> are separate species ( <i>L. africana</i> and <i>L. cyclotis</i> respectively).	
1997: what are the quad muscles	<i>Quadriceps femoris muscle</i>

The <b>quadriceps</b> femoris ( Latin for "four-headed <b>muscle</b> of the femur " ), also called simply the <b>quadriceps</b> , <b>quadriceps extensor</b> , <b>quads</b> , is a large <b>muscle</b> group that includes the four prevailing <b>muscles</b> on the front of the thigh .	
#2005: when was james madison in the house of representatives	James Madison
<b>James Madison, Jr.</b> ( March 16, 1751 ( O.S. March 5 ) – June 28, 1836 ) was an American statesman and political theorist, the fourth <b>President</b> of the United States (1809–1817).	
2007: what is and where is hydraulic fluid and used for	Hydraulic fluid
<b>Hydraulic fluids</b> , also called <b>hydraulic</b> liquids, are the medium by which power is transferred in <b>hydraulic</b> machinery .	
2008: Where Are Mahindra Tractors Made	Mahindra Tractors
<b>Mahindra Tractors</b> , the farm equipment division of <b>Mahindra &amp; Mahindra</b> , builds and sources <b>tractors</b> that are sold worldwide across six continents .	
2012: what is middle class in the us	American middle class
One of the first major studies of the <b>middle class</b> in <b>America</b> was White Collar: The <b>American Middle Classes</b> , published in 1951 by sociologist C. Wright Mills Mills Wright Mills .	
2013: what is the scientific name of the eastern tiger salamander?	Tiger Salamander
The proper common <b>name</b> is the <b>eastern tiger salamander</b> , to differentiate it from other closely related species .	
2020: what was the steelworkers strike	Steel strike of 1919
The <b>strike</b> began September 21, 1919, and <b>collapsed</b> on January 8, 1920.	
#2022: When was the first Mary Poppins book written	Mary Poppins
<b>Mary Poppins</b> is the title character of a series of children's <b>books written</b> by P. L. <b>Travers</b> .	
2024: what is sims language	Simslish
Initially, inspired by the Native American code talkers of World War II , <b>Sims</b> creator Will Wright and <b>language</b> expert Marc Gimbel suggested experimenting with the Navajo <b>language</b> to create Simlish .	
2037: who is elizabeth from general hospital who are the boys fathers	Elizabeth Webber
<b>Elizabeth</b> is part of two supercouple pairings, Lucky Spencer and <b>Elizabeth Webber</b> and Jason Morgan and <b>Elizabeth Webber</b> .	
2038: what is reagan known for	Presidency of Ronald Reagan
The <b>United States</b> presidency of Ronald <b>Reagan</b> , also known as the <b>Reagan</b> administration, was a Republican administration headed by Ronald <b>Reagan</b> from <b>January 20</b> , 1981, to <b>January 20</b> , 1989.	
2040: what is evoked otoacoustic emissions	Otoacoustic emission
Broadly speaking, there are two types of <b>otoacoustic emissions</b> : spontaneous <b>otoacoustic emissions</b> (SOAEs), which can occur without external stimulation, and <b>evoked otoacoustic emissions</b> (EOAEs), which require an <b>evoking stimulus</b> .	
2041: what is chep pallet	CHEP
<b>CHEP</b> offers wooden and plastic <b>pallets</b> , small display <b>pallets</b> , crates and IBC containers .	
#2046: when was scooby doo created	Scooby-Doo
This Saturday morning <b>cartoon</b> series featured four teenagers—Fred Jones , Daphne Blake , Velma Dinkley , and Norville "Shaggy" Rogers—and their talking brown Great Dane <b>dog</b> named Scooby-Doo , who solve <b>mysteries</b> involving supposedly <b>supernatural</b> creatures through a series of antics and missteps .	
2054: who is the CEO of FACEBOOK	Mark Zuckerberg
In 2004, at the age of twenty-three years, Zuckerberg became a billionaire as a result of <b>Facebook</b> and the number of <b>Facebook</b> users worldwide reached a total of one billion in 2012.	
2057: What Is The Largest Whale	Blue whale
The blue <b>whale</b> (Balaenoptera musculus) is a marine mammal belonging to the suborder of baleen <b>whales</b> (called Mysticeti ).	
2058: What Is Range in Math	Range (mathematics)
This is the current usage for <b>range</b> in computer science .	
Sometimes " <b>range</b> " refers to the codomain and sometimes to the image .	
Some books say that <b>range</b> of this function is its codomain, the set of all real numbers, reflecting that the function is real-valued .	
Other books say that the <b>range</b> is the function's image, the set of non-negative real numbers, reflecting that a number can be the output of this function if and only if it is a non-negative real number .	
In this case, the larger set containing the <b>range</b> is called the codomain .	
In mathematics, the <b>range</b> of a function refers to either the codomain or the image of the function, depending upon usage .	
2060: who is on the hundred dollar bill	United States one hundred-dollar bill
The <b>United States</b> one <b>hundred-dollar bill</b> (\$100) is a <b>denomination</b> of <b>United States</b> currency .	
2061: what is the ingredient in mustard	Mustard (condiment)
The other four <b>mustards</b> pictured are a simple table <b>mustard</b> with turmeric coloring (center left), a Bavarian sweet <b>mustard</b> (center-right), a Dijon <b>mustard</b> (lower-left), and a coarse French <b>mustard</b> made mainly from black <b>mustard</b> seeds (lower-right).	
2062: who are the members of the climax blues band?	Climax Blues Band
The <b>Climax Blues Band</b> (originally known as the <b>Climax Chicago Blues Band</b> ) were formed in Stafford , England in 1968.	
#2064: when was the web invented	World Wide Web
With a <b>web browser</b> , one can view <b>web</b> pages that may contain text, images, videos, and other multimedia , and navigate between them via hyperlinks .	
2066: what is kathmandu known for	Kathmandu
<b>Kathmandu</b> 's sister cities ( Lalitpur Patan ) and Bhaktapur are integral to <b>Kathmandu</b> 's cultural heritage, <b>tourism</b> industry, and <b>economy</b> ; therefore UNESCO's World Heritage Site lists all three cities' monuments and attractions together under one heading, " <b>Kathmandu</b> Valley: UNESCO World Heritage Site " .	
2071: where is mark sanchez from	Mark Sanchez
Despite a subpar performance, <b>Sanchez</b> led the Jets to the AFC <b>Championship Game</b> , a losing effort to the Indianapolis Colts , becoming the fourth rookie quarterback in NFL history to win his first playoff <b>game</b> and the second to win two playoff <b>games</b> .	
2080: where is cougar town filmed	Cougar Town
The show was created by Bill <b>Lawrence</b> and Kevin <b>Biegel</b> and is <b>produced</b> by Doozer (Lawrence's company) and Coquette <b>Productions</b> in association with <b>ABC Studios</b> .	
2089: where is green bay packers from	Green Bay Packers
The <b>Green Bay Packers</b> have won 13 league <b>championships</b> (more than any other team in the NFL), including nine NFL <b>championships</b> prior to the <b>Super Bowl era</b> and four <b>Super Bowl</b> victories—in 1967 ( <b>Super Bowl I</b> ), 1968 ( <b>Super Bowl II</b> ), 1997 ( <b>Super Bowl XXXI</b> ) and 2011 ( <b>Super Bowl XLV</b> ).	
#2091: when was the world of coca cola built	World of Coca-Cola
Its well-known advertising as well as a host of entertainment areas and attractions, and is <b>located</b> in Atlanta , Georgia (where the company's headquarters are <b>located</b> ) at Pemberton Place (named in honor of John Pemberton, the inventor of Coca-Cola ).	
2096: what is the highest mountain in america and where is it located?	Mount McKinley
<b>Mount McKinley</b> (also known as Denali taken from the Inuit Kotukon Athabaskan language meaning "The Great One") is the <b>highest mountain</b> peak in the United States and in North <b>America</b> with a summit <b>elevation</b> of above sea level .	
2097: what was the first honda car	Honda S600
Available as a roadster – bearing strong resemblance to the <b>Honda</b> S500 – and as a fastback coupé – introduced in March 1965 – the S600 was the <b>first Honda</b> available in two trim levels.	
2100: what is the population of san francisco	San Francisco
<b>San Francisco</b> ( ), officially the City and <b>County</b> of San Francisco , is the leading financial and cultural center of Northern California and the <b>San Francisco</b> Bay Area .	
2102: who is the writer of the beowulf poem?	Beowulf
<b>Beowulf</b> (, in Old English or ) is the conventional title of an Old English heroic epic poem consisting of 3182 alliterative long lines , set in <b>Scandinavia</b> , commonly cited as one of the most important works of Anglo-Saxon literature .	
2114: what is sherlock holmes job?	Sherlock Holmes
The <b>character</b> grew tremendously in popularity with the first <b>series</b> of <b>short</b> stories in The Strand Magazine , beginning with A Scandal in Bohemia in 1891; further <b>series</b> of <b>short</b> stories and two novels <b>published</b> in <b>serial</b> form appeared between then and 1927.	
2119: what is social security card used for	Social Security number
A Social <b>Security</b> number may be obtained by applying on Form SS-5, "Application for A Social <b>Security</b> Number Card".	
2120: what is the minimalist trend	Minimalism
<b>Minimalism</b> is any design or <b>style</b> in which the simplest and fewest elements are used to create the maximum effect .	
2126: who is the current Chief Justice of the U.S. supreme court?	Chief Justice of the United States
The <b>Chief Justice</b> of the United States is the head of the United States federal court system (the judicial branch of the federal government of the United States ) and the <b>chief</b> judge of the Supreme Court of the United States .	
#2129: when was the american labor union formed	Labor unions in the United States
The economist Joseph Stiglitz has asserted that, "Strong <b>unions</b> have <b>helped</b> to reduce inequality, whereas weaker <b>unions</b> have made it easier for CEOs , sometimes <b>working</b> with market forces that they have <b>helped</b> shape, to increase it."	
2132: what is the largest credit union	Navy Federal Credit Union
Navy Federal <b>Credit Union</b> (or Navy Federal) is a <b>credit union</b> headquartered in Vienna, Virginia , chartered and regulated under the authority of the National <b>Credit Union</b> Administration (NCUA) of the U.S. federal government.	
2135: what is definition of psychotic	Psychosis
Despite this, " <b>psychosis</b> " is generally given to noticeable deficits in normal behavior (negative signs) and more commonly to diverse <b>types</b> of <b>hallucinations</b> or <b>delusional beliefs</b> (e.g. grandiosity, <b>delusions</b> of persecution).	
2136: Who is General Grievous of Star Wars	General Grievous
<b>General Grievous</b> is a fictional <b>character</b> in the <b>Star Wars</b> universe .	
2139: where scotty mccreery from	Scotty McCreery
He also released a <b>Christmas album</b> , <b>Christmas</b> with Scotty McCreery , which has been certified gold.	
2151: what is my resting heart rate at age 24	Heart rate
<b>Tachycardia</b> is defined as a <b>resting heart rate</b> above 100 <b>bpm</b> , though persistent <b>rest rates</b> between 80-100 <b>bpm</b> , mainly if they are present during sleep, may be signs of hyperthyroidism or anemia (see below).	
2152: what is the name for an old horse-drawn vehicle	Carriage
Working <b>vehicles</b> such as the (four-wheeled) wagon and (two-wheeled) cart share important parts of the <b>history</b> of the carriage, as is the fast (two-wheeled) chariot .	
2154: where was the tsunami in 2005	2004 Indian Ocean earthquake and tsunami
The resulting <b>tsunami</b> was given various <b>names</b> , including the 2004 Indian Ocean <b>tsunami</b> , South Asian <b>tsunami</b> , Indonesian <b>tsunami</b> , and the Boxing Day <b>tsunami</b> .	
2155: what is the source of geothermal energy	Geothermal energy
The <b>geothermal</b> gradient , which is the difference in temperature between the core of the <b>planet</b> and its <b>surface</b> , drives a continuous conduction of thermal <b>energy</b> in the form of heat from the core to the <b>surface</b> .	
2158: who was ho chi minh in vietnam war	Ho chi minh
Hồ Chí Minh (Northern Vietnamese pronunciation : , Southern Vietnamese pronunciation : ), 19 May 1890 – 2 September 1969), born Nguyễn Sinh Cung and also known as Nguyễn Tất Thành and Nguyễn Ái Quốc, was a Vietnamese communist <b>revolutionary</b> leader who was prime minister (1945–1955) and president (1945–1969) of the Democratic Republic of Vietnam (North Vietnam).	
2167: who was charged with murder after the massacre at My Lai	My Lai Massacre
The My Lai Massacre ( , , or ) was the Vietnam War mass murder of between 347 and 504 unarmed civilians in South Vietnam on March 16, 1968, by United States Army soldiers of "Charlie" Company of 1st Battalion , 20th Infantry Regiment , 11th Brigade of the American Division .	
2173: what is metal music about	Heavy metal music
Underground scenes produced an array of more extreme, aggressive styles: thrash <b>metal</b> broke into the mainstream with bands such as Metallica , Megadeth , Slayer , and Anthrax , while other styles of the most extreme subgenres of <b>metal</b> like death <b>metal</b> and black <b>metal</b> remain subcultural phenomena.	
2175: what is the title of Hobbes main work	Thomas Hobbes
Thomas <b>Hobbes</b> of Malmesbury (5 April 1588 – 4 December 1679), in some older texts Thomas <b>Hobbes</b> of Malmesbury, was an English philosopher, best known today for his <b>work</b> on political philosophy .	
2188: what is the cabin pressure of us airlines	Cabin pressurization
The <b>cabin pressure</b> is regulated by the outflow valve.	
2191: what is button-down shirt?	Dress shirt
A <b>shirt</b> , dress <b>shirt</b> , button-front, button-down <b>shirt</b> , or button-up <b>shirt</b> is a garment with a collar , a full-length opening at the front from the collar to the hem, and sleeves with cuffs .	
#2192: when is administrative assistant day	Administrative Professionals' Day
Administrative Professionals Day (also known as Secretaries Day or Admin Day) is an unofficial secular holiday observed in several countries to recognize the work of secretaries , administrative assistants , receptionists , and other administrative support professionals.	
2204: what is steam by valve corporation	Steam (software)
The <b>Steam</b> logo is a stylised left-side fly-crank and rod from the Walschaerts valve gear of a steam locomotive .	
2209: what is civil engineering aBOUT	Civil engineering
It is traditionally broken into several sub-disciplines including environmental <b>engineering</b> , geotechnical <b>engineering</b> , geophysics , geodesy , control <b>engineering</b> , structural <b>engineering</b> , biomechanics , nanotechnology , transportation <b>engineering</b> , earth science , atmospheric sciences , forensic <b>engineering</b> , municipal or urban <b>engineering</b> , water resources <b>engineering</b> , materials <b>engineering</b> , coastal <b>engineering</b> , surveying , and construction <b>engineering</b> .	
2217: who was john f kennedy up against	United States presidential election, 1960
This election is notable as being the first time in U.S. history that two sitting U.S. Senators ( <b>Kennedy</b> and <b>Johnson</b> ) were elected to <b>president</b> and <b>vice-president</b> , a phenomenon that has been repeated once, by Barack <b>Obama</b> and Joe Biden in 2008 (in both cases, the <b>president</b> was the younger, more junior senator).	
2221: what is the controlled substance act known as	Controlled Substances Act
Two federal agencies, the Drug Enforcement Administration and the Food and Drug Administration , determine which substances are added to or removed from the various schedules , though the statute passed by Congress created the initial listing, and Congress has sometimes scheduled other substances through legislation such as the Hillory J. Farias and Samantha Reid Date-Rape Prevention Act of 2000, which placed gamma hydroxybutyrate in Schedule I. Classification decisions are required to be made on criteria including potential for abuse (an undefined term), currently accepted medical use in treatment in the United States , and international treaties.	
2226: what is the si unit of pressure	Pressure
While pressure may be measured in any unit of force divided by any unit of area, the SI unit of pressure (the newton per square metre ) is called the pascal (Pa) after the seventeenth-century philosopher and scientist Blaise Pascal .	
2227: what is the great basin area	Great Basin
It is noted for both its arid conditions and its <b>Basin</b> and <b>range</b> topography that varies from the North American low point at Badwater <b>Basin</b> to the highest point of the contiguous United States , less than away at the summit of Mount Whitney .	
#2230: when was fdr elected as president	Franklin D. Roosevelt
A dominant leader of the Democratic Party and the only American president elected to more than two terms , he built a New Deal Coalition that realigned American politics after 1932, as his domestic policies defined American liberalism for the middle third of the 20th century.	
2231: where was ms-13 originally from	MS-13
In the U.S., the MS-13 has an especially heavy presence in Los Angeles County and the San Francisco Bay Area in Northern California; the Washington, D.C. metropolitan areas of Fairfax County, Virginia, Montgomery County, Maryland, and Prince George's County, Maryland; Long Island, New York; the Boston, Massachusetts area; Charlotte, North Carolina; and Houston, Texas.	
2234: where was the super bowl in 1991	Super Bowl XXV
The Bills and their explosive no-huddle offense were making their first Super Bowl appearance after finishing the regular season with a 13-3 record, and leading the league in total points scored with 428.	
Super Bowl XXV was an American football game between the American Football Conference (AFC) champion Buffalo Bills and the National Football Conference (NFC) champion New York Giants to decide the National Football League (NFL) champion for the 1990 season .	
In advancing to their second Super Bowl , the Giants also posted a 13-3 regular season record, but with a ball-control offense and a defense that allowed a league low 211 points.	
2237: who are the characters in 90210 in season 3	90210 (season 3)
Christmas-themed episode "Holiday Madness", hit <b>season</b> highs in all key <b>demos</b> with 2.1 in The CW's target <b>demo</b> of women 18–34, a 1.4 in <b>adults</b> 18–34 and a 1.1 in <b>adults</b> 18–49.	
#2238: when was queen elizabeth ii married	Elizabeth II
There have been times of personal sorrow for her which include the death of her father at 56, the assassination of Prince Philip's uncle, Lord Mountbatten , the breakdown of her children's marriages in 1992 (a year deemed her annus horribilis ), the death in 1997 of her former daughter-in-law, Diana, Princess of Wales , and the deaths of her mother and sister in 2002.	
2249: who is the singer westlife	Westlife
Westlife sold over 50 million records worldwide, a total that included studio albums, singles, video releases, and compilation albums.	
2250: what is vitamin a for	Vitamin A
The associated acid (retinoic acid), a metabolite that can be irreversibly synthesized from <b>vitamin</b> A, has only partial <b>vitamin</b> A activity, and does not function in the retina for the visual cycle.	
#2254: when is it memorial day	Memorial Day
Memorial Day is not to be confused with Veterans Day ; Memorial Day is a day of remembering the men and women who died while serving, while Veterans Day celebrates the service of all U.S. military veterans, living or dead.	
2259: what is the function of the liver	Liver
The <b>liver</b> is necessary for survival; there is currently no way to compensate for the absence of <b>liver function</b> in the long term , although new <b>liver</b> dialysis techniques can be used in the short term .	
2261: what is the formula for calcium nitrate	Calcium nitrate

A variety of related salts are known including <b>calcium ammonium nitrate</b> dehydrate and <b>calcium potassium nitrate</b> dehydrate.	
2262: where is dia de los muertos celebrated	<i>Day of the Dead</i>
The holiday has spread throughout the world: In Brazil , <b>Dia de Finados</b> is a public holiday that many Brazilians celebrate by visiting cemeteries and <b>churches</b> .	
2263: where is the 2011 mlb all star game location	<i>2011 Major League Baseball All-Star Game</i>
The <b>2011 Major League Baseball All-Star Game</b> was the 82nd in-season exhibition game between the All-Stars of the National League (NL) and the American League (AL); the <b>leagues</b> composing <b>Major League</b> baseball .	
2267: where is the island New Guinea?	<i>New Guinea</i>
The Germans annexed the northern coast of the <b>eastern</b> half of the <b>island</b> as German <b>New Guinea</b> in their pre-World War I effort to establish themselves as a <b>colonial</b> power, whilst the south <b>eastern</b> portion was reluctantly <b>claimed</b> by Britain.	
2272: what is in milk	<i>Milk</i>
New Zealand , the European Union 's 27 member states, Australia , and the <b>United States</b> are the world's largest exporters of <b>milk</b> and <b>milk</b> products.	
2273: What is motorcycle speedway racing	<i>Motorcycle speedway</i>
There are now both <b>domestic</b> and international <b>competitions</b> in a number of <b>countries</b> including the <b>Speedway World Cup</b> whilst the highest overall scoring individual in the <b>Speedway Grand Prix</b> events is pronounced the <b>world champion</b> .	
2289: where are the most concentration of jews living	<i>Jews</i>
Converts to Judaism , whose status as <b>Jews</b> within the <b>Jewish</b> ethnos is equal to those born into it, have been absorbed into the <b>Jewish</b> people throughout the millennia.	
2293: who was the first european in the americas	<i>European colonization of the Americas</i>
In 1497, sailing from the <b>north</b> on behalf of England , John Cabot landed on the <b>North American coast</b> , and a year later, <b>Columbus</b> ' third voyage reached the <b>South American coast</b> .	
2295: where is testosterone produced	<i>Testosterone</i>
On average, in adult human males, the plasma concentration of <b>testosterone</b> is about 7–8 times as great as the concentration in adult human females' plasma, but as the metabolic consumption of <b>testosterone</b> in males is greater, the daily <b>production</b> is about 20 times greater in men.	
2299: what is cu the element	<i>Copper</i>
Its <b>compounds</b> are commonly encountered as copper(II) <b>salts</b> , which often impart <b>blue</b> or green colors to minerals such as azurite and turquoise and have been widely used historically as pigments.	
2301: what was the parthenon used for	<i>Parthenon</i>
In the 5th century AD, the <b>Parthenon</b> was converted into a <b>Christian church</b> dedicated to the <b>Virgin Mary</b> .	
#2307: when was jamestown colonized	<i>Jamestown, Virginia</i>
Historic Jamestowne , the archaeological site on Jamestown Island, is a cooperative effort by Jamestown National Historic Site (part of Colonial National Historical Park ), and Preservation Virginia .	
2308: what was the first year of kentucky derby	<i>Kentucky Derby</i>
The Kentucky Derby () is a Grade I stakes race for three-year-old Thoroughbreds , held annually in Louisville, Kentucky , United States , on the <b>first</b> Saturday in May, capping the two-week-long Kentucky Derby Festival .	
2311: who is on the \$10 bill	<i>United States ten-dollar bill</i>
The source of the face on the \$10 bill is John Trumbull 's 1805 portrait of <b>Hamilton</b> that belongs to the portrait collection of New York City Hall .	
2318: who is norah jones parents	<i>Norah Jones</i>
<b>Norah Jones</b> (born Geetali <b>Norah Jones Shankar</b> ; March 30, 1979) is an <b>American</b> singer-songwriter, pianist, and actress.	
2325: what is Carbon 14 dating is a type of?	<i>Carbon-14</i>
There are three naturally occurring isotopes of <b>carbon</b> on Earth: 99% of the <b>carbon</b> is <b>carbon-12</b> , 1% is <b>carbon-13</b> , and <b>carbon-14</b> occurs in trace amounts, i.e., making up about 1 part per trillion (0.000000001%) of the <b>carbon</b> in the atmosphere.	
#2331: when was the battle at tombstone fought	<i>Gunfight at the O.K. Corral</i>
The gunfight, believed to have lasted only about thirty seconds, was fought between the outlaw Cowboys <b>Billy Clainmore</b> , <b>Ike</b> and <b>Billy Clinton</b> , and Tom and Frank McLauri , and the opposing town <b>Marshal</b> Virgil Earp and his brothers Assistant Town <b>Marshal</b> Morgan and temporary lawman Wyatt , aided by Doc <b>Holiday</b> designated as a temporary <b>marshal</b> by Virgil.	
2335: what is the population of algoma wi for 2010	<i>Algoma, Wisconsin</i>
See also the <b>Town of Algoma</b> in Winnebago County, Wisconsin .	
2349: who is Dr. JB Danquah	<i>J. B. Danquah</i>
During his political career, he was one of the primary opposition <b>leaders</b> to Ghanaian president and independence <b>leader</b> Kwame Nkrumah .	
2356: where is kj 52 from	<i>KJ52</i>
The "KJ" part of his <b>name</b> refers to his old rap alias, "King J. Mac," a <b>name</b> which he later described in one of his podcasts as "horribly cheesy."	
2360: what is the prognosis of stomach cancer	<i>Stomach cancer</i>
<b>Stomach cancer</b> , or <b>gastric cancer</b> , refers to <b>cancer</b> arising from any part of the <b>stomach</b> .	
2367: what is el mate	<i>Mate (beverage)</i>
"Tea-bag" type infusions of <b>mate</b> (mate cocido ) have been on the market in Argentina, Paraguay and Uruguay for many years under such trade names as " Taragüí Vitality" in Argentina, "Pajarito" and "Kurupú" in Paraguay, and in Brazil under the name " <b>Mate</b> Leão".	
2369: what is batista doing now	<i>Dave Batista</i>
David Michael "Dave" Bautista, Jr. (born January 18, 1969), is an American mixed martial artist , bodybuilder , actor, and former professional wrestler best known for his time in World Wrestling Entertainment competing under the ring name <b>Batista</b> .	
2371: where is the arctic circle located on the earth	<i>Arctic circle</i>
In fact, because of atmospheric refraction and mirages , and because the sun appears as a disk and not a point, part of the midnight sun may be seen on the night of the northern summer solstice up to about 50 °() south of the <b>Arctic Circle</b> ; similarly, on the day of the northern winter solstice , part of the sun may be seen up to about 50° north of the <b>Arctic Circle</b> .	
2375: what is the erb's heart	<i>Erb's point (cardiology)</i>
Heart valves are labeled with "B", "T", "A", and "P".First <b>heart</b> sound: caused by atrioventricular valves - Bicuspid/Mitral (B) and Tricuspid (T).	
Second <b>heart</b> sound caused by semilunar valves - Aortic (A) and Pulmonary/Pulmonic (P).	
Front of thorax , showing surface relations of bones , lungs (purple), pleura (blue), and <b>heart</b> (red outline).	
#2377: when was pokemon first started	<i>Pokémon</i>
The term <b>Pokémon</b> , in addition to referring to the <b>Pokémon franchise</b> itself, also collectively refers to the 649 fictional species that have made appearances in <b>Pokémon</b> media as of the release of the fifth generation titles <b>Pokémon Black</b> 2 and <b>White</b> 2 ; with the upcoming releases of <b>Pokémon</b> X and Y , 6 new <b>Pokémon</b> have been featured in promotions for the <b>games</b> .	
2379: where are colors on stoplight	<i>Traffic light</i>
Traffic lights alternate the right of way accorded to road users by displaying lights of a standard <b>color</b> (red, yellow/amber , and green) following a universal <b>color</b> code .	
2392: what is mince meat made of	<i>Mince meat</i>
Variants of <b>mincemeat</b> are found in Australia , Brittany , Canada , northern Europe , Ireland , South Africa , the <b>United Kingdom</b> and the <b>United States</b> .	
2399: what was the GE building in rockefeller plaza called before.	<i>GE Building</i>
The <b>GE Building</b> is an Art Deco skyscraper that forms the centerpiece of <b>Rockefeller</b> Center in Midtown Manhattan , New York City , USA .	
2400: who are the two senators of louisiana	<i>List of United States Senators from Louisiana</i>
Louisiana was admitted to the Union on April 30, 1812, and elects <b>senators</b> to Classes 2 and Class 3 .	
2408: who is st patty?	<i>Saint Patrick's Day</i>
Saint Patrick's Day or the Feast of Saint Patrick (, "the Day of the Festival of <b>Patrick</b> ") is a cultural and religious holiday celebrated on 17 March.	
2427: WHAT IS THE LENGTH OF A NAUTICAL MILE	<i>Nautical mile</i>
The <b>nautical mile</b> is nearly equal to a minute of latitude on a chart, so a distance measured with a chart divider can be roughly converted to <b>nautical miles</b> using the chart's latitude scale.	
2431: what is firewire used for	<i>IEEE 1394</i>
Apple first included <b>FireWire</b> in some of its 1999 models, and most Apple computers since the year 2000 have included <b>FireWire</b> ports, though, as of 2013, nothing beyond the 800 version (IEEE-1394b).	
#2432: when was sry born	<i>Stevie Ray Vaughan</i>
As the younger <b>brother</b> of <b>Jimmie Vaughan</b> , Vaughan started playing the guitar at age seven and formed several bands that <b>occasionally</b> performed in local nightclubs.	
2433: what is stent surgery	<i>Stent</i>
The term may also refer to a tube used to temporarily hold such a natural conduit open to allow access for <b>surgery</b> .	
2439: what is tofu made of	<i>Tofu</i>
There are many different varieties of <b>tofu</b> , including fresh <b>tofu</b> and <b>tofu</b> that has been processed in some way.	
2440: Where is Bubbles the Chimp now	<i>Bubbles (chimpanzee)</i>
In 2003 the public learned that, like all captive chimpanzees, <b>Bubbles</b> had matured into a large and aggressive adult <b>chimp</b> unsuitable as a companion animal.	
2441: what is muse's lead singer's name	<i>Muse (band)</i>
Muse are known for their energetic and extravagant live performances and their fusion of many music genres , including space <b>rock</b> , progressive <b>rock</b> , alternative <b>rock</b> , heavy metal , classical music and electronica .	
#2451: when was the internet started	<i>History of the Internet</i>
The Internet was commercialized in 1995 when <b>NSFNET</b> was decommissioned, removing the last restrictions on the use of the <b>Internet</b> to carry <b>commercial</b> traffic.	
2455: where is rolling rock brewed	<i>Rolling Rock</i>
Rolling Rock is a 4.6% abv pale lager launched in 1939 by the Latrobe <b>Brewing</b> Company .	
#2456: when was Pope Benedict XVI elected?	<i>Pope Benedict XVI</i>
Prior to becoming <b>pope</b> , he was "a major figure on the <b>Vatican</b> stage for a quarter of a century" as "one of the most revered, influential and controversial members of the College of <b>Cardinals</b> "; he had an influence "second to none when it came to setting church priorities and directions" as one of <b>Pope John Paul II</b> 's closest confidants.	
2458: Who was John Adam's children	<i>John Adams</i>
John Adams (October 30, 1735 ( O.S. October 19, 1735) – July 4, 1826) was the second president of the <b>United States</b> (1797–1801), having earlier served as the <b>first vice president</b> of the <b>United States</b> .	
2459: who was on the 10 dollar bill	<i>United States ten-dollar bill</i>
The \$10 bill is the only U.S. paper currency in circulation in which the portrait faces to the left (the \$100,000 <b>bill</b> featured a portrait of Woodrow Wilson facing to the left, but was used only for intra-government transactions).	
2460: what is water jet propulsion	<i>Jetboat</i>
Unlike these previous <b>waterjet</b> developments, such as Campini's and the Hanley Hydrojet, Hamilton had a specific need for a <b>propulsion</b> system to operate in very-shallow <b>water</b> , and the <b>waterjet</b> proved to be the ideal solution.	
2463: who is E from entourage	<i>Eric Murphy</i>
Murphy is a fictional <b>character</b> on the comedy-drama television series <b>Entourage</b> .	
2472: what is puerto rico currency	<i>Currencies of Puerto Rico</i>
The Banco Español de Puerto Rico was renamed <b>Bank of Porto Rico</b> and issued bills equivalent to the <b>United States dollar</b> , creating the <b>Puerto Rican dollar</b> .	
2486: what is endodontic dentistry	<i>Endodontics</i>
Endodontists perform a variety of procedures including <b>endodontic</b> therapy (commonly known as "root canal therapy"), <b>endodontic</b> retreatment , surgery, treating cracked teeth , and treating dental trauma .	
2491: who is the junior senator of nc	<i>Kay Hagan</i>
When Hagan defeated Republican <b>incumbent</b> Elizabeth Dole in the 2008 United States <b>Senate</b> election , she became the first woman to defeat an <b>incumbent</b> woman in a <b>Senate</b> election.	
2492: who is heisman trophy named after	<i>Heisman Trophy</i>
The <b>Heisman Memorial Trophy</b> Award (usually known colloquially as the <b>Heisman Trophy</b> or the Heisman), is awarded annually to the player deemed the most outstanding player in collegiate <b>football</b> .	
2498: what is sadomasochism	<i>Sadomasochism</i>
Similarly, <b>sexual</b> sadism within the context of mutual consent should not be mistaken for acts of <b>sexual</b> violence or aggression.	
2501: where are Giant Panda Bears found?	<i>Giant Panda</i>
"black and white cat-foot", also known as the <b>giant panda</b> to distinguish it from the unrelated <b>red panda</b> , is a <b>bear</b> native to central-western and south western China .	
2505: what is the format of the canadian citizenship test	<i>Canadian Citizenship Test</i>
The Canadian Citizenship Test is a <b>test</b> , administered by <b>Citizenship</b> and Immigration <b>Canada</b> (CIC), that is required for all applicants for <b>Canadian citizenship</b> who are aged between 18 and 54 and who meet the basic requirements for <b>citizenship</b> .	
2510: What is up with Kent Hovind	<i>Kent Hovind</i>
Since January 2007, <b>Hovind</b> has been serving a ten-year <b>prison</b> sentence after being convicted of 58 <b>federal</b> counts, including 12 tax offenses, one count of obstructing <b>federal</b> agents, and 45 counts of structuring <b>cash</b> transactions .	
2518: where is good morning america studio	<i>Times Square Studios</i>
The <b>studio</b> is best known as the production home of ABC News ' <b>Good Morning America</b> (GMA), a <b>morning</b> news and talk program and segments for GMA on ABC News Now.	
2523: what is primary medicine	<i>Primary care</i>
Such a professional can be a <b>primary care physician</b> , such as a general <b>practitioner</b> or family physician , or depending on the locality, health system organization, and patient's discretion, they may see a pharmacist , a physician assistant , a <b>nurse practitioner</b> , a <b>nurse</b> (such as in the United Kingdom), a clinical officer (such as in <b>parts</b> of Africa), or an Ayurvedic or other traditional <b>medicine</b> professional (such as in <b>parts</b> of Asia).	
#2530: when was the first super bowl	<i>Super Bowl</i>
Super Bowl XLV, played in 2011, became the most-watched American television program in history, drawing an average audience of 111 million viewers and taking over the spot held by the previous year's <b>Super Bowl</b> , which itself had taken over the #1 spot held for twenty-eight <b>years</b> by the final episode of M*A*S*H .	
2532: what is the official language of america?	<i>Languages of the United States</i>
The situation is quite varied at the state and territorial levels, with some states mirroring the federal policy of <b>adopting</b> no <b>official language</b> in a <b>de jure</b> capacity, others <b>adopting English</b> alone, others <b>officially adopting English</b> as well as local <b>languages</b> , and still others <b>adopting</b> a policy of <b>de facto</b> bilingualism.	
2534: where is shropshire, uk	<i>Shropshire</i>
The two Shropshire unitary areas (covering all of the ceremonial county), together with the authorities covering the ceremonial county of Staffordshire, comprise the "Shropshire and Staffordshire" NUTS 2 region .	
Shropshire is one of England's most rural and sparsely populated counties, with the population density of the Shropshire Council area being just 91/km2 (337/sq mi).	
2539: what is the social norm approach?	<i>Social Norms Approach</i>
Despite the fact that college drinking is at elevated levels, the perceived amount almost always exceeds actual behavior. The <b>social norms approach</b> has shown signs of countering misperceptions, however research on <b>resulting changes</b> in behavior resulting from <b>changed</b> perceptions varies between mixed to conclusively nonexistent.	
2545: what is an sd memory card reader	<i>Memory card reader</i>
There are three categories of <b>card readers</b> sorted by the type and quantity of the <b>card</b> slots: single <b>card reader</b> (e.g. 1x SD-only), multi <b>card reader</b> (e.g. 9-in-1) and series <b>card reader</b> (e.g. 4x SD only).	
2547: What was Captain Ahab's Ship in the novel "Moby Dick"	<i>Moby-Dick</i>
Ishmael soon learns that <b>Ahab</b> has one purpose on this <b>voyage</b> : to seek out <b>Moby Dick</b> , a ferocious, enigmatic white sperm whale .	
#2563: when is halley's comet next	<i>Halley's Comet</i>
Halley's Comet or Comet Halley ( ), officially designated 1P/Halley, is the best-known of the short-period comets and is visible from Earth every 75–76 years.	
2566: what is soy made from	<i>Soybean</i>
Fat-free (defatted) soybean meal is a significant and cheap source of <b>protein</b> for animal feeds and many prepackaged meals ; soy vegetable oil is another <b>product</b> of processing the soybean crop.	
2569: what is linkedin used for	<i>LinkedIn</i>
LinkedIn reports more than 200 million acquired <b>users</b> in more than 200 countries and territories.	
2596: what is captcha code	<i>CAPTCHA</i>
Although most <b>CAPTCHAs</b> are letter pictures randomly generated, many of them have become difficult even for a human to read , so picture <b>CAPTCHAs</b> were created in which a human is shown a simple test to show a picture of a certain animal (given few animal pictures), which is simple for a human being to process, and therefore easy to pick, while a bot cannot process and solve the question because although it can analyze the picture, it cannot easily guess the animal.	
2611: who is in the group trinity 5 /	<i>Trin-i-tee 5</i>
In 1998, the <b>group</b> released their first album entitled <b>Trin-i-tee</b> 5 which debut at #1 on the Top Gospel charts .	
#2619: when is international men's day	<i>International Men's Day</i>
International Men's Day is celebrated in over 60 countries, including Trinidad and Tobago, Jamaica, Australia, India, China, United States, Romania, Singapore, Malta, United Kingdom, South Africa, Tanzania, <b>Zimbabwe</b> , Botswana, Seychelles, Burundi, Hungary, Ireland, Isle of Man, Ghana, Canada, Denmark, Norway, Austria, <b>Bosnia</b> and Herzegovina, Ukraine, France, Italy, Pakistan, Cuba, Antigua and Barbuda, St. Kitts and Nevis, St. Lucia , Grenada and Cayman Islands , on 19 November, and global support for the celebration is broad.	
#2635: when was the patriot act enacted	<i>Patriot Act</i>

On May 26, 2011, President Barack Obama used an Autopen to sign the **PATRIOT Sunsets Extension Act of 2011**, a four-year extension of three key provisions in the USA **PATRIOT Act** while he was in France: roving wiretaps, searches of business records (the "library records provision"), and conducting surveillance of "lone wolves"—individuals suspected of terrorist-related activities not linked to terrorist groups.

#2637: when was the tacoma bridge collapse?

Tacoma Narrows Bridge

Historically, the name "**Tacoma Narrows Bridge**" has applied to the original **bridge** nicknamed "Galloping Gertie", which opened in July 1940 but collapsed because of aerelastic flutter four months later, as well as the replacement of the original **bridge** which opened in 1950 and still stands today as the **westbound** lanes of the present-day twin **bridge** complex.

2640: what is brown flax

Flax

**Flax** (also known as common **flax** or linseed) (binomial name: Linum usitatissimum) is a member of the genus Linum in the family Linaceae.

#2669: when was the lady gaga judas song released

Judas (song)

Gaga performed "Judas" on a number of television shows, including The **Graham Norton Show**, Saturday Night Live, Good Morning Americas "Summer Concert Series", the French X Factor as well as on The Ellen DeGeneres Show.

2670: what is mpeg4 avc or sp/asp

H.264/MPEG-4 AVC

The **ITU-T H.264 standard** and the ISO/IEC MPEG-4 AVC standard (formally, ISO/IEC 14496-10 – MPEG-4 Part 10, Advanced Video Coding) are jointly maintained so that they have identical technical content.

2685: what is the definition of a map

Map

Although most commonly used to depict geography, **maps** may represent any space, real or imagined, without regard to context or scale; e.g. brain **mapping**, DNA **mapping** and extraterrestrial **mapping**.

2688: who is on blank 182 album cover

Enema of the State

Produced by Jerry **Finn**, the album was released on June 1, 1999 in the United States on MCA Records.

2691: what is the name of the late Major League old professor

Casey Stengel

After his **major league** career began, he acquired the nickname "**Casey**", which originally came from the initials of his hometown ("K. C."), which evolved into "**Casey**", influenced by the wide popularity of the poem **Casey at the Bat**.

#2693: when was the constitution written

Constitution

The **Constitution** of India is the longest **written constitution** of any sovereign **country** in the world, containing 448 articles, 12 schedules and 100 amendments, with 117,369 words in its English language version, while the **United States Constitution** is the shortest **written constitution**, at 7 articles and 27 amendments.

2697: what is the use of a sales invoice?

Invoice

In English, the context of the term **invoice** is usually used to clarify its meaning, such as "We sent them an **invoice**" (they owe us money) or "We received an **invoice** from them" (we owe them money).

2698: what was the actress who played the pink power ranger

Kimberly Hart

Kimberly is best remembered as the first **Pink Ranger** and first **Pink Ninja Ranger** from the first entry of the franchise Mighty Morphin **Power Rangers**.

2700: what is squash the sport

Squash (sport)

The **game** was formerly called **squash racquets**, a reference to the "**squashable**" soft ball used in the **game** (compared with the fatter ball used in its parent **game racquets** or rackets).

2702: where is the world cup in 2010

2010 FIFA World Cup

The **2010 FIFA World Cup** was the 19th **FIFA World Cup**, the world championship for men's national association football teams.

2707: what is hosting a website

Web hosting service

The most basic is web page and small-scale **file hosting**, where **files** can be uploaded via File Transfer Protocol (FTP) or Web interface.

#2708: when was the first nfl madden game released

Madden NFL

**Madden NFL** (known as John **Madden Football** before 1993) is an American **football** video **game series** developed by Electronic Arts Tiburon for EA Sports.

2738: what is the purpose of child support?

Child support

The 1992 United Nations Convention on the Rights of the **Child**, a binding convention signed by every member nation of the United Nations and formally ratified by all but Somalia and the United States, declares that the upbringing and development of **children** and a standard of **living** adequate for the **children's** development is a common responsibility of both parents and a fundamental human right for **children**, and asserts that the primary responsibility to provide such for the **children** rests with their parents.

2741: where was martin luther king shot?

Martin Luther King, Jr.

**Martin Luther King, Jr.** was established as a U.S. federal **holiday** in 1986.

2747: who is the mayor of chicago 2011

Chicago mayoral election, 2011

The city of **Chicago**, Illinois held a nonpartisan **mayoral** election on Tuesday, February 22, 2011.

#2749: when is the wv state fair

State Fair of West Virginia

The **State Fair of West Virginia** is an annual state fair for **West Virginia**.

#2755: when was the state of utah established

Utah

Approximately 80% of **Utah's** 2,817,222 people live along the Wasatch Front , centering on Salt **Lake** City , leaving vast expanses of the **state** nearly uninhabited and making the **population** the sixth most urbanized in the U.S. **Utah** is bordered by **Colorado** on the east, Wyoming on the northeast, Idaho on the north, Arizona on the south, and Nevada on the west.

#2762: when is world war hulk movie come

World War Hulk

The **series** consists of five main issues titled **World War Hulk**, with Greg Pak as writer and John Romita, Jr. as penciller, and three other limited **series** : **World War Hulk**: Front Line, **World War Hulk**: Gamma Corps, and **World War Hulk**: X-Men .

2763: WHERE WAS JOHN WAYNE BORN

John Wayne

Marion Mitchell Morrison (born Marion **Robert** Morrison; May 26, 1907 – June 11, 1979), better known by his stage name **John Wayne**, was an American film actor, director and producer.

2765: what is the difference between alpha lipoic acid and lipoic acid

Lipoic acid

The carbon atom at C6 is chiral and the molecule exists as two enantiomers (R)-(+)-lipoic acid (RLA) and (S)-(−)-lipoic acid (SLA) and as a racemic mixture (R/S)-lipoic acid (R/S-LA).

2767: what is the organic layer in an amoled screen

OLED

An OLED (organic light-emitting diode) is a light-emitting diode (LED) in which the emissive electroluminescent layer is a film of organic compound which emits light in response to an electric current.

2783: what is bourbon made of

Bourbon whisky

The name of the spirit derives from its historical association with an area known as Old **Bourbon**, around what is now **Bourbon** County, Kentucky (which, in turn, was named after the French House of **Bourbon** royal family).

2802: where is La Palma africa

La Palma

It was as a result of his visit to **La Palma** and Tenerife where he visited the Las Cañadas and Taburiente calderas, that the Spanish word for cauldron - " Caldera " - was introduced into the English language geological vocabulary.

2810: where was hillary clinton born

Hillary Rodham Clinton

That **election** marked the first time an American **First Lady** had run for public office; **Clinton** was also the first female senator to represent the state.

#2819: When was 27th amendment proposed

Twenty-seventh Amendment to the United States Constitution

The Twenty-seventh **Amendment** (Amendment XXVII) prohibits any law that increases or decreases the salary of members of the Congress from taking effect until the start of the next set of terms of office for Representatives.

2821: what is singapore's currency

Singapore dollar

The Monetary Authority of **Singapore** and the **Brunei Currency** and Monetary Board still maintain the historic exchangeability of their two currencies, the **Singaporean** dollar and the **Brunei** dollar, respectively.

2828: where was paul revere born

Paul Revere

Following the **war**, **Revere** returned to his **silversmith** trade and used the profits from his expanding business to finance his work in iron **casting**, bronze bell and cannon **casting**, and the forging of copper bolts and spikes.

2841: what is the type of democracy in which all citizens have the right to make government decisions

Democracy

One form of **democracy** is **direct democracy**, in which all eligible **citizens** have **direct** and **active** participation in the decision making of the government.

2847: what is ratchet from ratchet and clank

Ratchet & Clank

The franchise was created and developed by Insomniac Games and published by Sony Computer Entertainment for many different **PlayStation** consoles, such as **PlayStation** 2 and **PlayStation** 3 with the exclusion of Size Matters and Secret Agent **Clank**, which were developed by High Impact Games for the **PlayStation Portable**.

2848: when was the first real roller coaster

roller coaster

Some **roller coasters**, notably wild mouse **roller coasters**, run with single cars.

Will and testament

In the strictest sense, a "will" has historically been limited to **real property** while "testament" applies only to dispositions of **personal property** (thus giving rise to the popular title of the document as "Last Will and Testament"), though this distinction is seldom observed today.

Bracketing

Autobracketing is automatic **bracketing** by using a setting on the camera to take several bracketed shots (in contrast to the **photographer** altering the settings by hand between each shot).

Defragmentation

What the disk defragmenter tool

Disk drives

Defragmentation is advantageous and relevant to file systems on electromechanical **disk** drives.

Service-level agreement

2847: what is level of agreement mean

Level of agreement

As an example, **internet** service providers will commonly include service **level agreements** within the **terms** of their **contracts** with customers to define the level(s) of service being sold in plain language **terms**.

List of Harry Potter cast members

2883: who is the actor that plays harry potter

Harry Potter movies.

2884: what is the concept of "wellness" ?

Wellness (alternative medicine)

**Wellness** grew as a popular **concept** starting in the 19th century, just as the middle class began emerging in the industrialized world, and a time when a newly prosperous public had the time and the resources to pursue **wellness** and other forms of self-improvement.

American football positions

2899: what is the defensive line in football called

Defensive line

This has resulted in the development of three "platoons" of players, the offense (the team with the ball, who is trying to score), the defense (the team trying to prevent the other team from scoring, and to take the ball from them), and the special teams (who play in kicking situations).

Moons of Saturn

2903: what is the measurements of saturn's moons

Moons

They include the seven major satellites, four small **moons** which exist in a **trojan** orbit with larger **moons**, two mutually co-orbital **moons** and two **moons** which act as shepherds of **Saturn's F Ring**.

Bad Girls Club

2907: who are the girls from the bad girls club?

Bad Girls Club

If a "bad girl" breaks a rule, she is evicted from the show and, if it is early in the **season**, replaced by a new "bad girl".

Bernard Madoff

2916: what was bernie madoff selling

Bernard Madoff

Madoff founded the Wall Street firm **Bernard L. Madoff** Investment Securities LLC in 1960, and was its chairman until his arrest on December 11, 2008.

Invention of radio

2923: who is inventor of the radio

Others, notably Guglielmo Marconi , were concerned with practical improvements and the commercial application of **radio** to wireless telegraphy .

Scurvy

2925: what is scurvy disease

Other eponyms for **scurvy** include Moeller's disease and Cheadle's disease.

OpenID

2931: what is google openid

OpenID

An extension to the standard (the OpenID Attribute Exchange) facilitates the transfer of user attributes, such as name and gender, from the OpenID identity provider to the relying party (each relying party may request a different set of attributes, depending on its requirements).

Atherosclerosis

Arteriosclerosis is a general term describing any hardening (and loss of elasticity) of medium or large **arteries** (from the Greek arteria , meaning artery, and , meaning hardening); arteriosclerosis is any hardening (and loss of elasticity) of **arterioles** (small arteries); **atherosclerosis** is a hardening of an **artery** specifically due to an atherosomatic plaque.

USS John F. Kennedy (CV-67)

The ship is **named** after the 35th President of the United States, John F. Kennedy , and is nicknamed "Big John."

Douglas MacArthur

General of the Army Douglas MacArthur (26 January 1880–15 April 1964) was an American **general** and field marshal of the Philippine **Army** who was Chief of Staff of the **United States Army** during the 1930s and played a prominent **role** in the Pacific theater during World War II.

Mrs. Claus

2946: what is santa's wife's name

Mrs. Claus

Mrs. Claus is the wife of Santa Claus , the Christmas gift-bringer in North American and European Christmas tradition .

Saddle Creek

2949: when was Saddle Creek founded

Saddle Creek Records

Saddle Creek first appeared in print on a show flyer, offering to "Spend an evening with **Saddle Creek**" (later to be the title of the label's DVD.)

Psych

2958: Who is the highest scoring NBA player

List of National Basketball Association season scoring leaders

The National Basketball Association 's (NBA) **scoring** title is awarded to the **player** with the **highest** points per game average in a given season.

Direct Marketing

2961: what is direct marketing channel

Direct Marketing

Direct marketing relies on being able to **address** the members of a target **market** .

Kidney

2974: what is the name of the family who own the biltmore estates in nc

Biltmore Estate

Biltmore Estate is a large private **estate** and tourist attraction in Asheville , North Carolina .

Kidney

2979: where are the kidneys in your body

Kidneys

Common clinical conditions involving the **kidney** include the nephritic and nephrotic syndromes , renal cysts , acute **kidney** injury , **chronic kidney** disease , urinary tract infection , nephrolithiasis , and urinary tract obstruction .

White chocolate

2983: what is white chocolate made of

White chocolate

USA renewed the series for a seventh **season** , to include 16 episodes, on January 10, 2012 , and again for an eighth **season** , to include eight episodes, on December 19, 2012 .

Psych

2990: what is the main component of vaccines

Vaccine

Vaccines may be prophylactic (example: to prevent or ameliorate the effects of a future infection by any natural or "wild" pathogen ), or therapeutic (e.g. **vaccines** against cancer are also being investigated; see cancer **vaccine** ).

Alps

3005: what is the alpine mountain systems

Alps

The Alps are one of the great **mountain range** systems of Europe stretching approximately across eight Alpine countries from Austria and Slovenia in the **east** , Switzerland , Liechtenstein , Germany , France to the **west** and Italy and Monaco to the **south** .

Mary Matalin

3008: who is mary matalin married to

Mary Matalin

She appears in the award-winning documentary film Boogie Man: The Lee Atwater Story and also played herself, opposite her husband, James Carville , John Slattery , and Mary McCormack in the short lived HBO series K Street .

Basque language

3010: where is basque spoken

Basque

These **provinces** and many **areas** of Navarre are heavily populated by ethnic Basques, but Basque was, at least until the 1990s, not spoken as a native language in most of Alava, western parts of Biscay and central and southern **areas** of Navarre , either because it had been replaced by **Spanish** along the centuries, in some **areas**, or because it had never been spoken there, in other **areas**.

Dear John (2010 film)

3020: where is dear john filmed

Dear John is a 2010 American romantic drama - war film starring Amanda Seyfried and Channing Tatum .

Melissa & Joey

3034: what is melissa and joey about

Melissa & Joey

Melissa & Joey is an ABC Family original television series starring Melissa Joan Hart and Joey Lawrence .

Household income in the United States

3039: what is the average american income

Income

## 5 Approximate word clustering in selected Germanic languages

In this section, we present the approximate clustering algorithms for three representative Germanic languages: Danish, Dutch and German. All these three languages underwent orthographic reforms at certain points in the 20th century. We devise algorithms that accommodate to both the new and old orthographies.

It is worth mentioning that the three aforementioned Germanic languages contain some morphological features not found in English (an atypical Germanic language):

- Both Danish and German have highly inflected nouns and adjectives, and the same is true for the written standard of pre-1946 Dutch.
- Both Dutch and German verbs conjugate in the persons, and both these languages have verbs with separable prefixes.
- All these three languages allow/require a compound to be written as a single word without spaces or hyphens.

Up to minor changes, these descriptions for Dutch extend to Afrikaans (a variety of Dutch spoken in South Africa), while the features of Danish are also found in Norwegian and Swedish (two Scandinavian languages that share the same ancestor with Danish). Therefore, the algorithms presented here can be generalized to these languages with minimal modifications. We note, however, that Old Norse (the ancestral language of modern Scandinavian languages) and its conservative modern descendants (Faroese and Icelandic) have much higher morphological complexity, not yet fully captured by our current algorithms.

Like English, the Germanic strong verbs conjugate in a wide variety of “irregular” patterns. Our algorithms address the majority of these strong verbs, but do not cover all the irregularities in verb conjugation.

We will not impose the “heuristic capitalisation” procedure (Definition 4.1) on any of the non-English languages treated in this document. Accordingly, all the words in a vocabulary list must be reduced to lowercase form before they are submitted to any one of our modified Porter stemming algorithm for languages other than English.

### 5.1 Modified Porter stemming algorithm for Danish

Some adverbs, articles, prepositions and pronouns occur frequently in Danish texts, while conveying very little information in their own right. They are regarded as stop words.

**Definition 5.1** (Danish stop words). If a word belongs to the following list<sup>49</sup>:

*ad, af, al, aldrig, alle, allerede, alligevel, alt, altfor, altid, anden, andet, andre, andres, at, atter, bag, bagved, bare, blandt, blev, blevet, bliv, blive, bliver, blot, bort, borte, burde, burdet, bør, både, baade, da, de, dem, den, denne, dennelunde, dens, der, deres, derfor, derfra, derhenne, derpå, derpaa, dersom, dertil, des, desto, det, dets, dette, dig, din, dine, disse, dit, dog, du, efter, efterhånden, efterhaanden, eftersom, egen, eget, egne, eller, en, én, end, endnu, endvidere, ene, eneste, enhver, enhvert, ens, enten, er, et, ethvert, flere, flest, fleste, for, foran, forbi, fordi, forudsat, fra, før, førend, ganske, gennem, gjorde, gjort, gör, gøre, ham, han, hans, har, havde, have, heller, hellere, hen, hende, hendes, henne, her, herfra, herhenne, hermed, hertil, hinanden, hinsides, hos, hun, hvad, hvadenten, hvem, hver, hverken, hvert, hvilke, hvilket, hvis, hvor, hvordan, hvorfor, hvorhen, hvorhenne, hvormed, hvorvidt, i, iblandt, idet, iflad, ifølge, igen, igennem, ihvorvel, ikke, imellem, imens, imod, ind, inde, inden, indtil, ingen, intet, ja, jeg, jer, jeres, jo, kan, kun, kunde, kunne, kunnet, langs, lidt, ligesom, ligeså, ligesa, man, mange, mangen, med, medens, megen, meget, mellem, men, mens, mere, mest, mig, min, mine, mit, mod, må, maa, måske, maaske, måtte, maatte, måttet, ned, nede, nej, nemlig, nogen, noget, nogle, nok, nu, næmlig, nær, nærværd, næsten, når, naar, närsomhelst, naarsomhelst, ofte, oftere, oftest, ofteste, og, også, ogsaa, om, omend, omkring, omme, op, oppe, os, osv, over, overfor, ovre, på, paa, selv, selvom, siden, sidst, sidste, sig, sikken, sin, sine, sit, skal, skulde, skulle, snart, som, sommetider, straks, så, saa, sådan, saadan, sådanne, saadanne, sådant, saadant, såfremt, saafremt, således, saaledes, såvel, saavel, thi, til, tilbage, tit, trods, uagtet, ud, ude, uden, under, undertiden, undtagen, var, vi, vil, vilde, ville, villet, vor, vore, vores, vort, vær, være, været,*

then we consider it a Danish stop word. All the Danish stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list. □

<sup>49</sup>Our list of Danish stop words is based on [snowball.tartarus.org/algorithms/danish/stop.txt](http://snowball.tartarus.org/algorithms/danish/stop.txt), together with forms derived from substitutions  $\ddot{a} \rightarrow aa$  (accommodating to pre-1948 Danish orthography).

### 5.1.1 Effective spelling and essential root

**Algorithm 5.2** (Danish effective spelling). Set  $\mathbf{V} = (a|e|i|o|u)$ . For a Danish word  $\hat{\sigma}$ , its effective spelling  $\text{EffSpell}(\hat{\sigma})$  is constructed in the following steps:

- (1) Convert to lowercase, and replace<sup>50</sup>

dug	dvikl	elizaX~	fattig	formue	forster	fætter	hoved	jone	klæd
δuy	dbikl	ελιζ	πριργ	φορτ	φορστρ	νεψφε	ηεαδ	γονε	κλαδ
kone	langsom	læk~	mund	mål	ordn~	peg	rolig	sidde~	sider
koñe	σλωωτ	δλак	μυñδ	σμαλ	ωρδн	ποιντ	ρωλиг	σιτε	παյер
sl(og å)	soldat	stolt	sukk	sum~	sund	synd	sød	uværdig~	miss
ζληα	σολδατ	σтолтζ	σукк	σум	γσυñδ	τσуы	σωіт	ωοθλσ	zμіσ
									s(a i)d
									σиτе

- (2) Replace

'	χ̄(u)fuldX	X $\epsilon$ (g l v)ent	ee	gods~	guv	hær	hår	janu~	kon~	kær~	land	lind	misχ̄(t)~
∅	χ̄	Xent	eē	γοδσ	γυν	ηοστ	χαρ	јну	con	κær	land	lind	μσχ̄
mæssig	part	præ~	ros~	sol~	sonet~	ss	uge	varm~	χværdigX	yv	ø~		
∅	part	πρæ	ροσ	σολ	σονετ	β	ωεκ	ναρμ	χ̄	οј	ζxo		
(liden mindre mindst mindste smaa små)	ro(∅ en)	seen(∅ de)	~χ̄(a e i o u)tes	~isme(∅ en)	~ist(∅ en er erne)(∅ s)								
lille	ρωλиг	σσωω	χ̄es	e	e								

- (3) Replace

f(a æ)d	(here herre)~	(ond)værre værst)~	X $\epsilon$ (tagg tigd tigg)~
fað	zherre	onxved	bX
X $\epsilon$ lege(∅ dr t)~	ald~	b(o ø)g~	bild
zzX	gamle	kbog	bill~
	βild	tbill	bland~
	mblanc	ebord	bord~
		bring	bragt~
		βu	bu
		fbode	bøde~
		ldam	dam~
		danzts	dans~
		δich	dig~
		sdum	dum~
dyd~	d(å aa)r~	fald~	far~
vdyd	δcar	efald	rfor
		fruct	frugt~
		lst	først~
		zxfor	før~
		γlim	glim~
		γod	god~
		kgrad	grad~
		rgrød	grød~
		hgud	gud~
		hgast	gæst~
		ghils	hils~
hind~	hum	hund~	hus
χind	huu	dxund	hus
		χor	xiv
		ckald	kend~
		kend	kold~
		comp	komp~
		wkone	kone~
		gkun	kun~
		dkor	kør~
		λav	lav~
		λok	luk~
		xlys	lys~
		rlob	løb~
mind~	mini~	mørk~	m(å aa)d~
μind	μini	dmork	ny(∅ e ere est et t)~
		wmad	nå~
			rede~
			r(y ø)g~
			sagχ̄(d t)
			se~
			sir
skænd~	smal~	smul~	som~
zskand	zmal	xmul	svær
		zom	syg~
		dfsvar	søg
		σyg	søek
		ζon	søn
		zil	til~
			tro(∅ ede en ens et r)s)~
			tung~
ty~	undre~	v(å aa)gn~	yngre~
θy	wundre	wvagn	ung
		oabp	oabp
		jår	jår
		ga	ga
		γod	bed(re st ste)
			f(å aa)(∅ et r)
			fik
mr	mrs	mød(∅ t te tes)	møde(∅ r rne rnes rs s t ts)
zηeppe	zφip	μod	μod
			sag
			sagten
			se(∅ r s t)
s(å aa)(∅ s)	ti	ved	æld(re st ste)
σσaw	10ti	vid	gammel
			~st(od (å aa)(∅ elig(∅ e t) et r))
			stand
~X $\epsilon$ (br m)(oder or)(∅ en ens s)	Xoð	~X $\epsilon$ (br m)ødre(∅ ne nes s)	~far(∅ en ens s)
		Xoð	fað
			ligger

- (4) Do  $st\chī(o|ø) \sim \rightarrow s\chīt\chī$ ,  $vi \rightarrow vzi$ ,  $y \rightarrow u$ ,  $æ \rightarrow a$ ,  $ø \rightarrow o$ ,  $(å|aa) \rightarrow a$  and undouble characters (i.e. do  $mm \rightarrow m$  etc.);

- (5) If the result so far is an empty string, stop here. Otherwise, break down the result into  $\hat{\sigma}_1\hat{\sigma}_2$ , where  $\ell(\hat{\sigma}_1) = 1$ , and work on  $\hat{\sigma}_2$  in four sequential steps:

(5.1) Do  $el \rightarrow l$ ,  $ē \rightarrow e$ ,  $riv \rightarrow irv$ ;

(5.2) Do  $\chī\mathbf{V}re \rightarrow \chīer$ ;

<sup>50</sup>To avoid confusion of Greek nu with Latin vee, we write the former as ψ in substitution rules.

- (5.3) Do  $\text{inter} \rightarrow \text{inter}$ ;
- (5.4) Call the result so far as  $\hat{\sigma}'_2$ . Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ .

**Definition 5.3** (Danish protected range). Set  $\mathbf{V} = (a|e|i|o|u)$ ,  $\mathbf{V}_m^* = (a|i|o|u)_m$  and  $\mathbf{C}_{m_0} = \overline{(a|e|i|o|u)}_{m_0}$ . Let  $\hat{\sigma}$  be the effective spelling of a Danish word, its protected range  $\text{ProtRg}(\hat{\sigma}) = \max\{\lambda_1(\hat{\sigma}), \lambda_2(\hat{\sigma})\}$  is determined by two non-negative integers  $\lambda_1(\hat{\sigma})$  and  $\lambda_2(\hat{\sigma})$  specified through the following procedures:

- Look for the string pattern  $(\emptyset|ab|af|be|in|\pi\rho a)\mathbf{C}_{m_0}(e|\mathbf{V}_m^*)\overline{(a|e|i|o|r|u)}_{[0,1]}\sim$  in the string  $\hat{\sigma}$ :<sup>51</sup>
  - If the string pattern above is found, the last position occupied by such a string defines  $\lambda_1(\hat{\sigma})$ ; otherwise, set  $\lambda_1(\hat{\sigma}) = 0$ ;
  - Do  $\sim(\emptyset|t)(\emptyset|s) \rightarrow \emptyset$  on  $\hat{\sigma}$ , and call the result  $\hat{\sigma}'$ ;
  - Look for the pattern  $(a|o|s|\beta|t|u)$  in the string  $\hat{\sigma}'$ ;
  - If a letter in the pattern above is found, the last position occupied by such a letter defines  $\lambda_2(\hat{\sigma})$ ; otherwise, set  $\lambda_2(\hat{\sigma}) = 0$ .
- 

**Algorithm 5.4** (Danish essential root). Let  $\hat{\sigma}$  be the effective spelling of a Danish word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:

(1) Break down  $\hat{\sigma} = \hat{\sigma}_1\hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}^{[\text{ProtRg}(\hat{\sigma})]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma})$  is equal to the protected range of  $\hat{\sigma}$ .

(2) On  $\hat{\sigma}_2$ , perform the following substitutions in a sequel:<sup>52</sup>

- (2.1) Do  $\sim e(\emptyset|s) \rightarrow \emptyset$ ;
- (2.2) Do  $\sim e(d|e|n|r|s|t)_{m_0} \rightarrow \emptyset$ .

The result after these two steps of operations is called  $\hat{\sigma}'_2$ .

(3) Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ , before doing  $\sim s(\emptyset|t) \rightarrow \emptyset$ ;

(4) Do  $\sim \hat{x}^\epsilon(b|d|g|k)t \rightarrow \hat{x}$ ,  $\sim er \rightarrow e$ ,  $\sim gd \rightarrow g$ .

*Example 5.4.1.* In many cases, the essential root of a Danish word is already a very close approximation to its etymological stem. For instance, we may consider the inflected forms of two Danish nouns *datter* “daughter” and *mand* “man”, as tabulated below:

$\hat{\sigma}$	$\text{EffSpell}(\hat{\sigma})$	$\text{EssRoot}(\text{EffSpell}(\hat{\sigma}))$	$\hat{\sigma}$	$\text{EffSpell}(\hat{\sigma})$	$\text{EssRoot}(\text{EffSpell}(\hat{\sigma}))$
<i>datter</i>	<i>dater</i>	<i>dat</i>	<i>mand</i>	<i>mand</i>	<i>mand</i>
<i>datteren</i>	<i>dateren</i>	<i>dat</i>	<i>manden</i>	<i>manden</i>	<i>mand</i>
<i>datterens</i>	<i>daterens</i>	<i>dat</i>	<i>mandens</i>	<i>mandens</i>	<i>mand</i>
<i>datters</i>	<i>daters</i>	<i>dat</i>	<i>mands</i>	<i>mands</i>	<i>mand</i>
<i>døtre</i>	<i>doter</i>	<i>dot</i>	<i>mænd</i>	<i>mand</i>	<i>mand</i>
<i>døtrene</i>	<i>doterne</i>	<i>dot</i>	<i>mændene</i>	<i>mandene</i>	<i>mand</i>
<i>døtrenes</i>	<i>doternes</i>	<i>dot</i>	<i>mændenes</i>	<i>mandenes</i>	<i>mand</i>
<i>døtres</i>	<i>doters</i>	<i>dot</i>	<i>mænds</i>	<i>mands</i>	<i>mand</i>

The declensions of *datter* “daughter” exhibit vowel alternations in their essential roots. This phenomenon, together with the vowel alternations in Danish irregular verbs, will be handled in §5.1.2.

### 5.1.2 Admissible mutation and approximate clustering

The Danish vowel blotting mechanism is slightly different from the English counterpart (Algorithm 4.9).

**Algorithm 5.5** (Danish vowel blotting). Set  $\mathbf{V} = (a|e|i|o|u)$  and  $\mathbf{C}_{m_0} = \overline{(a|e|i|o|u)}_{m_0}$ . For a token string  $\hat{\sigma}$ , its blotted form  $\text{BlotV}_1(\hat{\sigma})$  is constructed as follows:

- If the string pattern  $(\emptyset|ab|af|be)\mathbf{C}_{m_0}\mathbf{V}\sim$  can be found in the string  $\hat{\sigma}$ , then the last position occupied by such a pattern is replaced by the letter “a”.
- Otherwise, leave the string  $\hat{\sigma}$  intact.

<sup>51</sup>See Definition 3.3 for the multiplicity notation  $(\hat{\sigma}_1|\cdots|\hat{\sigma}_n)_{[0,1]}$ .

<sup>52</sup>In other words, the core algorithm for essential root extraction runs as follows: keep the last “strong” vowel *a*, *i*, *o* or *u* in non-final position, plus one subsequent letter; delete final *a*; erase the final appearance of *e* and all the letters thereafter.

*Example 5.5.1.* As a continuation of Example 5.4.1, we point out that  $\text{BlotV}_1(\text{EssRoot}(\text{EffSpell}(\hat{\sigma})))$  evaluates to *dat* (resp. *mand*) for all the declensions of *datter* (resp. *mand*). If we consider  $\text{BlotV}_1(\text{EffSpell}(\hat{\sigma}))$  instead, the output is identical to  $\text{EffSpell}(\hat{\sigma})$  for every declined form of *mand*, but involves a substitution  $o \rightarrow a$  for the case of *datter*.

*Example 5.5.2.* A few more examples involving inflected forms of the Danish irregular verb *finde* “find”:

$\hat{\sigma}$	<i>fandt</i>	<i>find</i>	<i>finder</i>	<i>funden</i>	<i>fundet</i>	<i>fundne</i>
$\text{EffSpell}(\hat{\sigma})$	<i>fandt</i>	<i>find</i>	<i>finder</i>	<i>funden</i>	<i>fundet</i>	<i>fundne</i>
$\text{BlotV}_1(\text{EffSpell}(\hat{\sigma}))$	<i>fandt</i>	<i>fan</i>	<i>fander</i>	<i>fanden</i>	<i>fandet</i>	<i>fandne</i>
$\text{EssRoot}(\text{EffSpell}(\hat{\sigma}))$	<i>fand</i>	<i>find</i>	<i>find</i>	<i>fund</i>	<i>fund</i>	<i>fund</i>
$\text{BlotV}_1(\text{EssRoot}(\text{EffSpell}(\hat{\sigma})))$	<i>fand</i>	<i>fan</i>	<i>fand</i>	<i>fand</i>	<i>fand</i>	<i>fand</i>

Similar to what we did in §4.2.2 for the case of English, we will construct a “heredity test function”  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  from a “simple heredity test function” in Algorithm 5.6 and a set of “admissible suffix mismatch and vowel alternation” rules in Algorithm 5.8.

**Algorithm 5.6** (Simple heredity test). *Suppose that  $\hat{\alpha}$  and  $\hat{\beta}$  are both lowercase strings. Set  $\mathbf{V} = (a|e|i|o|u)$ . The Boolean-valued function  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if  $\hat{\alpha}$  contains at least one instance of  $\mathbf{V}$ , AND at least one of the following three conditions holds:<sup>53</sup>*

- (i)  $\hat{\alpha} = \hat{\beta}$ ;
- (ii)  $\hat{\beta} = \hat{\alpha}t$ ;
- (iii)  $\ell(\hat{\beta}) > \ell(\hat{\alpha}) > \frac{\ell(\hat{\beta})}{2}$  AND  $\hat{\alpha} = \hat{\beta}^{[\ell(\hat{\alpha})]}$  AND  $\hat{\beta}^{\{\ell(\hat{\alpha})+1\}} = (e|s)$ . (See Definition 3.1 for the notations  $\hat{\beta}^{[n]}$  and  $\hat{\beta}^{\{n\}}$ .)

In what follows, we shall define  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta})$  and  $\text{SuffixSW}(\hat{\alpha}, \hat{\beta})$  as what is done in the English case (Algorithm 4.11), but the definitions of  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{NW}^*(\hat{\alpha}, \hat{\beta})$  and  $\text{SuffixSW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{SW}^*(\hat{\alpha}, \hat{\beta})$  will be subtly different.

**Algorithm 5.7** (Roots and Suffixes by NW and SW). *For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the function values  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta})$  and  $\text{NW}^*(\hat{\alpha}, \hat{\beta})$  are determined through the following procedure:*

- Use the Needleman–Wunsch algorithm to align the sequences as  $\text{NW}(\hat{\alpha}, \hat{\beta})$ .
- If  $\text{NW}(\hat{\alpha}, \hat{\beta})$  ends with a mismatch (shown in brackets in Examples 3.7.1 and 3.7.2), use this mismatch to define  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta})$ ; otherwise, define  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = [\emptyset, \emptyset]$ .
- Construct  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$  from all the (not necessarily contiguous) matching string portions in  $\text{NW}(\hat{\alpha}, \hat{\beta})$ , joined together by  $\# = \boxed{\text{U+0023}}$ :
- If the first mutation bracket in  $\text{NW}(\hat{\alpha}, \hat{\beta})$  does not appear at the end [which in turn, makes a non-void contribution to  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta})$ ], then define  $\text{NW}^*(\hat{\alpha}, \hat{\beta})$  as this first mismatch; otherwise, define  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \emptyset$ ;

The function values  $\text{RootSW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{SuffixSW}(\hat{\alpha}, \hat{\beta})$  and  $\text{SW}^*(\hat{\alpha}, \hat{\beta})$  are determined similarly, through the SW function.

*Example 5.7.1.* Let  $\hat{\alpha} = \text{fandt}$ ,  $\hat{\beta} = \text{find}$ , then we have  $\text{NW}(\hat{\alpha}, \hat{\beta}) = \text{SW}(\hat{\alpha}, \hat{\beta}) = f[a, i]nd[t, \emptyset]$ , leading to  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = \text{SuffixSW}(\hat{\alpha}, \hat{\beta}) = [t, \emptyset]$ ,  $\text{RootNW}(\hat{\alpha}, \hat{\beta}) = \text{RootSW}(\hat{\alpha}, \hat{\beta}) = \#nd$ ,  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \text{SW}^*(\hat{\alpha}, \hat{\beta}) = [a, i]$ .

**Algorithm 5.8** (Admissible suffix mismatch and vowel alternation). *For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function*

$$\text{AdmMut}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta}))$$

*returns **TRUE** if*

$$\text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = [(\emptyset | (e|nd|t)_m | ((bar|dom|hed|ig|lig|skab)\mathbf{X}), (\emptyset | (e|nd|t)_m | ((bar|dom|hed|ig|lig|skab)\mathbf{X}))]$$

**AND** at least one of the following two conditions holds:

- (i)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \emptyset$  AND  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$  contains at least one instance of  $\mathbf{V} = (a|e|i|o|u)$ ;
- (ii)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = ([a, i] | [a, o] | [a, u] | [e, i] | [i, u] | [o, u])$ .

<sup>53</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical OR.

Similarly, one can evaluate another Boolean-valued function

$$\text{AdmMut}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta}))$$

by trading all the occurrences of NW in the statements above with SW.

**Algorithm 5.9** (Heredity test function). *For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  returns TRUE if at least one of the following three conditions holds:*

- (i)  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta}) = \text{TRUE}$ ;
- (ii)  $\ell(\text{RootNW}(\hat{\alpha}, \hat{\beta})) \geq \frac{\max\{\ell(\hat{\alpha}), \ell(\hat{\beta})\}}{2}$  AND  $\text{AdmMut}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta})) = \text{TRUE}$ ;
- (iii)  $\ell(\text{RootSW}(\hat{\alpha}, \hat{\beta})) \geq \frac{\max\{\ell(\hat{\alpha}), \ell(\hat{\beta})\}}{2}$  AND  $\text{AdmMut}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta})) = \text{TRUE}$ .

Our approximate clustering algorithm for Danish differs from the English counterpart (Algorithm 4.15) in one specific detail: in Algorithm 5.10 below, vowel blotting is applied to both the effective spelling and the essential root in two stages, while Algorithm 4.15 has restricted the use of vowel blotting to the essential root.

**Algorithm 5.10** (Approximate clustering of Danish words). *The approximate clustering of a list of Danish words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  is completed in two stages:*

- (1) *We sort the list  $\{(\hat{\alpha}_1, \text{EffSpell}(\hat{\alpha}_1), \text{BlotV}_1(\text{EffSpell}(\hat{\alpha}_1))), \dots, (\hat{\alpha}_N, \text{EffSpell}(\hat{\alpha}_N), \text{BlotV}_1(\text{EffSpell}(\hat{\alpha}_N)))\}$  alphabetically according to the third component (with higher priority) and the second component (with lower priority). If two consecutive neighbors in the alphabetized list  $\{(\hat{\alpha}_{(1)}, \text{EffSpell}(\hat{\alpha}_{(1)}), \text{BlotV}_1(\text{EffSpell}(\hat{\alpha}_{(1)}))), \dots, (\hat{\alpha}_{(N)}, \text{EffSpell}(\hat{\alpha}_{(N)}), \text{BlotV}_1(\text{EffSpell}(\hat{\alpha}_{(N)})))\}$  satisfy  $\text{HrdTest}(\text{EffSpell}(\hat{\alpha}_{(n)}), \text{EffSpell}(\hat{\alpha}_{(n+1)})) = \text{FALSE}$ , where  $n \in \mathbb{Z} \cap [1, N]$ , then add a demarcation line between these two entries. In this way, the alphabetized list is divided into separate groups of words tagged with their effective spellings:  $\{g_1 = \{(\hat{\alpha}_{(1,1)}, \text{EffSpell}(\hat{\alpha}_{(1,1)}), \text{BlotV}_1(\text{EffSpell}(\hat{\alpha}_{(1,1)}))), \dots, (\hat{\alpha}_{(1,n_1)}, *, *)\}, \dots, g_M = \{(\hat{\alpha}_{(M,1)}, *, *), \dots, (\hat{\alpha}_{(M,n_M)}, *, *)\}\}$ ,<sup>54</sup> where each sublist also preserves alphabetic order. In particular, we have  $\hat{\alpha}_{(1,1)} = \hat{\alpha}_{(1)}$  and  $\hat{\alpha}_{(M,n_M)} = \hat{\alpha}_{(N)}$ .*
- (2) *For each group of words  $g_m = \{(\hat{\alpha}_{(m,1)}, *, *), \dots, (\hat{\alpha}_{(m,n_m)}, *, *)\}$  where  $m \in \mathbb{Z} \cap [1, M]$ , we augment it into a tagged entry  $G_m = (g_m, \text{EffSpell}(\hat{\alpha}_{(m,1)}), \text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})), \text{BlotV}_1(\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))))$ . The list  $\{G_1, \dots, G_M\}$  is sorted alphabetically, with respect to  $\text{BlotV}_1(\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})))$  (with highest priority),  $\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))$  (with medium priority), and  $\text{EffSpell}(\hat{\alpha}_{(m,1)})$  (with lowest priority). If two consecutive neighbors in the alphabetized list  $\{G_{(1)} = (g_{(1)}, \hat{\gamma}'_{(1)}, \hat{\gamma}''_{(1)}, \hat{\gamma}'''_{(1)}), \dots, G_{(M)} = (g_{(M)}, \hat{\gamma}'_{(M)}, \hat{\gamma}''_{(M)}, \hat{\gamma}'''_{(M)})\}$  satisfy*

$$\text{HrdTest}(\hat{\gamma}''_{(m)}, \hat{\gamma}''_{(m+1)}) = \text{FALSE}$$

AND

$$\text{HrdTest}(\hat{\gamma}'_{(m)}, \hat{\gamma}'_{(m+1)}) = \text{FALSE}$$

where  $m \in \mathbb{Z} \cap [1, M]$ , then add a demarcation line between these two entries. In this manner, the alphabetized list  $\{G_{(1)}, \dots, G_{(M)}\}$  is divided into separate groups of tagged clusters  $\{\Gamma_1 = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, \Gamma_K = \{G_{(K,1)}, \dots, G_{(1,m_K)}\}\}$ . Finally, the output list of word clusters  $\{\check{\Gamma}_1, \dots, \check{\Gamma}_K\}$  is constructed by discarding all the tags (effective spellings, essential roots, vowel blotted forms) from each  $\Gamma_k$ , where  $k \in \mathbb{Z} \cap [1, K]$ .

*Example 5.10.1.* As an illustration of our clustering algorithm, we pick the following families of Danish words, where an “approximate translation” in English, enclosed in quotation marks, is appended to the end of each family.<sup>55</sup>

and, anden, andens, ands, ænder, ænderne, ændernes, ænders — “duck”;  
 Anders, Andersen, Andersens — “Andrew”;  
 barn, barnet, barnets, barns, børn, børnene, børnenes, børns — “child”;  
 bonde, bonden, bondens, bondes, bønder, bønderne, bøndernes, bønders — “farmer”;  
 datter, datteren, datterens, datters, døtre, døtrene, døtrenes, døtres — “daughter”;  
 ferie, ferien, feriens, ferier, ferierne, ferierne, feriers, feries — “holiday”;  
 fængsel, fængselerne, fængselet, fængselets, fængsels, fængsler, fængslernes, fængslers, fængslet, fængslets — “jail”;  
 himle, himlen, himlene, himlenes, himlens, himmel, himmelen, himmelnens, himmels — “heaven”;  
 hus, hus’, huse, husende, husene, husenes, huser, huses, husets, hust, huste — “house”;

<sup>54</sup>Hereafter, we use an asterisk (\*) to abbreviate components that are clear from context.

<sup>55</sup>We note that within each family, words may be related to each other through both inflections (adjective comparisons, noun declensions, verb conjugations) and derivations (such as compound nouns formed from several components, verbs derived from nouns, etc.). Therefore, the “approximate translation” only matches certain members in the specified word family.

*husven, husvennen, husvennens, husvenner, husvennerne, husvennernes, husvenners, husvens* — “family friend”; *hænder, hænderne, hændernes, hænders, Haand, hånd, hånden, hånds* — “hand”; *kø, køen, køens, kør, kørne, kørernes, kørs, køs* — “queue”; *køb, købene, købenes, køber, købet, købets, købs, købt, købte, købtes* — “purchase”; *København, Københavns* — “Copenhagen”; *mand, manden, mandens, mands, mænd, mændene, mændenes, mænds* — “man”; *prins, prinse, prinselig, prinsen, prinsens, prinser, prinserne, prinsernes, prinsers, prinses* — “prince”; *prinsesse, prinsessekrone, prinsessen, prinsessens, prinsesser, prinsesserne, prinsessernes, prinsessers, prinsesses* — “princess”; *ske, skeen, skeens, skeer, skeerne, skeernes, skeers, skes* — “spoon”; *ske, sker, sket, skete* — “happen”; *stor, store, stort, større, størst* — “big”; *ung, ungdom, unge, ungt, yngre, yngst* — “young”; *ven, vennen, vennens, venner, vennerne, vennernes, venners, vens* — “friend”; *venlig, venligere, venligst, venligste, venligt* — “friendly”; *vent, vente, ventede, ventedes, venter, ventes, ventet* — “wait”; *vin, vine, vinen, vinene, vinenes, vinens, vines, vins* — “wine”; *vinter, vinteren, vinterens, vinters, vintre, vintrene, vintrenes, vintres* — “winter”; *vinterferie, vinterferien, vinterferiens, vinterferier, vinterferierne, vinterferierne, vinterferierne, vinterferiers, vinterseries* — “winter holiday”; *ø, øen, øens, øer, øerne, øernes, ørs, øs* — “island”; *å, åen, åens, åer, åerne, åernes, åers, ås* — “creek”; *år, årene, årenes, året, årets, års* — “year”.

Applying Algorithm 5.10 to this list of words, we obtain the following clustering results:

{*and, anden, andens, Anders, Andersen, Andersens, ands, ænder, ænderne, ændernes, ænders*}, {*barn, barnet, barnets, barns, børn, børnene, børnenes, børns*}, {*bonde, bonden, bondens, bondes, bønder, bønderne, bøndernes, bønders*}, {*datter, datteren, datterens, datters, døtre, døtrene, døtrenes, døtres*}, {*ferie, ferien, feriens, ferier, ferierne, feriernes, feriers, series*}, {*fængsel, fængselerne, fængselet, fængselets, fængsels, fængsler, fængslernes, fængslers, fængslet, fængslets*}, {*himle, himlen, himlene, himlenes, himlens, himles, himmel, himmelen, himmelens, himmels*}, {*hus, hus’, huse, husende, husene, husenes, huser, huses, huset, husets, hust, huste, husven, husvennen, husvennens, husvenner, husvennerne, husvennernes, husvenners, husvens*}, {*hænder, hænderne, hændernes, hænders, Haand, hånd, hånden, hånds*}, {*kø, køen, køens, kør, kørne, kørernes, kørs, køs*}, {*køb, købene, købenes, køber, købet, købets, købs, købt, købte, købtes*}, {*København, Københavns*}, {*mand, manden, mandens, mands, mænd, mændene, mændenes, mænds*}, {*prins, prinse, prinsen, prinsens, prinser, prinserne, prinsernes, prinses, prinses, prinsesse, prinsessekrone, prinsessen, prinsessens, prinsesser, prinsesserne, prinsessernes, prinsessers, prinsesses*}, {*prinselig*}, {*ske, skeen, skeens, skeer, skeerne, skeernes, skeers, sker, skes, sket, skete*}, {*stor, store, stort, større, størst*}, {*ung, ungdom, unge, ungt, yngre, yngst*}, {*ven, veninde, veninden, venindens, veninder, veninderne, venindernes, veniders, venindes, vennen, vennens, venner, vennerne, vennernes, venners, vens*}, {*venlig, venligere, venligst, venligste, venligt*}, {*vent, vente, ventede, ventedes, venter, ventes, ventet*}, {*vin, vine, vinen, vinene, vinenes, vinens, vines, vins*}, {*vinter, vinteren, vinterens, vinters, vintre, vintrene, vintrenes, vintres*}, {*vinterferie, vinterferien, vinterferiens, vinterferier, vinterferierne, vinterferierne, vinterferierne, vinterferierne, vinterferiers, vinterseries*}, {*ø, øen, øens, øer, øerne, øernes, ørs, øs*}, {*å, åen, åens, åer, åerne, åernes, åers, ås*}, {*år, årene, årenes, året, årets, års*}.

*Example 5.10.2.* To further test our algorithm against Danish irregular verbs, we throw the following list:

*afbrudt, afbryd, afbryde, afbryder, afbrød* — “interrupt”; *bar, bær, bære, bærer, båret* — “carry”; *besat, besatte, besæt, besætte, besætter* — “occupy”; *beskrev, beskreven, beskrevet, beskrevne, beskriv, beskriver* — “describe”; *blotlagde, blotlagt, blotlæg, blotlægge, blotlægger* — “expose”; *bræk, brække, brækkede, brækker, brækket* — “break”; *drak, drik, drikker, drukken, drukket, drukne* — “drink”; *fandt, find, finder, funden, fundet, fundne* — “find”; *fik, få, fået, får* — “get”; *flyv, flyve, flyver, fløj, fløjet* — “fly”; *gav, gavs, giv, give, givende, givendes, giver, gives, givet, givets* — “give”; *gentag, gentagen, gentager, gentaget, gentagne, gentog* — “repeat”;

*græd, græde, græder, grædt* — “weep”;  
*hjulp, hjulpen, hjulpet, hjulpne, hjælp, hjælpe, hjælper* — “help”;  
*kom, komme, kommen, kommende, kommer, kommet, komne* — “come”;  
*lad, lader, ladet, ladt, lod* — “allow”;  
*lig, ligge, ligger, ligget, lå* — “lie (be in horizontal position)”;  
*løb, løbe, løber, løbet* — “run”;  
*optræd, optræde, optræder, optrådt, optrådte* — “appear”;  
*ryg, ryge, ryger, røg, røget* — “smoke”;  
*sang, sunget, syng, syngende, synger* — “sing”;  
*trak, trukken, trukket, trukne, træk, trækende, trækker* — “pull”;  
*æd, æde, æder, ædt, åd* — “eat (like an animal)”

into our algorithm, and obtain

{afbrudt, afbryd, afbryde, afbryder, afbrød}, {bar, bær, bære, bærer, båret}, {besat, besatte, besæt, besætte, besætter}, {beskrev, beskreven, beskrevet, beskrevne, beskriv, beskriver}, {blotlagde, blotlagt, blotlæg, blotlægge, blotlægger}, {braek, brække, brækkede, brækker, brækket}, {drak, drik, drikker, drukken, drukket, drukne}, {fandt, find, finder, funden, fundet, fundne}, {fik, få, fæt, får}, {flyv, flyve, flyver, fløj, fløjet}, {gav, gavs, giv, give, givende, givendes, giver, gives, givet, givets}, {gentag, gentagen, gentager, gentaget, gentagne, gentog}, {græd, græde, græder, grædt}, {hjulp, hjulpen, hjulpet, hjulpne, hjælp, hjælpe, hjælper}, {kom, komme, kommen, kommende, kommer, kommet, komne}, {lad, lader, ladet, ladt, lod}, {lig, ligge, ligger, ligget, lå}, {løb, løbe, løber, løbet}, {optræd, optræde, optræder, optrådt, optrådte}, {ryg, ryge, ryger, røg, røget}, {sang, sunget, syng, syngende, synger}, {trak, trukken, trukket, trukne, træk, trækende, trækker}, {æd, æde, æder, ædt, åd}.

*Example 5.10.3.* Combining the inputs from the last two examples, we obtain the following output:

{afbrudt, afbryd, afbryde, afbryder, afbrød}, {and, anden, andens, Anders, Andersen, Andersens, ands, ænder, ænderne, ændernes, ænders}, {bar, bær, bære, bærer, båret}, {barn, barnet, barnets, barns, børn, børnene, børneenes, børns}, {besat, besatte, besæt, besætte, besætter}, {beskrev, beskreven, beskrevet, beskrevne, beskriv, beskriver}, {blotlagde, blotlagt, blotlæg, blotlægge, blotlægger}, {bonde, bonden, bondens, bondes, bønder, bønderne, bøndernes, bønders}, {braek, brække, brækkede, brækker, brækket}, {datter, datteren, datterens, datters, døtre, døtrene, døtrenes, døtres}, {drak, drik, drikker, drukken, drukket, drukne}, {fandt, find, finder, funden, fundet, fundne}, {ferie, ferien, feriens, ferier, ferierne, feriernes, feriers, series}, {fik, få, fæt, får}, {flyv, flyve, flyver, fløj, fløjet}, {fængsel, fængselerne, fængselet, fængselets, fængsels, fængsler, fængslerne, fængslers, fængslet, fængslets}, {gav, gavs, giv, give, givende, givendes, giver, gives, givet, givets}, {gentag, gentagen, gentager, gentaget, gentagne, gentog}, {græd, græde, græder, grædt}, {himle, himlen, himlene, himlens, himles, himmel, himmelen, himmelens, himmels}, {hjulp, hjulpen, hjulpet, hjulpne, hjælp, hjælpe, hjælper}, {hus, hus'}, huse, husende, husene, husenes, huser, huses, huset, hust, huste, husven, husvennen, husvennens, husvenner, husvennerne, husvennerne, husvenners, husvens}, {hænder, hænderne, hændernes, hænders, Haand, hånd, hånden, hånds}, {kom, komme, kommen, kommende, kommer, kommet, komne}, {kø, køen, køens, kør, kørne, køernes, køers, køs}, {køb, købene, købenes, køber, kobet, købets, købs, købt, købte, købtes}, {København, Københavns}, {lad, lader, ladet, ladt, lod}, {lig, ligge, ligger, ligget, lå}, {løb, løbe, løber, løbet}, {mand, manden, mandens, mands, mænd, mændene, mændenes, mænds}, {optræd, optræde, optræder, optrådt, optrådte}, {prins, prinse, prinsen, prinsens, prinser, prinserne, prinsernes, prinsers, prinses, prinsesse, prinsessekrone, prinsessen, prinsessens, prinsesser, prinsesserne, prinsessernes, prinsessers, prinsesses}, {prinselig}, {ryg, ryge, ryger, røg, røget}, {sang, sunget, syng, syngende, synger}, {ske, skeen, skeens, skeer, skeerne, skeernes, skeers, sker, skes, sket, skete}, {stor, store, stort, større, størst}, {trak, trukken, trukket, trukne, træk, trækende, trækker}, {ung, ungdom, unge, ungt, yngre, yngst}, {ven, veninde, veninden, venindens, veninder, veninderne, venindernes, veninders, venindes, vennen, vennens, venner, vennerne, vennernes, venners, vens}, {venlig, venligere, venligst, venligste, venligt}, {vent, vente, ventede, ventedes, venter, ventes, ventet}, {vin, vine, vinen, vinene, vinenes, vinens, vines, vins}, {vinter, vinteren, vinterens, vinters, vintre, vintrene, vintrenes, vintres}, {vinterferie, vinterferien, vinterferiens, vinterferier, vinterferierne, vinterferierne, vinterferiernes, vinterferiers, vinterseries}, {æd, æde, æder, ædt, åd}, {ø, øen, øens, øer, øerne, øernes, øers, øs}, {å, åen, åens, åer, åerne, åernes, års}, {år, årene, årenes, året, årets, års}.

### 5.1.3 Heuristic detection of compounds

In Danish (as well as German and Dutch), compounds are usually spelt as a single word, without spaces or hyphens between the constituting components. For example, the Danish compound *vinterferie* “winter holiday” is formed by joining *vinter*

“winter” and *ferie* “holiday”. Aside from direct concatenation, Danish compounds may also involve an intervening letter *e* or *s* [46, §12.2.2.2].

Automatic detection of the “hidden word boundaries” within Germanic compounds is a non-trivial task. However, if we limit our scope to compounds with only two constituting components that both appear in isolations in a vocabulary list, there are heuristic methods to break down these binary compound words. We will integrate such a heuristic compound detection method with Danish word clustering, later in Algorithm 5.13 (which is a variation on Algorithm 5.10), where any detected binary compound is “dissolved” into its two constituting components (“head” and “tail”) and filed as two copies, one clustered with the “head” word, the other with the “tail” word.

Before stating Algorithm 5.13 in full, we need some preparations (Definition 5.11 and 5.12) to sort out potential binary compounds and their constituting components.

**Definition 5.11** (String minus operation). Suppose that  $\ell(\hat{\beta}) \geq \ell(\hat{\alpha})$  and  $\hat{\beta} = \hat{\beta}^{[\ell(\hat{\alpha})]} \hat{\tau}$ , then  $\hat{\beta} \ominus \hat{\alpha} = \hat{\tau}$  defines the “string minus operation”. In other words,  $\hat{\beta} \ominus \hat{\alpha}$  is a string that leaves out the first  $\ell(\hat{\alpha})$  characters in  $\hat{\beta}$ .  $\square$

**Algorithm 5.12** (Heuristic identification of Danish binary compounds). Let  $\Lambda^{\hat{\rho}} = \{\hat{\rho}_1, \dots, \hat{\rho}_Q\}$  be a list of distinct Danish essential roots (without vowel blotting) that contain at least one instance of  $\mathbf{V} = (a|e|i|o|u)$  and DO NOT match the following string patterns:

$$(\pi\rho a|a|a\bar{s}ag|ak|al|an\mathbf{X}|ar|at|be|ben|dag|du|frem|lo|mi|mod|na|nar|\\ ned|om|op|ot|ov|rak|re|rig|rod|ru|sam|tag|u|ud|und|vfor\mathbf{X}|\mathbf{X}ant|zil|\mathbf{X}).$$

The output of the function  $\text{CpdDet}(\Lambda^{\hat{\rho}})$  is obtained through the following procedures:

- (1) Construct a list  $\{(\hat{\rho}_1, \lambda_1^{\hat{\rho}}), \dots, (\hat{\rho}_Q, \lambda_Q^{\hat{\rho}})\}$  where  $\lambda_q^{\hat{\rho}} = \{\hat{\rho}_{(q,1)}, \dots, \hat{\rho}_{(q,n_q)}\}$  is a subset of  $\Lambda^{\hat{\rho}}$  whose members all match the string pattern  $\hat{\rho}_{q\sim}$ , for  $q \in \mathbb{Z} \cap [1, Q]$ .
- (2) Expand the aforementioned entry  $(\hat{\rho}_q, \lambda_q^{\hat{\rho}})$  into a list of triplets  $\{(\hat{\rho}_{(q,1)}, \hat{\rho}_q, \hat{\rho}_{(q,1)} \ominus \hat{\rho}_q), \dots, (\hat{\rho}_{(q,n_q)}, \hat{\rho}_q, \hat{\rho}_{(q,n_q)} \ominus \hat{\rho}_q)\}$  for every  $q \in \mathbb{Z} \cap [1, Q]$  such that  $\lambda_q^{\hat{\rho}} \neq \emptyset$ . Collect all these triplets as one runs through the list  $\{(\hat{\rho}_1, \lambda_1^{\hat{\rho}}), \dots, (\hat{\rho}_Q, \lambda_Q^{\hat{\rho}})\}$ . The list of these triplets  $\{(\hat{\rho}_{\langle q \rangle}, \hat{\eta}_{\langle q \rangle}, \hat{\rho}_{\langle q \rangle} \ominus \hat{\eta}_{\langle q \rangle}), \dots, (\hat{\rho}_{\langle Q' \rangle}, \hat{\eta}_{\langle Q' \rangle}, \hat{\rho}_{\langle Q' \rangle} \ominus \hat{\eta}_{\langle Q' \rangle})\}$  contains potentially valid decompositions of compounds.
- (3) Screen the aforementioned list of triplets as follows: for every  $q' \in \mathbb{Z} \cap [1, Q']$ , if  $(\hat{\rho}_{\langle q' \rangle}, \hat{\eta}_{\langle q' \rangle}, \hat{\tau}_{\langle q' \rangle} = \hat{\rho}_{\langle q' \rangle} \ominus \hat{\eta}_{\langle q' \rangle})$  satisfies

$$\ell(\hat{\rho}_{\langle q' \rangle} \ominus \hat{\eta}_{\langle q' \rangle}) \geq 2 \quad \text{AND} \quad \hat{\rho}_{\langle q' \rangle} \ominus \hat{\eta}_{\langle q' \rangle} = \mathbf{X}_1(a|e|i|o|u)\mathbf{X}_2,$$

then construct  $\hat{\tau}_{\langle q' \rangle}^*$  by performing  $(e|s|t) \sim \rightarrow \emptyset$  on  $\hat{\tau}_{\langle q' \rangle}$  and  $\hat{\tau}_{\langle q' \rangle}^{**}$  by doing  $(er|es) \sim \rightarrow \emptyset$  on  $\hat{\tau}_{\langle q' \rangle}$ , before generating a list  $\lambda_{\langle q' \rangle}^{\hat{\tau}}$  by members of  $\Lambda^{\hat{\rho}}$  that match the pattern  $(\hat{\tau}_{\langle q' \rangle}|\hat{\tau}_{\langle q' \rangle}^*|\hat{\tau}_{\langle q' \rangle}^{**})$ ; otherwise, set  $\lambda_{\langle q' \rangle}^{\hat{\tau}} = \emptyset$ .

- (4) Collect all the triplets  $(\hat{\rho}_{\langle q' \rangle}, \hat{\eta}_{\langle q' \rangle}, \lambda_{\langle q' \rangle}^{\hat{\tau}})$  where  $\lambda_{\langle q' \rangle}^{\hat{\tau}}$  is non-void and  $\hat{\tau}_{\langle q' \rangle}$  DOES NOT match the following string patterns:

$$(dom\mathbf{X}|hed\mathbf{X}|lig\mathbf{X}|om).$$

This list of triplets  $\text{CpdDet}(\Lambda^{\hat{\rho}})$  contains the heuristic decompositions of all the identified binary compounds.

**Algorithm 5.13** (Approximate clustering of Danish words with heuristic detection of compounds). The approximate clustering of a list of Danish words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  respecting compounding is completed in four stages:

- (1) Do as in Algorithm 5.10(1).
- (2) Do as in Algorithm 5.10(2). Save both the tagged clusters  $\{\Gamma_1 = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, \Gamma_K = \{G_{(K,1)}, \dots, G_{(1,m_K)}\}\}$  and the list of word clusters  $\{\check{\Gamma}_1, \dots, \check{\Gamma}_K\}$  for further use.
- (3) Construct a tagged list of word clusters  $\{(\check{\Gamma}_1, \Lambda_1^{\hat{\rho}}), \dots, (\check{\Gamma}_K, \Lambda_K^{\hat{\rho}})\}$  where  $\Lambda_k^{\hat{\rho}}$  is the union of all Danish essential roots (without vowel blotting) available to  $\Gamma_k$ , for  $k \in \mathbb{Z} \cap [1, K]$ . Set  $\Lambda^{\hat{\rho}} = \Lambda_1^{\hat{\rho}} \cup \dots \cup \Lambda_K^{\hat{\rho}}$ , and evaluate  $\text{CpdDet}(\Lambda^{\hat{\rho}})$ .
- (4) The first component  $\hat{\rho}_{\langle q'' \rangle}$  of each triplet  $(\hat{\rho}_{\langle q'' \rangle}, \hat{\eta}_{\langle q'' \rangle}, \lambda_{\langle q'' \rangle}^{\hat{\tau}})$  in  $\text{CpdDet}(\Lambda^{\hat{\rho}})$  is called a “dissolvable compound”, the second component  $\hat{\eta}_{\langle q'' \rangle}$  a “heuristic head”, and the first member in the third component  $\lambda_{\langle q'' \rangle}^{\hat{\tau}}$  a “heuristic tail”. In the tagged list  $\{(\check{\Gamma}_1, \Lambda_1^{\hat{\rho}}), \dots, (\check{\Gamma}_K, \Lambda_K^{\hat{\rho}})\}$ , every entry containing a “dissolvable compound” is removed, and regrouped with the entries matching its “heuristic head” and “heuristic tail”. Finally, remove all tags.

*Example 5.13.1.* We may revisit Example 5.10.3 with the algorithm above. The results are given below:

{afbrudt, afbryd, afbryde, afbryder, afbrød}, {and, anden, andens, Anders, Andersen, Andersens, ands, ænder, ænderne, ændernes, ænders}, {bar, bær, bære, bærer, båret}, {barn, barnet, barnets, barns, børn, børnene, børnenes, børns}, {besat, besatte, besæt, besætte, besætter}, {beskrev, beskreven, beskrevet, beskrevne, beskriv, beskriver}, {blotlagde, blotlagt, blotlæg, blotlægge, blotlægger}, {bonde, bonden, bondens, bondes, bønder, bønderne, bøndernes, bønders}, {bræk, brække, brække, brække, brækket}, {datter, datteren, datterens, datters, døtre, døtrene, døtrenes, døtres}, {drak, drik, drikker, drukken, drukket, drukne}, {fandt, find, finder, funden, fundet, fundne}, {ferie, ferien, feriens, ferier, ferierne, feriernes, feriers, series, vinterferie, vinterferien, vinterferiens, vinterferier, vinterferierne, vinterferiernes, vinterferiers, vinterseries}, {fik, få, fæt, får}, {flyv, flyve, flyver, fløj, fløjet}, {fængsel, fængselerne, fængselet, fængselets, fængsels, fængsler, fængslernes, fængsler, fængslets}, {gav, gavs, giv, give, givende, givendes, giver, gives, givet, givets}, {gentag, gentagen, gentager, gentaget, gentagne, gentog}, {græd, græde, græder, grædt}, {himle, himlen, himlene, himlenes, himlens, himles, himmel, himmelen, himmelens, himmels}, {hjalp, hjulpen, hjulpel, hjulpe, hjælp, hjælpe, hjælper}, {hus, hus}, {huse, husende, husene, husen, huses, huset, hust, huste, husven, husvennen, husvennens, husvenner, husvennerne, husvennerne, husvenners, husvens}, {hænder, hænderne, hændernes, hænders, Haand, hånd, hånden, hånds}, {kom, komme, kommen, kommende, kommer, kommet, komne}, {kø, køen, køens, kør, kørne, kørernes, kørs, køs}, {køb, købene, købenes, køber, købet, købets, købs, købt, købte, købtes}, {København, Københavns}, {lad, lader, ladet, ladt, lod}, {lig, ligge, ligget, lå}, {løb, løbe, løber, løbet}, {mand, manden, mandens, mands, mænd, mændene, mændenes, mænds}, {opræd, opræde, opræder, oprådt, oprådte}, {prins, prinse, prinselig, prinsen, prinsens, prinser, prinserne, prinsers, prinses, prinsesse, prinsesekrone, prinsessen, prinsessens, prinsesser, prinsesserne, prinsessernes, prinsessers, prinsesses}, {ryg, ryge, ryger, røg, røget}, {sang, sunget, syng, syngende, synger}, {ske, skeen, skeens, skeer, skeerne, skeernes, skeers, sker, skes, sket, skete}, {stor, store, stort, større, størst}, {trak, trukken, trukket, trukne, trek, trækende, trækker}, {ung, ungdom, unge, ungt, yngre, yngst}, {ven, veninde, veninden, venindens, veninder, veninderne, venindernes, venindes, venlig, venligere, venligst, venligste, venligt, vennen, vennens, venner, vennerne, vennernes, venners, vens}, {vent, vente, ventede, ventedes, venter, ventes, ventet}, {vin, vine, vinen, vinene, vinenes, vines, vins}, {vinter, vinteren, vinterens, vinterferie, vinterferien, vinterferiens, vinterferier, vinterferierne, vinterferierne, vinterferiers, vinterferies, vinters, vintre, vintrene, vintrenes, vintres}, {æd, æde, æder, ædt, åd}, {ø, øen, øens, øer, øerne, øernes, øers, øs}, {å, åen, åens, åer, åerne, åernes, åers, ås}, {år, årene, årenes, året, årets, års}.

Here, the word family generated by *vinterferie* “winter holiday” is correctly redistributed to the constituting components *vinter* “winter” and *ferie* “holiday”. In the meantime, there are no inappropriate compound detections, such as grouping *bære* “carry” with *barn* “child”.

*Example 5.13.2.* In Fig. S4, we further apply the aforementioned word clustering algorithm to topic extraction and machine translation (see Table S1 for text source).

The Danish word *ved* can mean “at” or “know” (present tense), which can be either a function word or a highly frequent content word (from the English perspective). We have chosen not to include *ved* in our list of Danish stop words (Definition 5.1), but have made no efforts to resolve its polysemy. Thus, in the second cluster of Fig. S4a, we have conflated several concepts that are spelt as *ved* (“at” or “know”) and *vide* (“know” or “wide”).

Moreover, the English topic *de* (part of the proper name “de Bourgh”) cannot be correctly translated by our algorithm, because the list of Danish stop words (Definition 5.1) includes *de* “the (plural), they, those”. Nevertheless, we consider a link between English *de* and Danish *Bourgh* an excusable error (marked in amber color in Fig. S4b, since these two words always stay together in the entire text—with nearly identical kinetic behavior).

It should be noted that Jane Austen used the word *related* mostly in the sense of *told* in her novel, which is correctly translated into *fortalte* in Danish.

MR VED SAGT ELIZABETH KOMME FÅ DARCY GÖDT LIGE MRS SE BENNET GÅ BINGLEY GIVE TAGE JANE  
SØSTER TÅNKET SAMMEN LANGE MISS TID NØGENHED ONSKE STØRDE HØRT TRODE GLÆDE SVAREDE FAMILIE STÅB MIND WICKHAM FINDE COLLINS VIS  
DAGEN GANG HELE LYDIA TO SIKKER HÅBE MAND ØJEBLIKT ØJEBLIKT LADE REJST STED TALE VEN SYNE BREV MORREN FAR DÆTRER KERE UNGE SKONT LADY  
MORBROTH CATHERINE KENDTE GRUND MÅDE HUSFØR SPÆRGÅS HOJST GIFTED DAMER SAGT OMPLERKSOMHED FOELSELSPRÆGET SANDSKØN SKRÆVE GARDINER HOLD TALTE  
AR. ORD MULIGUD UDTRYK FALDS VIRKELIG LONDON UDRØR HERRER MORBROTH FREM REG HJEMM HØFLIGT PÅFOLGEN VEJEN LIZZY SATTE SELSKAB AFTENEN LIV BESØG HENGJUVENDE  
FORSTÅ Mennesker LYKKE LONGBOURN HURTIGT LYKKELIG LOVE PIGEEN VENLIG VENDTE KVINDEN NOGENSENDE CHARLOTTE MODER KÆRLIGHED BEGGE LUCAS NVÆRN BEKENDTSAHL IMIDLERDT NEFRÆ MENTE  
NÅEG KORT SIR FORSTE SKYLDEN BESLUTTE GENSHED SÆRLIG SAMTALE PÅSYDE TIDSPOINT ÅH SMÆRE SØM BEHAGELIG SMILE ALENE LEST TIDLIGERE TILFOJEDE KITTY ARE NETHERFIELD ELSKVÆRTIGS SÆRBURDE

**SPILLE** FORENT SKYDDE BESØKED  
**ALVORLIG** BERENDTSKAB STOLTSED LOBET MENER SLAGS STUEN VIST SANDED STADIG OBERST  
**GERNE** ENERHED ASSESS ROG<sup>1</sup> PØDAGOE TROST<sup>2</sup> TILLES<sup>3</sup> FAST PAR FALLS<sup>4</sup> TILDE<sup>5</sup> HØRT<sup>6</sup> PLAN<sup>7</sup> LYS<sup>8</sup>  
**BESTEMT** TABEGL<sup>9</sup> VENTE<sup>10</sup> VENINDES<sup>11</sup> TIME<sup>12</sup> FORSKØL<sup>13</sup> SIDE<sup>14</sup> FORSGØR<sup>15</sup> STØRST<sup>16</sup> FORTELL<sup>17</sup> GÆRD<sup>18</sup>  
**FORBEDRET** AGTEGL<sup>19</sup> MÅL<sup>20</sup> FORM<sup>21</sup> AIJERTE<sup>22</sup> OVERASKET<sup>23</sup> FORT<sup>24</sup> ESKER<sup>25</sup> OPFATTELL<sup>26</sup> FULDSTED<sup>27</sup> FEJL<sup>28</sup> SKIFFE<sup>29</sup> MERKE<sup>30</sup>  
**MØRGEN** ROVDEDE<sup>31</sup> YNGSTE<sup>32</sup> FOLES<sup>33</sup> ISER<sup>34</sup> BØRS<sup>35</sup> BØRN<sup>36</sup> IVRIG<sup>37</sup> ART<sup>38</sup> KØMPLIMENT<sup>39</sup> FORBØDE<sup>40</sup> OVERBEVIS<sup>41</sup> HALV<sup>42</sup>  
**UGRUNN** MENTER<sup>43</sup> SLUG<sup>44</sup> OUTAY<sup>45</sup> UDGET<sup>46</sup> FORSKELIG<sup>47</sup> ØVERSTED<sup>48</sup> NET<sup>49</sup> ALDRIT<sup>50</sup> FØR<sup>51</sup> SVIGER<sup>52</sup> TAV<sup>53</sup>  
**DEB** VERDEN JAMES<sup>54</sup> UDSK<sup>55</sup> MINDE<sup>56</sup> PØDE<sup>57</sup> NIECE<sup>58</sup> INTERESE<sup>59</sup> BRUD<sup>60</sup> BUND<sup>61</sup> PENG<sup>62</sup> LYTTEN<sup>63</sup> KØS<sup>64</sup> OPTAGT<sup>65</sup>  
**INDVENDINGEN** FORHE<sup>66</sup> RESPECT<sup>67</sup> BEVEGELSE<sup>68</sup> TJENER<sup>69</sup> SPANSIS<sup>70</sup> TIDELIG<sup>71</sup> BÆRTING<sup>72</sup> BÅND<sup>73</sup> TESTIMONIUM<sup>74</sup> MØSTER<sup>75</sup> GIFT<sup>76</sup>  
**URHÆGLIG** SITUATION<sup>77</sup> ØGE<sup>78</sup> OVER<sup>79</sup> FOLK<sup>80</sup> PLÅGE<sup>81</sup> NEVO<sup>82</sup> ANKOSKT<sup>83</sup> VERELSE<sup>84</sup> UNDVALD<sup>85</sup> NYGRÆDDER<sup>86</sup> ELDE<sup>87</sup> GLÆDE<sup>88</sup> APPAS<sup>89</sup> VÆR<sup>90</sup> UDVIDE<sup>91</sup> EMET<sup>92</sup> BRENDE<sup>93</sup>  
**OPGÅNNING** SPØRG<sup>94</sup> FORM<sup>95</sup> YKL<sup>96</sup> TÅND<sup>97</sup> TILDE<sup>98</sup> PLÅGE<sup>99</sup> RYGT<sup>100</sup> OVERÅSLÆ<sup>101</sup> TENDER<sup>102</sup> SAUTING<sup>103</sup> ØPENNING<sup>104</sup> MØDT<sup>105</sup> DERBYSHIRE<sup>106</sup> BRIGHTON<sup>107</sup> BORTSET<sup>108</sup> ÅNGA<sup>109</sup> SYND<sup>110</sup> BEREDE<sup>111</sup>  
**FINE** ABRÅ<sup>112</sup> MARIA<sup>113</sup> LØFT<sup>114</sup> FOR<sup>115</sup> ENGLELIGHED<sup>116</sup> ELEGANT<sup>117</sup> BETRAYING<sup>118</sup> BØRDEM<sup>119</sup> GÆRD<sup>120</sup> GÆRD<sup>121</sup> GØDE<sup>122</sup> FORBLØFF<sup>123</sup> PAR<sup>124</sup> HALMT<sup>125</sup> FORBLØFF<sup>126</sup> VOL<sup>127</sup> ØM<sup>128</sup> SKØNNED<sup>129</sup> BEMÆR<sup>130</sup> VOL<sup>131</sup> TILST<sup>132</sup> SØDANVÆR<sup>133</sup> NO<sup>134</sup> LUN<sup>135</sup> INVITATION<sup>136</sup>

(a)

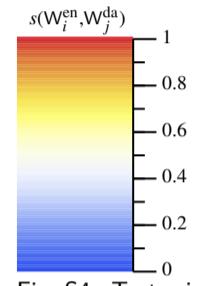
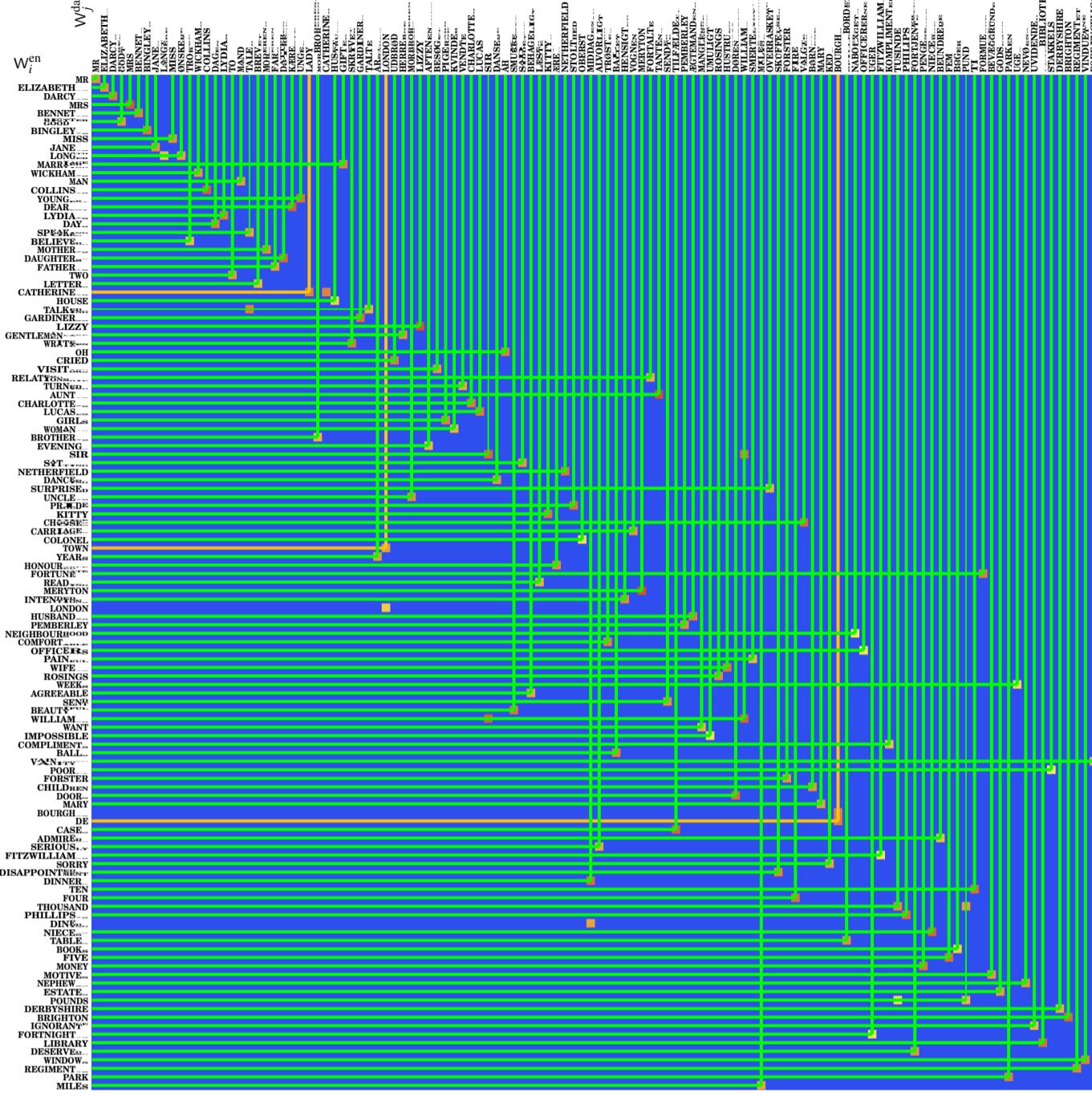


Fig. S4. Text mining in Danish. (a) Statistically identified topics ( $n_{ii} \geq 20$ ) in a Danish version of *Pride and Prejudice*, with the same color encoding scheme as Fig. S3. (b) Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{da}})$  between selected topics in English and Danish versions of *Pride and Prejudice*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) rank of the similarity score for the optimal node. Green (resp. amber) cross-hair indicates an exact (resp. a close but non-exact) match.

## 5.2 Modified Porter stemming algorithm for German

**Definition 5.14** (German stop words). If the lowercase<sup>56</sup> form of a word belongs to the following list<sup>57</sup>:

ab, aber, all, alle, allein, allem, allen, aller, alles, als, also, am, an, ander, andere, anderem, anderen, anderer, anderes, anderm, andern, anders, aneinander, auch, auf, aufeinander, aufs, aus, auseinander, ausser, außer, außerdem, außerdem, bei, beide, beidem, beiden, beider, beides, beieinander, beim, bevor, bin, binnen, bis, bisher, bislang, bist, da, dabei, dafür, dagegen, dahin, damit, danach, dann, daran, darauf, darf, darfst, darin, darüber, darum, das, dass, daß, dasselbe, davon, dazu, dein, deine, deinem, deinen, deiner, deines, deins, dem, deme, demselben, den, denen, denn, dennoch, denselben, der, deren, derer, dern, dero, derselbe, derselben, des, deshalb, deß, desselben, dessen, desto, dich, die, dies, diese, dieselbe, dieselben, diesem, diesen, dieser, dieses, dir, doch, dort, dorther, dorthin, du, durch, durcheinander, dürfe, dürfen,dürfend,dürfest,dürfet,dürft,dürfte, durften,dürften,dürftest,dürftest,dürftet,ebenso, ein, einander, eine, einem, einen, einer, eines, einig, einige, einigem, einigen, einiger, einiges, einmal, empor, er, es, etwas, euch, euer, eure, euren, eurer, eures, falls, für, füreinander, ganz, ganze, ganzem, ganzen, ganzer, ganzes, gar, gedurft, gegen, gegeneinander, gegenüber, gehabt, gekonnt, gemocht, gemusst, gemußt, genug, gerade, geradezu, gesollt, getan, gewesen, gewollt, geworden, hab, habe, haben, habend, habest, habet, habt, hast, hat, hatte, hätten, hatten, hattest, hättest, hattet, hättet, her, herab, heran, herauf, heraus, herbei, herein, herüber, hervor, hier, hierher, hierhin, hin, hinab, hinauf, hinaus, hindurch, hinein, hingegen, hinten, hinter, hintereinander, hinüber, hinunter, hinweg, hinzu, ich, ihm, ihn, ihnen, ihr, ihre, ihrem, ihren, ihrer, ihres, im, immer, immerhin, in, indem, indes, indessen, ineinander, ins, iregendetwas, iregendwo, irgend, irgendein, irgendeine, irgendeinem, irgendeinen, irgendeiner, irgendeines, irgendemand, irgendemandem, irgendemanden, irgendwann, irgendwas, irgendwelche, irgendwelchem, irgendwelchen, irgendwelcher, irgendwelches, irgendwer, irgendwie, irgendwohin, ist, ja, je, jede, jedem, jeden, jeder, jedes, jedoch, jemand, jemandem, jemanden, jemandes, jene, jenem, jenen, jener, jenes, jetzt, kann, kannst, kaum, kein, keine, keinem, keinen, keiner, keines, könne, können, könnend, könnest, könnet, könnt, konnte, könnte, konnten, könnten, konntest, könntest, könntet, könntet, los, machen, mag, magst, man, manche, manchem, manchen, mancher, manches, manchmal, mehr, mehre, mehrere, mehres, mein, meine, meinem, meinen, meiner, meines, meins, mich, mir, mit, miteinander, mochte, möchte, mochten, möchten, mochtest, möchtest, mochtet, möchtet, möge, mögen, mögend, mögest, möget, mögt, muss, muß, müsse, müssen, müssend, müsstest, müsstet, müsst, müßt, müßt, müßt, müsste, müßte, müsstet, müßten, müsstet, müßten, müßt, müsstest, müßtest, müsstest, müßtest, müßtet, müßtet, müßtet, na, nach, nachdem, nacheinander, nebeneinander, nein, nicht, nichts, nie, niemals, niemand, niemandem, niemanden, niemandes, niemands, nimmer, noch, nun, nur, ob, oben, obgleich, obwohl, obzwar, oder, oft, öfter, öftesten, ohne, pro, schon, sehr, sei, seien, seiet, sein, seine, seinem, seinen, seiner, seines, seins, seist, seit, seitdem, selbst, sich, sie, sind, so, sobald, sogar, solange, solch, solche, solchem, solchen, solcher, solches, soll, solle, sollen, sollend, sollest, sollet, sollst, sollt, sollte, sollten, solltest, solltet, sondern, sonst, soviel, soweit, sowie, sowohl, stets, tat, tätet, taten, täten, tatest, tätest, tatet, tätet, trotz, trotzdem, tu, tue, tuen, tuend, tuest, tuet, tun, ture, tust, tut, über, überdies, übereinander, um, umsonst, und, uns, unse, unsem, unsen, unser, unsere, unserem, unser, unserer, unseres, unserm, unsern, unses, unsre, unsrem, unsren, unsrer, unsres, unter, untereinander, usw, viel, viele, vielen, vieles, vielleicht, vielmehr, vom, von, voneinander, vor, vorbei, vordem, voreinander, vorüber, während, wann, war, wär, ward, warden, wardest, wardet, wardst, wäre, waren, wären, wärest, wäret, warst, wärst, wart, warum, was, weg, weil, weiter, welche, welchem, welchen, welcher, welches, wem, wen, wenig, wenige, wenigen, weniger, wenigsten, wenigstens, wenn, wenngleich, wer, werd, werde, werden, werdend, werdest, werdet, weshalb, wessen, weswegen, wie, wieder, wieso, wieviel, wiewohl, will, willst, wir, wird, wirst, wo, woher, wohin, wohl, wolle, wollen, wollend, wollst, wollet, wollt, wollte, wollten, wolltest, wolltet, worden, wurde, würde, wurden, würden, wurdest, würdest, wurdet, würdet, ziemlich, zu, zudem, zueinander, zum, zur, zusammen, zurück, zuvor, zwar, zwischen,

then we consider it a German stop word. All the German stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list.  $\square$

The post-1996 German orthography requires writing *ss* in place of  $\beta$  in certain words. This new rule affects the stop word *daß*, whose new spelling *dass* has been added to our stop word list. We need to add *schon* “already” to the list of stop words, to avoid conflation with *schön* “beautiful” during approximate clustering (which ignores all umlauts). Although the words *schon* and *schön* were indeed etymologically related (reminiscent of the adverbial and adjectival uses of the word *pretty* in English), their current meanings are well separated.

<sup>56</sup>For the matching of German stop words, lowercase forms are enforced before comparison. Therefore, both *sie* “she, they” and *Sie* “you (polite form)” appearing in a German text are considered stop words, even though the capitalized form *Sie* does not appear in our list.

<sup>57</sup> Our list of German stop words is based on [snowball.tartarus.org/algorithms/german/stop.txt](http://snowball.tartarus.org/algorithms/german/stop.txt), with extensive additions that roughly match their counterparts in English.

For clustering German words, it is necessary that we ignore umlauts, so that *Apfel* “apple” is clustered with *Äpfel* “apples” etc. We also need to convert all the occurrences of  $\beta$  to *ss* (which matches the orthographic practice for German spelling in Switzerland), in order to accommodate to texts written before the German spelling reform in 1996: *Kuß* “kiss” (noun, pre-1996), *Kuss* “kiss” (noun, post-1996), *Küsse* “kisses” (noun, pre- and post-1996). This necessarily carries certain risk, as *Buße* “penance” (pre- and post-1996) will be conflated with *Busse* “buses” (pre- and post-1996).

All German nouns are capitalized, while adjectives and verbs derived from nouns are not: cf. *Deutsch* “German” (noun, “the German language”) and *deutsch* “German” (adjective). It is thus advisable to ignore capitalization during clustering, contrary to the practice in English.

### 5.2.1 Effective spelling and essential root

**Algorithm 5.15** (German effective spelling). Set  $\mathbf{V} = (a|e|\dot{e}|i|\dot{i}|o|\dot{o}|u)$  and  $\mathbf{V}^* = (a|\dot{e}|i|\dot{i}|o|\dot{o}|u)$ . For a German word  $\hat{\sigma}$ , its effective spelling  $\text{EffSpell}(\hat{\sigma})$  is constructed in the following steps.<sup>58</sup>

(1) Convert to lowercase, and replace

$(\emptyset \text{ver})st(a e i o \ddot{u})rb\sim$		$\hat{x}\text{kunft}$		$\hat{x}sicht$		$\hat{x}tr(\ddot{a} a)\text{cht}$		$a\text{bwesen}(d h)\mathbf{X}$				
$\sigma\tau\epsilon\rho b$		$\hat{x}\hat{x}^+\text{kunft}$		$\hat{x}\times_2 s i \chi t$		$\hat{x}^+\hat{x}\text{tracht}$		$\alpha\beta\sigma\epsilon n\tau$				
acht~	$\ddot{a}hn$	$a\text{nwesen}(d h)\mathbf{X}$	$b(\ddot{u} u)\text{ch}$	$bald\sim$	$benn\sim$	$bessie\sim$	$blut\sim$	$born\sim$	$brot$	$dame$	$dämmere$	
8e	$\sigma\mu\ddot{a}hn$	$\pi\text{resen}t$	$\beta u k$	$\beta a\lambda\delta$	$\beta e n n$	$\beta e o\sigma\theta ie$	$\beta \lambda w o\delta$	$\beta r u n n$	$\beta r o t$	$\ddot{a}ame$	$\ddot{a}awmep$	
dent~	$diana\sim$	$dicht$	$ehe\sim$	$erinn\sim$	$f(\ddot{u} u)(ss \beta)$	$f e s s e l$	$f e s s l$	$f e s t \sim$	$feuer$	$g a r d i \sim$		
$\delta e\tilde{v}t\tau$	$\delta i a\tilde{v}na$	$\theta i k t$	$\mu r e$	$\mu \tilde{v}e m p n$	$\varphi i t \sigma \zeta$	$\varphi e \sigma e l$	$\varphi e \sigma l$	$\varphi e \sigma t$	$\varphi e u e p$	$\gamma a r \tilde{v} i$		
$g e f (\ddot{a} a) h r$	$g e g e n w a r t \sim$	$g e h e i m$	$g e s t a l t$	$g e w a l t \sim$	$g o t t \sim$	$h(\ddot{a} a ie) l t$	$h(\ddot{o} o) h e$	$h a l l$				
$\delta a\tilde{v}ypr$	$\pi r e s e n t$	$\sigma e k r e t$	$\sigma \eta a i p$	$\varphi o p c t$	$\gamma o t t$	$\chi o \lambda \delta$	$h o c h e$	$h o \sigma a l l$				
$h a n n a h \sim$	$h e l e n \sim$	$h e l l$	$h u l l$	$i n n e \sim$	$i r r$	$j o n e \sim$	$l e a h \sim$	$l e d i g$	$l i e b s t$	$l i e d$	$l o n d \sim$	
$\chi a\tilde{v}n\tilde{v}a\chi$	$\varepsilon \lambda e \tilde{v} \eta$	$h \chi e l l$	$\omega r a \pi$	$i \tilde{v} n e$	$i u p r$	$\gamma v o n e$	$\lambda e \alpha a \eta$	$\sigma i \tilde{v} y$	$l i e b$	$\lambda i e \delta$	$\lambda o n \delta$	
$m(ama utti)\sim$	$m(\ddot{u} u)n d \sim$	$m e i l e$	$m i n u t$	$m i \ddot{u} d$	$o a \sim$	$o b e r s t$	$o r t \sim$	$p f e r d$	$p l ö t z$	$p o r t r ä t$	$q u e l l$	
$m u t t e r$	$\mu i \tilde{v} \delta$	$\mu e i l e$	$\mu i \tilde{v} \omega t$	$\mu t i \rho \delta$	$\omega a$	$\omega b e r s t$	$w o r t$	$\chi o \rho s$	$\sigma i \ddot{u} \delta e \tilde{v}$	$p o r t r \tilde{v} \delta t$	$q v e l l$	
$r(\ddot{o} o)t \sim$	$r e c h t$	$r e e \sim$	$r e i c h$	$r e i s$	$r e l i g$	$r i (s s  \beta) \sim$	$r o b e r t$	$r o c h e s \sim$	$s(a \ddot{a})(s s  \beta)$			
$r o \theta$	$\rho i \gamma \chi \tau$	$\rho e e$	$\rho p e i c h$	$\rho e e i \sigma s$	$\rho e l i y$	$\rho i \varphi \tau$	$\rho o \beta e r t$	$\rho \omega x \epsilon \sigma$		$s i t z$		
$s c h u l t e r \sim$	$s o m m e r \sim$	$s p i e l$	$s t i l l$	$t a l e n t$	$t h e m$	$t o t \sim$	$v e r s u c h \sim$	$w e g$	$w e l k$	$w e l l \sim$	$w o h l$	$a c h$
$s c h u l \delta e r$	$\sigma o \mu m e r$	$\pi \lambda a l$	$\sigma t i l l$	$\tau a \lambda e \tilde{v} t$	$\theta e \mu$	$\sigma t e r b$	$\nu e r s u \chi$	$\omega e y$	$\omega i \theta p k$	$\omega \tilde{v} \delta a l$	$wh o l$	$\alpha a \chi$
$b e s e s s e n$	$de$	$g a t t e (\emptyset n)$	$g a t t i n (\emptyset n e n)$	$g e s e s s e n$	$g i b t$	$m r s$	$t o d (\emptyset e e n e s s)$		$\sim l e u t e (\emptyset n)$			
$b e s i t z e n$	$\delta e e$	$e h e m a n n$	$e h e f r a u$	$s i t z e n$	$g e b e n$	$z \phi p a u$	$\sigma t e r b$		$m a n n$			

(2) Replace

$(\emptyset ge)gönn\sim$		$(\emptyset ge)h i n d \sim$		$(b \ddot{o}s b o s h e i t b o s h a f t)\mathbf{X}$		$\hat{x}^*(a e i k n o p s t u y)t i s c h \mathbf{X}$						
$\emptyset$	$g o e n n$	$\chi i n d$		$b b o e s$		$\hat{x} i k$						
$X^e(a i)l i s t$	$X^e(g e  k s c h s s  \beta)l i n g$	$\ddot{a}$	$a b s e n t$	$b e k \sim$	$b e s t \hat{x}^*(e) \sim$	$b e t t$	$b r a n n t \sim$	$ch$	$e b \sim$	$ei$		
$X l$	$X l a n g$	$a$	$\alpha \beta \sigma e n t$	$b e k$	$b e s t \hat{x}$	$z b e t t$	$b r e n n$	$\zeta$	$x e b$	$\hat{e}$		
$er \tilde{v} \sigma t$	$f i n g e r$	$f i \tilde{u} n f$	$g e g e s s$	$g e s u n d \sim$	$h ü b e$	$j a n u \sim$	$l a n g s a m \sim$	$le b$	$le i b \sim$	$l ö s$	$m ä d$	$m a n i e r$
$er \tilde{v} \sigma t$	$\varphi i n y e r$	$5 e$	$e s s$	$y e s u n d$	$h e b e$	$j n u$	$\gamma l a n g s a m$	$le \beta$	$le i b$	$\lambda e i \sigma t$	$ma k d$	$mu a n i e r$
$monat$	$ö$	$o b \sim$		$p a s s a g$	$r u h$		$s c h l e c h t$		$s c h o c k$		$s e t z$	$s i t z$
$\mu o n a t$	$o$	$o b$		$p a \sigma \sigma a g$	$\rho u h$		$s c h l e c h t$		$s c h o c k$		$\zeta e t z$	$\sigma i t z$
$sonn$	$\beta$	$tr ä n$	$tz$	$\ddot{u}$	$\ddot{u} \ddot{e} r$		$\ddot{u} \ddot{e} r r e d \sim$		$v e r k ü n d$		$v i e r \sim$	$v u l g$
$\sigma o \tilde{v} \tilde{v}$	$ss$	$\tau e a \rho$	$z t z$	$u$	$x \ddot{u} \ddot{e} r$		$\pi e \sigma \beta$		$\pi \rho o \tilde{v} \sigma$		$4 e$	$x v x u l g$
$wahr$	$w a s s e r$	$w e l t$	$w e r t$	$w ü n s c h$	$z o l l$	$g e b r o c h e n$	$l a b e$	$l a s t$	$m i s s (\emptyset e s)$		$\sim i s (m u s  t i s c h)\mathbf{X}$	
$v w a h r$	$w a \sigma e r$	$x w e l t$	$f w e r t$	$w u n s \zeta$	$\tau o \lambda \lambda$	$b r e c \tilde{e} n$	$\lambda a \beta$	$l e s e n$	$\mu i \sigma$		$\emptyset$	

(3) Replace

<sup>58</sup>To avoid confusion of Greek nu with Latin vee, we write the former as  $\tilde{v}$  in substitution rules.

$(\emptyset be)str\sim$	$(\emptyset ge)last\sim$	$(\emptyset ge)lieb\sim$	$(bra\hat{c}te gebra\hat{c}t)$	$(lass liess)$	$abend\sim$	$ahn\sim$						
$\sigma tr$	$xlast$	$lxieb$	$bringe$	$lass$	$azbend$	$xzahn$						
$alt\sim$	$anz\hat{\chi}\#(u)$	$arb$	$be\sigmaitz$	$bes\hat{\chi}\#(s t)\sim$	$bet\sim$	$bild$	$bitter\sim$	$bus\sim$	$\zeta s\hat{\chi}\#(t)$	$da\hat{c}t$	$dank$	$dien$
$a\lambda\tau$	$ant\hat{\chi}$	$arb$	$\varphi o\sigma tu\tilde{v}$	$ges\hat{\chi}$	$\beta\beta et$	$\beta\beta il$	$\beta\beta itter$	$zbus$	$x\hat{\chi}$	$denk$	$thank$	$dien$
$dumm$	$erb$	$fass$	$f\hat{e}\sim$	$folg$	$froh$	$fund\hat{\chi}\#(e)$	$gast$	$gel(a i u)ng\sim$	$gew(a i o)nn$	$grund$	$gruss$	
$\delta umm$	$e\beta\beta e$	$\varphi ass$	$fv\hat{e}$	$vfolg$	$vroh$	$\varphi und\hat{\chi}$	$\gamma ast$	$ge\lambda ing$	$wi\tilde{v}$	$grund$	$kgruss$	
$gut$	$hand$	$hart$	$hu\sim$	$ie$	$in^X\epsilon(d g k)\cdot$	$kind$	$klag$	$komi$	$komp$	$kumm$	$kund$	$kuns$
$besser$	$han\delta$	$xhart$	$shu$	$\hat{i}$	$\hat{i}nX$	$gkind$	$klag$	$komi$	$comp$	$xkummmxkund$	$gkuns$	$\lambda\hat{e}d$
$lock$	$lub$	$lust$	$mutter$	$nahr$	$name$	$nummer$	$ohr\sim$	$ort\sim$	$platz$	$raf\sim$	$re\check{c}n$	$rest\sim$
$lock$	$\lambda ub$	$zlust$	$\mu other$	$\check{v}ahr$	$\check{v}ame$	$xnunmmer$	$o\hat{h}p$	$ort$	$pplatz$	$\rho af$	$\rho e\check{c}n$	$pest$
$s\zeta l$	$segn$	$setz$	$si\check{c}t$	$sieg$	$sohn$	$stamm$	$stand$	$stemm$	$su\check{c}$	$tan\check{z}$	$tis\check{c}$	$tugend$
$s\zeta \lambda$	$szegn$	$zsetz$	$\sigma i\check{c}t$	$sxi\check{e}g$	$szohn$	$\sigma tamm$	$steh$	$\sigma ztemm$	$zsu\check{c}$	$ta\check{v}z$	$tis\chi$	$xtugend$
$voll$			$w(\hat{e} ie)s\hat{\chi}\#(s)$		$w\hat{e}n$		$wind$		$winter$		$uhr\sim$	$verg\sim$
$ffoll$			$zxw\hat{e}\hat{s}\hat{\chi}$		$xv\hat{e}n$		$wind$		$wwinter$		$wort$	$wusst$
$geh(\emptyset e en est st t)$		$bat(\emptyset en est et)$		$gebeten$		$herr(\emptyset en n)$			$kam(\emptyset e en est et st t)$		$mr$	
$ging$		$bitten$		$bitten$		$hherr$			$komm$		$hherr$	
$oh$		$sass(\emptyset en est et t)$		$segen(\emptyset s)$		$tante(\emptyset n)$		$\sim cs$	$\sim ement(\emptyset s)$	$\sim enz(\emptyset en)$		$\sim w(\hat{e} ie)s$
$xcoh$		$sitz$		$szegn$		$\tau ante$		$x$	$\hat{i}ren$	$\emptyset$		$zxw\hat{e}s$

- (4) If the result so far is an empty string, stop here. Otherwise, break down the result into  $\hat{\sigma}_1\hat{\sigma}_2$ , where  $\ell(\hat{\sigma}_1) = \max\{3, \ell(\hat{\sigma})\}$ , obtain  $\hat{\sigma}'_1$  by doing  $kam\sim\rightarrow z kam$  on  $\hat{\sigma}_1$  before deducing  $\hat{\sigma}'_2$  from  $\hat{\sigma}_2$  in two sequential steps:

- (3.1) Do  $el\rightarrow l$ ;  
 (3.2) Do  $kationX\rightarrow z\hat{i}ren$ ,  $\sim ln\rightarrow l$ ,  $(ffol|haf\ddot{t}|ist|massig)X\rightarrow\emptyset$ .

Concatenate  $\hat{\sigma}'_1$  and  $\hat{\sigma}'_2$ .

- (5) Do  $s\zeta lo\zeta\sim\rightarrow s\zeta los$ ,  $\sim\hat{\chi}\#(a|e|\hat{e}|i|\hat{o}|s|u)s\rightarrow\hat{\chi}$ ,  $\sim\hat{\chi}\#V(\zeta en|ling)(\emptyset|e|en|s)\rightarrow\hat{\chi}$ ,  $\sim(h|k)\hat{e}t(\emptyset|en)\rightarrow\emptyset$ ,  $\sim ik(\emptyset|en|er(\emptyset|n|s))\rightarrow\emptyset$ ,  $\sim l\hat{e}n(\emptyset|s)\rightarrow\emptyset$ ,  $\sim V\mathbf{X}_{\hat{\chi}\#}V^*st\rightarrow V\mathbf{X}\hat{\chi}$  on  $\hat{\sigma}'_1\hat{\sigma}'_2$ .

**Definition 5.16** (German protected range). Set  $\mathbf{V}^* = (a|\hat{e}|i|\hat{o}|u)$ ,  $\mathbf{C}_{m_0} = \overline{(a|e|\hat{e}|i|\hat{o}|u)}_{m_0}$  and  $\mathbf{C}_m = \overline{(a|e|i|\hat{o}|u|\hat{e}|i)}_m$ . Let  $\hat{\sigma}$  be the effective spelling of a German word, its protected range  $\text{ProtRg}(\hat{\sigma}) = \max\{\lambda_1(\hat{\sigma}), \lambda_2(\hat{\sigma})\}$  is determined by two non-negative integers  $\lambda_1(\hat{\sigma})$  and  $\lambda_2(\hat{\sigma})$  specified through the following procedures:

- Look for the string pattern  $((\emptyset|a|be|e|\hat{e}|er|ge|i|mit|na\zeta|u|vor|x\ddot{u}ber|zer|zu)\mathbf{C}_{m_0}\mathbf{V}^*\mathbf{C}_m e)\overline{(a|e|\hat{e}|i|\hat{o}|r|u)}_{m_0}\sim$  in the string  $\hat{\sigma}$ ;
- If the string pattern above is found, the last position occupied by such a string defines  $\lambda_1(\hat{\sigma})$ ; otherwise, set  $\lambda_1(\hat{\sigma}) = 0$ ;
- Look for the pattern  $(a|o|u)$  in the string  $\hat{\sigma}$ ;
- If a letter in the pattern above is found, the last position occupied by such a letter defines  $\lambda_2(\hat{\sigma})$ ; otherwise, set  $\lambda_2(\hat{\sigma}) = 0$ .  $\square$

Unlike Danish, but similar to Dutch, our German essential root is free from common separable verb prefixes (see step (3) in the algorithm below).

**Algorithm 5.17** (German essential root). Set  $\mathbf{V} = (a|e|\hat{e}|i|\hat{o}|u)$  and  $\mathbf{C} = \overline{(a|e|\hat{e}|i|\hat{o}|u)}$ . Let  $\hat{\sigma}$  be the effective spelling of a German word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:

- (1) Break down  $\hat{\sigma} = \hat{\sigma}_1\hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}^{[\text{ProtRg}(\hat{\sigma})]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma})$  is equal to the protected range of  $\hat{\sigma}$ .

- (2) On  $\hat{\sigma}_2$ , perform the following substitutions in a sequel:<sup>59</sup>

- (2.1) Do  $\sim s\rightarrow\emptyset$ ;  
 (2.2) Do  $\sim e(d|e|m|n|r|st|t)_{m_0}\rightarrow\emptyset$ .

<sup>59</sup>In other words, the core algorithm for essential root extraction runs as follows: keep the last “strong” vowel *a*, *i*, *o* or *u* in non-final position, plus one subsequent letter; delete final *a*; erase the final appearance of *e* and all the letters thereafter.

The result after these two steps of operations is called  $\hat{\sigma}'_2$ .

- (3) Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ , before doing

$$(\emptyset|ab|an|auf|aus|bē|durç|ēn|h(er|in)(\emptyset|ab|an|auf|aus|ēn|xüber)|mit|naç|um|vor|wêter|w(i|i)der|xüber|zu(\emptyset|ruck))(\emptyset|ge|zu)\mathbf{X}\mathbf{V} \sim \rightarrow \mathbf{X}\mathbf{V};$$

- (4) Do  $(\hat{\chi}_{\times 2}t|\sim\hat{\chi}_{\times 2}) \rightarrow \hat{\chi}$ ;<sup>60</sup> (Check Example 3.4.2 to see the notation  $\hat{\chi}_{\times 2}$  for double letters.)

- (5) Do  $\sim\mathbf{C}t \rightarrow \mathbf{C}$ ,  $nd \rightarrow \tilde{v}d$ .

### 5.2.2 Admissible mutation and approximate clustering

Like English, German has hundreds of irregular verbs in daily use. Unlike English, these German verbs may not only have irregular past tense and past participles, but also have irregular present tense conjugations.

The German vowel blotting mechanism is similar to the Danish version (Algorithm 5.5).

**Algorithm 5.18** (German vowel blotting). Set  $\mathbf{V}_m = (a|e|\hat{e}|i|\hat{i}|o|u)_m$  and  $\mathbf{C}_{m_0} = \overline{(a|e|\hat{e}|i|\hat{i}|o|u)}_{m_0}$ . For a token string  $\hat{\sigma}$ , its blotted form  $\text{BlotV}_1(\hat{\sigma})$  is constructed as follows:

- If the string pattern  $(\emptyset|a|be|e|\hat{e}|i|o|u)\mathbf{C}_{m_0}\mathbf{V}_m \sim$  can be found in the string  $\hat{\sigma}$ , then the last position occupied by such a pattern is replaced by the letter “a”.
- Otherwise, leave the string  $\hat{\sigma}$  intact.

Similar to what we did in §5.1.2 for the case of Danish, we will construct a bivariate Boolean-valued function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  on a “simple heredity test function” in Algorithm 5.19, and a set of “admissible suffix mismatch and vowel alternation” rules in Algorithm 5.20.

**Algorithm 5.19** (Simple heredity test). Suppose that  $\hat{\alpha}$  and  $\hat{\beta}$  are both lowercase strings. Set  $\mathbf{V} = (a|e|\hat{e}|i|\hat{i}|o|u)$ . The Boolean-valued function  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if  $\hat{\alpha}$  contains at least one instance of  $\mathbf{V}$ , **AND** at least one of the following four conditions holds:<sup>61</sup>

- $\hat{\alpha} = \hat{\beta}$ ;
- $\hat{\beta} = \hat{\alpha}d$ ;
- $\hat{\beta} = \hat{\alpha}t$ ;
- Appending the last character of  $\hat{\alpha}$  to itself, one obtains  $\hat{\beta}$ , i.e.  $\hat{\alpha}\Omega(\hat{\alpha}) = \hat{\beta}$ ;
- $\ell(\hat{\beta}) > \ell(\hat{\alpha}) \geq \frac{\ell(\hat{\beta})}{2}$  **AND**  $\hat{\alpha} = \hat{\beta}^{[\ell(\hat{\alpha})]}$  **AND**  $\hat{\beta}^{[\ell(\hat{\alpha})+1]} = (e|n|s)$ . (See Definition 3.1 for the notations  $\hat{\beta}^{[n]}$  and  $\hat{\beta}^{\{n\}}$ .)

In what follows, we define  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{NW}^*(\hat{\alpha}, \hat{\beta})$  and  $\text{SuffixSW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{SuffixSW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{SW}^*(\hat{\alpha}, \hat{\beta})$  as what is done in the Danish case (Algorithm 5.7).

**Algorithm 5.20** (Admissible suffix mismatch and vowel alternation). For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function

$$\text{AdmMut}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta}))$$

returns **TRUE** if<sup>62</sup>

$$\begin{aligned} (\text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = [(\emptyset|(d|e|n|t)_m|(\Omega(\text{RootNW}(\hat{\alpha}, \hat{\beta}))|bar|er|ig|in|\hat{r}|isç|liç|sam|st|tum|ung)\mathbf{X})), \\ (\emptyset|(d|e|n|t)_m|(\Omega(\text{RootNW}(\hat{\alpha}, \hat{\beta}))|bar|er|ig|in|\hat{r}|isç|liç|sam|st|tum|ung)\mathbf{X}))] \end{aligned}$$

$$\text{OR } \text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = [ih, og] \text{ OR } \text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = [og, ih])$$

**AND** at least one of the following three conditions holds.<sup>63</sup>

- $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \emptyset$  **AND**  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$  contains at least one instance of  $(a|e|\hat{e}|i|\hat{i}|o|u|z)$ ;

<sup>60</sup>In Mathematica codes, the substitution rule reads  $\{x\_ \sim x\_ \sim "t" \mid \text{WordBoundary} :> x\}$ .

<sup>61</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical **OR**.

<sup>62</sup>Depending on the programming language chosen, the string *ih* may be sorted before or after *og*.

<sup>63</sup>We note that  $\text{NW}^*(\hat{\alpha}, \hat{\beta})$  and  $\text{NW}^*(\hat{\beta}, \hat{\alpha})$  differ only in the order of the two components in the bracket.

- (ii)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = ([a, e][a, i][a, \hat{i}][a, \hat{i}][a, o][a, u][ah, i][ah, o][au, \hat{i}][au, o][e, i][e, \hat{i}][e, o][\hat{e}, i][\hat{e}, \hat{i}][eh, i][eh, o][i, o][\hat{i}, i][\hat{i}, o][\hat{i}, u][o, u]);$
- (iii)  $\text{NW}^*(\hat{\beta}, \hat{\alpha}) = ([a, e][a, i][a, \hat{i}][a, \hat{i}][a, o][a, u][ah, i][ah, o][au, \hat{i}][au, o][e, i][e, \hat{i}][e, o][\hat{e}, i][\hat{e}, \hat{i}][eh, i][eh, o][i, o][\hat{i}, i][\hat{i}, o][\hat{i}, u][o, u]).$

Similarly, one can evaluate another Boolean-valued function

$$\text{AdmMut}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta}))$$

by trading all the occurrences of NW in the statements above with SW.

**Algorithm 5.21** (Heredity test function). *In what follows,  $\ell([\hat{\sigma}, \hat{\tau}]) = \min\{\ell(\hat{\sigma}), \ell(\hat{\tau})\}$ . For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  returns TRUE if at least one of the following three conditions holds:*

- (i)  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta}) = \text{TRUE};$
- (ii)  $\ell(\text{RootNW}(\hat{\alpha}, \hat{\beta})) + \ell(\text{SuffixNW}(\hat{\alpha}, \hat{\beta})) + \ell(\text{NW}^*(\hat{\alpha}, \hat{\beta})) \geq \frac{\max\{\ell(\hat{\alpha}), \ell(\hat{\beta})\}}{2}$  AND  $\text{AdmMut}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta})) = \text{TRUE};$
- (iii)  $\ell(\text{RootSW}(\hat{\alpha}, \hat{\beta})) + \ell(\text{SuffixSW}(\hat{\alpha}, \hat{\beta})) + \ell(\text{SW}^*(\hat{\alpha}, \hat{\beta})) \geq \frac{\max\{\ell(\hat{\alpha}), \ell(\hat{\beta})\}}{2}$  AND  $\text{AdmMut}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta})) = \text{TRUE}.$

**Algorithm 5.22** (Approximate clustering of German words). *The algorithm is essentially the same as Algorithm 5.10, except that German rules (instead of Danish rules) apply to all the tags (effective spelling, essential root, vowel blotting etc.).*

*Example 5.22.1.* As an illustration of our clustering algorithm, we pick the following families of German words, where an “approximate translation” in English, enclosed in quotation marks, is appended to the end of each family:

*Apfel, Äpfel, Äpfelchen, Äpfeln, Apfels* — “apple”;  
*Arbeit, arbeite, arbeiten, Arbeiten, arbeitest, arbeitet, arbeitete, arbeiteten, gearbeitet* — “work”;  
*Arbeiter, Arbeitern, Arbeiters* — “worker”;  
*Arbeitsplatz, Arbeitsplätze, Arbeitsplätze, Arbeitsplätzen, Arbeitsplatzes* — “workplace”;  
*Arbeitstag, Arbeitstage, Arbeitstagen, Arbeitstages, Arbeitstags* — “workday”;  
*Bus, Busse, Bussen, Busses* — “bus”;  
*Buße, Bußen* — “penance”;  
*dunkel, dunkelste, dunkelstem, dunkelsten, dunkelster, dunkelstes, dunkle, dunklem, dunklen, dunkler, dunklere, dunklerem, dunkleren, dunklerer, dunkleres, dunkles* — “dark”;  
*Ei, Eie, Eier, Eiern, Eies* — “egg”;  
*Eigelb, Eigelbe, Eigelben, Eigelbs* — “egg yolk”;  
*Eis, Eise, Eises* — “ice”;  
*frei, freie, freiem, freien, freier, freiere, freierem, freieren, freierer, freieres, freies, freiste, freistem, freisten, frei-  
sten, freister, freistes* — “free”;  
*Freiheit, Freiheiten* — “freedom”;  
*gehaust, haus, Haus, Häuschen, hause, Hause, hausen, hausend, Häuser, Hauses, haust, hauste, hauste, hausted* — “house”;  
*geküsst, Kuss, Kuß, Küsschen, Kusse, küsse, Küsse, küssen, Küssen, Kusses, küsstest, küsst, küsstest, küsst-  
tet* — “kiss”;  
*gelb, gelbe, gelbem, gelben, gelber, gelbere, gelberem, gelberen, gelberer, gelberes, gelbes, gelbste, gelbstem,  
gelbsten, gelbster, gelbstes* — “yellow”;  
*gross, groß, große, großem, großen, großer, größere, größerem, größeren, größerer, größeres, großes, größte,  
größtem, größten, größter, größtes* — “big”;  
*Öl, Öle, Ölen, Öles, Öls* — “oil”;  
*Platz, Platze, Plätze, Plätzten, Platzes* — “place”;  
*studier, studiere, studiere, studiere, studieren, studierend, studierest, studieret, studierst, studiert, studiert, stu-  
dierte, studierte, studierten, studiertest, studiertest, studiertet* — “study”;  
*Tag, Tage, Tagen, Tages, Tags* — “day”.

Applying Algorithm 5.22 to this list of words, we obtain the following clustering results:

{Apfel, Äpfel, Äpfelchen, Äpfeln, Apfels}, {Arbeit, arbeite, arbeiten, Arbeiten, Arbeit, Arbeitern, Arbeiters, arbeitest, arbeitet, arbeitete, arbeiteten, gearbeitet}, {Arbeitsplatz, Arbeitsplätze, Arbeitsplätze, Arbeitsplätzen, Arbeitsplatzes}, {Arbeitstag, Arbeitstage, Arbeitstagen, Arbeitstages, Arbeitstags}, {Bus, Busse, Buße, Bussen, Bussen, Busses}, {dunkel, dunkelste, dunkelstem, dunkelsten, dunkelster, dunkelstes, dunkle, dunklem, dunklen, dunkler, dunklere, dunklerem, dunkleren, dunklerer, dunkleres, dunkles}, {Ei, Eie, Eier, Eiern, Eies}, {Eigelb, Eigelbe, Eigelben, Eigelbs}, {Eis, Eise, Eises}, {frei, freie, freiem, freien, freier, freiere, freierem, freieren, freierer, freieres, freies, Freiheit, Freiheiten, freiste, freistem, freisten, freister, freistes}, {gehaust, haus, Haus, Häuschen, hause, Hause, hausen, hausend, Häuser, Hauses, haust, hauste, hausted}, {geküsst, Kuss, Kuß, Küsschen, Kusse, küsse, Küsse, küssten, Küssen, Kusses, küsstest, küsst, küsstest, küsstet}, {gelb, gelbe, gelbem, gelben, gelber, gelbere, gelberem, gelberen, gelberer, gelberes, gelbes, gelbste, gelstem, gelbsten, gelbster, gelbstes}, {gross, groß, große, großem, großen, großer, größere, größerem, größeren, größerer, größeres, großes, größte, größtem, größten, größter, größtes}, {Öl, Öle, Ölen, Öles}, {Platz, Platze, Plätze, Plätzchen, Platzes}, {studier, studiere, studieren, studierend, studierest, studieret, studierst, studiert, studierten, studiertest, studiertet}, {Tag, Tage, Tagen, Tages, Tags}.

*Example 5.22.2.* To further test our algorithm against German irregular verbs with various vowel alternation patterns across different tenses and moods, we throw the following list:

befahl, befähle, befahlen, befählen, befählest, befählet, befahlst, befahlt, befähle, befehlen, befahlend, befahlest, befahlet, befählt, befiehlst, befiehlt, befohlen — “command”;  
begann, begänne, begannen, begännen, begännest, begännet, begannst, begannt, beginn, beginne, beginnen, beginnend, beginnest, beginnet, beginnst, beginnt, begönne, begonnen, begönnen, begönnest, begönnet — “begin”;  
beiß, beiße, beißen, beißend, beißest, beißet, beißt, biss, bisse, bissen, biskest, bisst, gebissen — “bite”;  
blas, blase, blasen, blasend, blasest, blaset, blast, bläst, bliß, bliese, bliesen, bliest, bliest, geblasen — “blow”;  
brach, bräche, brachen, brächen, brächest, brächet, brachst, bracht, breche, brechen, brechend, brechest, brechet, brecht, brich, brichst, bricht, gebrochen — “break”;  
brät, brate, braten, bratend, bratest, bratet, brätst, briet, briete, brieten, brietest, brietet, gebraten — “fry”;  
fahr, fahre, fahren, fahrend, fahrest, fahret, fährst, fahrt, fährt, fuhr, führe, fuhren, führen, führrest, fühhret, fuhrst, fuhr, gefahren — “drive”;  
fall, falle, fallen, fallend, fallest, fallet, fällt, fällt, fiel, fiele, fielen, fielest, fielet, fielst, fielt, gefallen — “fall”;  
fand, fände, fanden, fänden, fandest, fändest, fandet, fändet, finde, finden, findend, findest, findet, gefunden — “find”;  
fechte, fechten, fechtend, fechtest, fechtet, ficht, fichtst, focht, föchte, fochten, fochtest, föchtest, fochtet, föchtet, gefochten — “fence”;  
flieg, fliege, fliegen, fliegend, fliegest, flieget, fliegt, flog, flöge, flogen, flögen, flög, flög, flög, flög, flög, geflogen — “fly”;  
fraß, fräße, fraßen, fräßen, fräbest, fräbet, fraßt, fresse, fressen, fressend, fressest, fresset, fressst, friss, frisst, gefressen — “feed”;  
galt, gälte, galten, gälten, galtest, gältest, galtet, gältet, gegolten, gelte, gelten, geltend, geltet, geltet, gilt, giltst, gölte, gölten, göltest, göltet — “count”;  
geglichen, gleich, gleiche, gleichen, gleichend, gleichest, gleichet, gleichst, gleich, glich, gliche, glichen, gleichest, gleichst, gleich — “resemble”;  
gegriffen, greif, greife, greifen, greifend, greifest, greifet, greift, griff, griffe, griffen, griffest, griffet, griffst, griff — “grab”;  
gehangen, häng, hänge, hängen, hängend, hängest, hänget, hängt, hing, hinge, hingen, hingest, hinget, hingst, hing — “hang”;  
gehoben, haben, heb, hebe, heben, hebend, hebest, hebet, hebst, hebt, hob, höbe, hoben, höben, höbest, höbet, hobst, hobt, hübe, hüben, hübest, hübet — “heave”;  
gekrochen, kriech, krieche, kriechen, kriechend, kriechest, kriechet, kriechst, kriecht, kroch, kröche, krochen, kröchen, kröch, kröch, kröch, kröch — “creep”;  
gelaufen, lauf, laufe, laufen, laufend, laufest, laufet, läuft, lief, liefe, liefen, liefest, liefet, lieft, lieft — “walk”;  
gelesen, las, läse, lasen, läsen, läsest, läset, last, lese, lesen, lesend, lesest, leset, lest, lies, liest — “read”;  
genommen, nahm, nähme, nahmen, nähmen, nähmest, nähmet, nahmst, nahmt, nehme, nehmen, nehmend, nehmest, nehm, nehmst, nimmt, nimmt — “take”;  
geritten, reite, reiten, reitend, reitest, reitet, ritt, ritte, ritten, rittest, rittet — “ride”;  
gerufen, rief, riefe, riesen, riesen, riesest, rieset, rieft, rieft, ruf, rufe, rufen, rufend, rufest, rufet, rufst, ruf — “shriek”;

*geschlafen, schlaf, schlafe, schlafen, schlafend, schlafest, schlafet, schläfst, schlaft, schläßt, schlief, schliefé, schließen, schliefest, schliefet, schliefst, schliefs — “sleep”;*  
*geschlossen, schließ, schließe, schließen, schließend, schließest, schließet, schließt, schloss, schlösse, schlossen, schlössen, schlössest, schlössest, schlösset, schlösset, schlosst — “close”;*  
*geschrieben, schreibe, schreiben, schreibend, schreibest, schreibet, schreibst, schreibt, schrieb, schriebe, schrieben, schriebest, schriebet, schriebst, schriebt — “write”;*  
*gesoffen, sauf, saufe, saufen, saufend, saufest, saufet, säufst, sauft, soff, söffe, soffen, söffen, söffest, söffet, soffst, sofft — “drink (of an animal), booze”;*  
*gestoßen, stieß, stieße, stießen, stießest, stießet, stießt, stoß, stoße, stoßen, stoßend, stoßest, stoßet, stoßt, stößt — “shove”;*  
*getroffen, traf, träfe, trafen, träfen, träfest, träfet, trafst, trافت, treffen, treffend, treffest, treffet, trefft, trifft, trifft — “meet”;*  
*gezogen, zieh, ziehe, ziehen, ziehend, ziehest, ziehet, zieht, zog, zöge, zogen, zögen, zögest, zöget, zogst, zogt — “pull”*

into our algorithm, and obtain

{befahl, befähle, befahlen, befählest, befählet, befahlst, befehlt, befehle, befehlten, befehlend, befehlest, befehlet, befehlt, befehlt, befehlt, befohlen}, {begann, begänne, begannen, begännen, begännest, begännet, begannst, begannt, beginn, beginne, beginnen, beginnend, beginnest, beginnet, beginnst, beginnt, begönne, begonnen, begönnen, begönnest, begönnet}, {beiß, beiße, beißen, beißend, beißest, beißet, beißt, biss, bisse, bissen, bissest, bisset, bisst, gebissen}, {blas, blase, blasen, blasend, blasest, blaset, blast, bläst, blies, bliese, bliesen, bliest, bliest, bliest, geblasen}, {brach, bräche, brachen, brächen, brächest, brächet, brachst, bracht, breche, brechen, brechend, brechest, brechet, brecht, brich, brichst, bricht, gebrochen}, {brät, brate, braten, bratend, bratest, bratet, brätst, briet, briete, brieten, brietest, briitet, gebraten}, {fahr, fahre, fahren, fahrend, fahrest, fahret, fährst, fahrt, fährt, fuhr, führe, fuhren, führen, führer, führst, führer, gefahren}, {fall, falle, fallen, fallend, fallest, fallet, fällt, fällt, fiel, fiele, fielen, fielest, fielet, fielst, fielt, gefallen}, {fand, fände, fanden, fänden, fandest, fändest, fandet, fändet, finde, finden, findend, findest, findet, gefunden}, {fechte, fechten, fechtest, fechtest, fechtet, ficht, fichtst, focht, föchte, föchten, fochtest, föchtest, föchtest, gefochten}, {flieg, fliege, fliegen, fliegend, fliegest, flieget, fliegst, fliegt, flog, flöge, flogen, flögen, flögest, flöget, flogst, flogt, geflogen}, {fraß, fräße, fraßen, fräßen, fräbest, fräbet, fraßt, fresse, fressen, fressend, fressest, fresset, friss, frisst, gefressen}, {galt, gälte, galten, gälten, galtest, gältest, galtet, gälter, gegolten, gelte, gelten, geltend, geltet, geltet, gilt, giltst, gölte, gölten, göltest, göltet}, {geglichen, gleich, gleiche, gleichen, gleichend, gleichest, gleichet, gleichst, gleich, glich, gliche, glichen, glichest, glichst, glicht}, {gegriffen, greif, greife, greifen, greifend, greifest, greifet, greift, griff, griffe, griffen, griffest, griffet, griffst, grift}, {gehangen, häng, hänge, hängen, hängend, hängest, hängst, hängt, hing, hinge, hingen, hingest, hinget, hingst, hingt}, {gehoben, haben, heb, hebe, heben, hebend, hebest, hebet, hebst, heb, hob, höbe, hoben, höben, höbest, höbet, hobst, hobt, hübe, hüben, hübest, hübet}, {gekrochen, kriech, krieche, kriechen, kriechend, kriechest, kriechet, kriechst, kriecht, kroch, kröche, krochen, kröchen, kröchest, kröchet, krochst, krocht}, {gelaufen, lauf, laufe, laufen, laufend, laufest, laufet, läufst, lauft, lief, liefe, liefen, liefest, liefet, liefst, lieft}, {gelesen, las, läse, lasen, läsen, läsest, läset, last, lese, lesen, lesend, lesest, leset, lest, lies, liest}, {genommen, nahm, nähme, nahmen, nähmen, nähmest, nähmet, nahmst, nahmt, nehme, nehmen, nehmend, nehmest, nehmt, nimm, nimmst, nimmt}, {geritten, reite, reiten, reitend, reitest, reitet, ritt, ritte, ritten, rittest, ritett}, {gerufen, rief, riefe, riefen, riefest, riefet, rieft, rieft, ruf, rufe, rufen, rufend, rufest, rufet, rufst, ruf}, {geschlafen, schlaf, schlafe, schlafen, schlafend, schlafest, schlafet, schläfst, schläft, schlief, schliefen, schliefest, schliefet, schliefst, schliefs}, {geschlossen, schließ, schließe, schließen, schließend, schließest, schließet, schließt, schloss, schlösse, schlossen, schlössen, schlössest, schlössest, schlösset, schlösset, schlosst}, {geschrieben, schreibe, schreiben, schreibend, schreibest, schreibet, schrieb, schriebe, schrieben, schriebest, schriebet, schriebt}, {gesoffen, sauf, saufe, saufen, saufend, saufest, saufet, säufst, sauft, soff, söffe, soffen, söffen, söffest, söffet, soffst, sofft}, {gestoßen, stieß, stieße, stießen, stießest, stießet, stießt, stoß, stoße, stoßen, stoßend, stoßest, stoßet, stoßt, stößt}, {getroffen, traf, träfe, trafen, träfen, träfest, träfet, trafst, trافت, treffen, treffend, treffest, treffet, trefft, trifft, trifft}, {gezogen, zieh, ziehe, ziehen, ziehend, ziehest, ziehet, zieht, zog, zöge, zogen, zögen, zögest, zöget, zogst, zogt}

It might be noted that our input list above also includes some irregular subjunctive forms, which are rarely encountered in modern documents.

### 5.2.3 Heuristic detection of compounds

The following algorithm for heuristic detection of German compounds differs from the Danish version (Algorithm 5.12) only in some specific details. To make the context clear, we still state the algorithm in full. (In what follows, the string minus operation  $\hat{\beta} \ominus \hat{\alpha}$  is prescribed by Definition 5.11.)

**Algorithm 5.23** (Heuristic identification of German binary compounds). *Let  $\Lambda^{\hat{\beta}} = \{\hat{\rho}_1, \dots, \hat{\rho}_Q\}$  be a list of distinct German essential roots (without vowel blotting) that contain at least one instance of  $\mathbf{V} = \{a|e|\hat{e}|i|\hat{i}|o|u\}$  and DO NOT match the following string patterns:*

*(ab|al|am|ar|arg|bar|be|beg|bek|bek|én|fal|ge|gros|isç|los|mal|man|miss|naç|nis|sam|sçaf|sçäft|tal|u|un|ung|ver|vers|xüber).*

*The the output of the function CpdDet( $\Lambda^{\hat{\beta}}$ ) is obtained through the following procedures:*

- (1) *Construct a list  $\{(\hat{\rho}_1, \lambda_1^{\hat{\beta}}), \dots, (\hat{\rho}_Q, \lambda_Q^{\hat{\beta}})\}$  where  $\lambda_q^{\hat{\beta}} = \{\hat{\rho}_{(q,1)}, \dots, \hat{\rho}_{(q,n_q)}\}$  is a subset of  $\Lambda^{\hat{\beta}}$  whose members all match the string pattern  $\hat{\rho}_{q\sim}$ , for  $q \in \mathbb{Z} \cap [1, Q]$ .*
- (2) *Expand the aforementioned entry  $(\hat{\rho}_q, \lambda_q^{\hat{\beta}})$  into a list of triplets  $\{(\hat{\rho}_{(q,1)}, \hat{\rho}_q, \hat{\rho}_{(q,1)} \ominus \hat{\rho}_q), \dots, (\hat{\rho}_{(q,n_q)}, \hat{\rho}_q, \hat{\rho}_{(q,n_q)} \ominus \hat{\rho}_q)\}$  for every  $q \in \mathbb{Z} \cap [1, Q]$  such that  $\lambda_q^{\hat{\beta}} \neq \emptyset$ . Collect all these triplets as one runs through the list  $\{(\hat{\rho}_1, \lambda_1^{\hat{\beta}}), \dots, (\hat{\rho}_Q, \lambda_Q^{\hat{\beta}})\}$ . The list of these triplets  $\{(\hat{\rho}_{(1)}, \hat{\eta}_{(1)}, \hat{\rho}_{(1)} \ominus \hat{\eta}_{(1)}), \dots, (\hat{\rho}_{(Q')}, \hat{\eta}_{(Q')}, \hat{\rho}_{(Q')} \ominus \hat{\eta}_{(Q')})\}$  contains potentially valid decompositions of compounds.*
- (3) *Screen the aforementioned list of triplets as follows: for every  $q' \in \mathbb{Z} \cap [1, Q']$ , if  $(\hat{\rho}_{(q')}, \hat{\eta}_{(q')}, \hat{\tau}_{(q')}) = \hat{\rho}_{(q')} \ominus \hat{\eta}_{(q')}$  satisfies*

$$\ell(\hat{\rho}_{(q')} \ominus \hat{\eta}_{(q')}) \geq 2 \quad \text{AND} \quad \hat{\rho}_{(q')} \ominus \hat{\eta}_{(q')} = \mathbf{X}_1(a|e|\hat{e}|i|\hat{i}|o|u)\mathbf{X}_2,$$

*then construct  $\hat{\tau}_{(q')}^*$  by performing  $(e|n|s)\sim \rightarrow \emptyset$  on  $\hat{\tau}_{(q')}$  and  $\hat{\tau}_{(q')}^{**}$  by doing  $(en|er|es|ge)\sim \rightarrow \emptyset$  on  $\hat{\tau}_{(q')}$ , before generating a list  $\lambda_{(q')}^{\hat{\tau}}$  by members of  $\Lambda^{\hat{\beta}}$  that match the pattern  $(\hat{\tau}_{(q')}|\hat{\tau}_{(q')}^*|\hat{\tau}_{(q')}^{**})$ ; otherwise, set  $\lambda_{(q')}^{\hat{\tau}} = \emptyset$ .*

- (4) *Collect all the triplets  $(\hat{\rho}_{(q')}, \hat{\eta}_{(q')}, \lambda_{(q')}^{\hat{\tau}})$  where  $\lambda_{(q')}^{\hat{\tau}}$  is non-void and  $\hat{\tau}_{(q')}$  DOES NOT match the following string patterns:*

$$(h\acute{e}t\mathbf{X}|liç\mathbf{X}|k\acute{e}t\mathbf{X}|um|ung\mathbf{X}).$$

*This list of triplets CpdDet( $\Lambda^{\hat{\beta}}$ ) contains the heuristic decompositions of all the identified binary compounds.*

**Algorithm 5.24** (Approximate clustering of German words with heuristic detection of compounds). *The procedure runs essentially the same way as Algorithm 5.24, except that German rules replace Danish rules.*

*Example 5.24.1.* Testing the algorithm above against the combined inputs from Examples 5.22.1 and 5.22.2, we obtain the following result:

{Apfel, Äpfel, Äpfelchen, Äpfeln, Apfels}, {Arbeit, arbeite, arbeiten, Arbeiten, Arbeiter, Arbeitern, Arbeiters, arbeitest, arbeitet, arbeitete, arbeiteten, Arbeitsplatz, Arbeitsplätze, Arbeitsplätzen, Arbeitsplatzen, Arbeitsplatzes, Arbeitstag, Arbeitstage, Arbeitstagen, Arbeitstages, Arbeitstags, gearbeitet}, {Arbeitsplatz, Arbeitsplätze, Arbeitsplätzen, Arbeitsplätzen, Arbeitsplatzes, Platz, Platze, Plätze, Plätzchen, Platzes}, {Arbeitstag, Arbeitstage, Arbeitstagen, Arbeitstages, Arbeitstags, Tag, Tage, Tagen, Tages, Tags}, {befahl, befähle, befahlen, befählen, befählest, befählet, befahlst, befählt, befehle, befehlen, befehlend, befehlest, befehlt, befehlst, befehlt, befohlen}, {begann, begänne, begannen, begännen, begännest, begännet, begannst, begannt, beginn, beginne, beginnen, beginnend, beginnest, beginnet, beginnst, beginnt, begönne, begonnen, begönnen, begönnest, begönnet}, {beiß, beiße, beißen, beißend, beißest, beißet, beißt, biss, bisse, bissen, bissest, bisset, gebissen}, {blas, blase, blasen, blasend, blasest, blaset, blast, bläst, bliese, bliesen, bliest, bliest, geblasen}, {brach, bräche, brachen, brächen, brächest, brächet, brachst, bracht, breche, brechen, brechend, brechest, brechet, brecht, brich, brichst, bricht, gebrochen}, {brät, brate, braten, bratend, bratest, bratet, brätst, briet, briete, brieten, briest, brietet, gebraten}, {Bus, Busse, Buße, Bussen, Busses}, {dunkel, dunkelste, dunkelstem, dunkelsten, dunkelster, dunkelstes, dunkle, dunklem, dunklen, dunkler, dunklere, dunklerem, dunkleren, dunklerer, dunkleres}, {Ei, Eie, Eier, Eiern, Eies, Eigelb, Eigelbe, Eigelben, Eigelbs}, {Eigelb, Eigelbe, Eigelben, Eigelbs, gelb, gelbe, gelbem, gelben, gelber, gelbere, gelberem, gelberen, gelberer, gelberes, gelbes, gelbste, gelbstem, gelbsten, gelbster, gelbstes}, {Eis, Eise, Eises}, {fahr, fahre, fahren, fahrend, fahrest, fahret, fährst, fahrt, fährt, fuhr, führe, fuhren, führen, fühest, führet, führst, fahrt, gefahren}, {fall, falle, fallen, fallend, fallest, fallet, fällst, fallt, fällt, fiel, fiele, fielen, fielest, fielet, fielst, fielt, gefallen}, {sand, sände, sanden, fänden, fandest, fändest, fandet, fändet, finde, finden, findend, findest, findet, gefunden}, {fechte, fechten, fechtend, fechtest, fechtet, ficht, fichtst, focht, föchte, fochten, föchten, fochtest, föchtest, fochtet, gefochten}, {flieg, fliege, fliegen, fliegend, fliegest, flieget, fliegt, flog, flöge, flögen, flög, flöget, flog, geflogen}, {fraß, fräße,

*fraßen, fräßen, fräβest, fräβet, fraßt, fresse, fressen, fressend, fressest, fresset, fressst, friss, frisst, gefressen}, {frei, freiem, freien, freier, freiere, freierem, freiereren, freierer, freieres, Freiheit, Freiheiten}, {freie, freiste, frei-stem, freisten, freister, freistes}, {galt, gälte, galten, gälten, galtest, gältest, galtes, gältes}, {geltend, geltend, geltet, geltet, gilt, giltst, gölte, göltens, göltet, göltet}, {geglichen, gleich, gleiche, gleichen, gleichend, gleichest, gleichet, gleichst, gleich, glich, gliche, glichen, glichest, glichet, glichst, glich}, {gegriffen, greif, greife, greifen, greifend, greifest, greifet, greift, griff, griffe, griffen, griffest, griffet, griffst, grift}, {gehangen, häng, hänge, hängen, hängend, hängest, hänget, hängt, hängt, hing, hinge, hingen, hingest, hinget, hingst, hingt}, {gehaust, haus, Haus, Häuschen, hause, Hause, hausen, hausend, Häuser, Hauses, haust, hauste, hausted}, {gehoben, haben, heb, hebe, heben, hebend, hebest, hebet, hebst, heb, hob, höbe, hoben, höben, höbest, höbet, hobst, hobt, hübe, hüben, hübest, hübet}, {gekrochen, kriech, krieche, kriechen, kriechend, kriechest, kriechet, kriechst, kriecht, kroch, kröche, krochen, kröchest, kröchet, krochst, krocht}, {geküsst, Kuss, Kuß, Küschen, Kusse, küsse, Küsse, küssen, Küssen, Kusses, küsst, küsst, küsstest, küsstet}, {gelaufen, lauf, laufe, laufen, laufend, laufest, laufet, läuft, lief, liefe, liefen, liefest, liefet, liefst, lieft}, {gelesen, las, läse, lesen, läsen, lässt, lässt, last, lese, lesen, lesend, lesest, leset, lest, lies, liest}, {genommen, nahm, nähme, nahmen, nähmen, nähmest, nähmet, nahmst, nahmt, nehme, nehmen, nehmend, nemest, nemet, nehmst, nimm, nimmst, nimmt}, {geritten, reite, reiten, reitend, reitest, reitet, ritt, ritte, ritten, rittest, rittet}, {gerufen, rief, riefe, riefen, riefest, riefet, rieft, rieft, ruf, rufe, rufen, rufend, rufest, rufet, rufst, ruft}, {geschlafen, schlaf, schlafe, schlafen, schlafend, schlafest, schlafet, schläfst, schläft, schlief, schliefen, schliefest, schliefet, schliefst, schliefst}, {geschlossen, schließ, schließe, schließen, schließend, schließest, schließet, schließt, schloss, schlösse, schllossen, schlössen, schlössest, schlösset, schlösset, schlösset, schlösset}, {geschrieben, schreibe, schreiben, schreibend, schreibest, schreibet, schreibst, schrieb, schriebe, schrieben, schriebest, schriebet, schriebst, schriebt}, {gesoffen, sauf, saufe, saufen, saufend, saufest, saufet, säufst, sauft, säuft, soff, söffe, soffen, söffen, söffest, söffet, soffst, sofft}, {gestoßen, stieß, stieße, stießen, stießest, stießet, stießt, stoß, stoße, stoßen, stoßend, stoßest, stoßet, stoßt, stößt}, {getroffen, traf, trafe, traßen, träfest, träfet, trafst, trast, treffe, treffen, treffend, treffest, treffet, trifft, trifft, trifft}, {gezogen, zieh, ziehe, ziehen, ziehend, ziehest, ziehet, ziehst, zieht, zog, zöge, zogen, zögen, zögést, zöget, zogst, zogt}, {gross, groß, große, großem, großen, großer, größere, größerem, größerer, größerer, größeres, großes, größte, größtem, größten, größter, größtes}, {Öl, Öle, Ölen, Öles, Öls}, {studier, studiere, studieren, studierend, studierest, studieret, studierst, studiert, studierte, studierten, studiertest, studiertet}.*

Here, the compounds *Arbeitsplatz* “workplace”, *Arbeitstag* “workday” and *Eigelb* “egg yolk” are correctly dissolved into their respective constituting components. In the meantime, we note that *bisse* “bite (a subjunctive form)” is not confused with the noun *Busse* “buses (nominative, genitive or accusative)”, because the vowel alternation pattern [i, u] is disallowed in German tags (effective spellings or essential roots), unlike the case in English. Actual German verb conjugations do exhibit alternations between the vowels *i* and *u*, when the vowel in question is followed by *nd*, *ng* or *nk*. Our Algorithm 5.15(3) takes care of this by relabeling *i* in such scenarios as *î*, so that the vowel alternation pattern [î, u] will be later permitted in Algorithm 5.20.

*Example 5.24.2.* In Fig. S5, we further apply the aforementioned word clustering algorithm to topic extraction and machine translation (see Table S1 for text sources).

According to German orthography, diacritical marks in the letters *ä*, *ö*, *ü* are ignored during alphabetical sorting. This rule, if implemented literally, sometimes reduces readability of certain stacked vowels in our word clusters. Therefore, in Fig. S5, we make a compromise: the letters *ä*, *ö*, *ü* are sorted after *z* during alphabetization, as in Swedish orthography.

It should be noted that we have used a German version for *Jane Eyre* that contains pre-1902 spellings. For example, in Fig. S5a’–B’, we see *Thee* “tea” and *Thür* “door” instead of their modern counterparts *Tee* and *Tür*. Fortunately, our stemming algorithm is not adversely affected by such archaic orthography.

MR ELIZABETH SAGTE GUTE DARCY S HEN GESETZ GEGEN KÖMME LIEBE  
 MRS BENNET GAB BINGLEY WEISSEN HALTEN GEPRÄCHEN JANE L LÄSSEN SCHWESTER  
 NEHMEN MISS FELLEN MANN MACHTE STÄDTE TAG NÄCHSTEN WICKHAM FÄNDEN  
 GEDANKEN FREUDEN LÄSSEN FRAU FRAGE EIGENEN HOFFEYAHN  
 AUGEN RIEP WUNSCH VATER BRÖHNS FAMILIE LADY vONTE ANTWORTETE SCHREIBEN  
 SCHLIESSLICH DAMEN CATHERINE SPIEL AVSDRUCK STIMMEN GLAUBE BTTE LEBEN  
 GLÜCK BALD ERWARTETE GLEICH LONDON ZIMMER SICHER AUFMERKSAMKEIT ERZÄHLT BESUCH WEIT  
 KURZES VERLÄSSEN KENNEN WIRKLICH ÜBERZEUG LIZZY WORT RECHT ABEND TANTE MAL NATÜRLICH  
 SETZTE VORLÄPPEN LETZTES SCHÖNHEIT NAHE NACHSTEN GLÜCKLICH LONGBOURN BEMERKUNG  
 VERHALTEN JAHR GENUA MADCHEN WIEDERHÖLLE FREUDE GERN LEID WOHNEN WEISE ZWEI LUCAS TANZEN  
 REISE FRÜHER ABSICHT STOLZ EINLADUNG HOFFNUNG VÖLLIG NETHERFIELD MOGLICH LEICHT ÜBERRASCHEN KITTY ASTRAEG  
 VERSUCHEN LÄSER WEGE HERZEN ERNST LÄCHELN REICH VERSTEHEN BESONDERS BENEHMEN ERWIDERE EHRE OBERST KIND SCHWÄGEN  
 ERFAHRENEN ERGÜGEN KAPITÄN ABSCHIED UNTERHALTEN SORGE NAMEN ERINNERN MEINUNG SÄMMLN BESITZ VORSTELLEN ENTSTREBT  
 SIR VERBAUS VERSPRECHEN GELEGENHEIT ERKLÄRUNG REDEN HALBE WANDTE URTEIL SCHLECHT  
 ERSCHINNEN ZWEIFEL THEM DREI STUNDE ÜBERBLAUF KUTSCHE BEOKCHTEN AFANG MNHEM HAND PAA MORGEN BESTIMMT BALL S99  
 GESICHT FREI FOLGEN ENTSCHEIDUNG INSICHT GEFÄHR ENTSCHEIDUNG ESSEN ANGST ROSINS SCHICK ALLEGHENY VERLIEBT  
 FREUNDLICHE TRUST BEGANN WELT SOFORT JEDENFALLS ALTBURG SCHLICHT VERSICHERN ERWÄHNT VERHEIRATEL TUR SCHWÄGER WILLIAM FORSTER BARONIN FEST  
 ERFAHRENEN VERMUTET SEITE VOLLE BESCHAFTIG ARM ANGEGELENTET KOPF FRM EREIGNIS EMPFANG MINITES PERCH  
 ALLERGÜGEN PLÄNE BESUCHER BAUSCHEN OFFIZIERE BEGRÜNLICHE BIRCH BEZEICHNUNG STÄNDIG VERBUND  
 FORT DE ALTEN LÄUTE HAUTPLATZ REIZEN SPARSAM VERDIENST VERBRACHT NEUGER CHARAKTER MARY HERTFORDSHIRE ANKUNFT ZWEITES  
 TISCH COUSINE FÄLT EHE EHE DINGE OH ALTER RAUME FEHLAUFZUGEN MUHE ART VERHÄLTEN SYTIZ TIEF PERSOON  
 HURST DAUCH ZÄHLER HÄNDERL FERTIG FERDNER FERDNER FERDNER FERDNER FERDNER FERDNER FERDNER  
 HEUTE UNTERBRUCHEN BESTÄTIGT GENERALIN UNSO MEIST ZERBET GAST FUNF BEDARF GELD VONANGFÜHL  
 BUCH GESTERN HOMM PÄPPEL FERDNER FERDNER FERDNER FERDNER FERDNER FERDNER FERDNER FERDNER  
 MINDE BRIGHTON KLIMA KLEINER ANPREDNEN KLEINER STÄRK ARBELT SAMLE FRIEDE BERICHT WILIG  
 VERLEGENHEIT SPAR BESUCHER AUFMERKSAMKEIT RECHT HUNTSCHEN KRONEN KLEINER KLEINER KLEINER  
 WEIT GESCHICHT TESTIM BIBLIOTHEK VERTRÄG KOSTE UNTERSUCHT GESUNDHEIT DOMIN VERTRETEN VERSTÄNDIG  
 BESTIMP BELEIDIG BEHÄNDIG VERBUND VERDACH SOHN HOFFENTLICH STRAKE ZIEHEN PARTNER N ENTHUSIASM  
 VERGANGENHEIT VERBUND VERDACH GESUNDHEIT DOMIN VERTRETEN VERSTÄNDIG  
 (a)

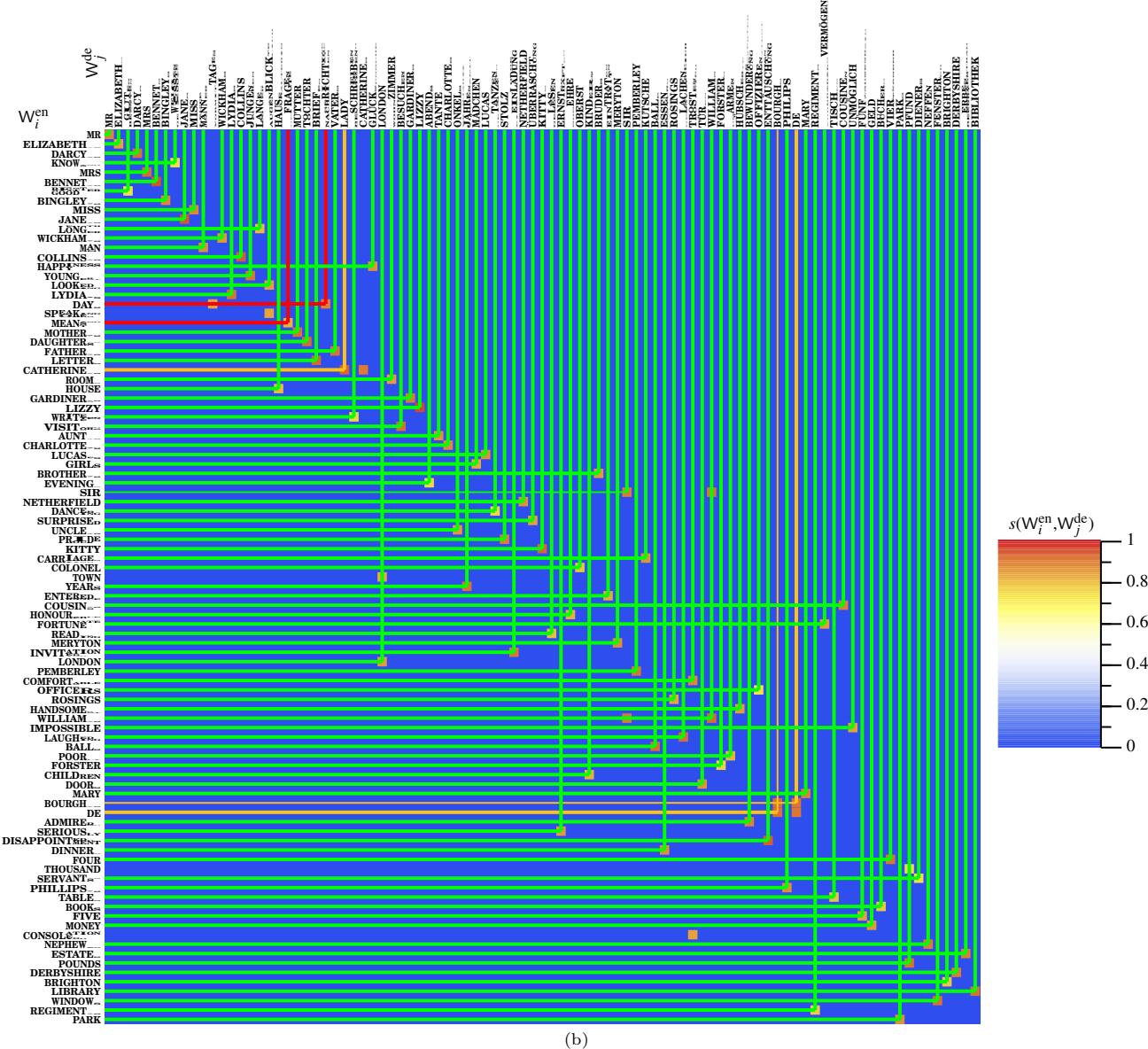


Fig. S5. Text mining in German. (a) Statistically identified topics ( $n_{ii} \geq 20$ ) in a German version of *Pride and Prejudice*, with the same color encoding scheme as Fig. S3. (b) Semantic similarities  $s(W_i^{en}, W_j^{de})$  between selected topics in English and German versions of *Pride and Prejudice*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) ranking of the similarity score for the optimal node. Green (resp. red) cross-hair indicates a correct (resp. incorrect) match. Amber cross-hair marks a link between distinct concepts that share the same hypernyms.

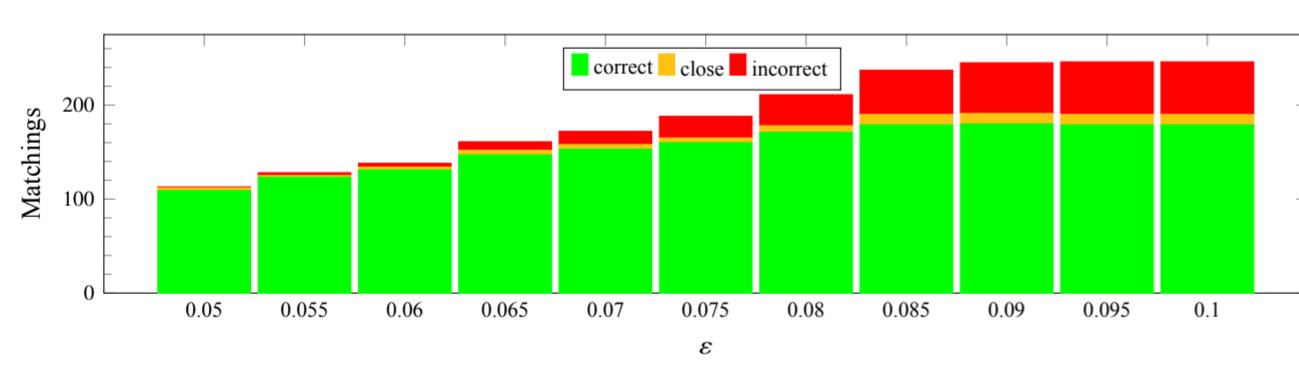
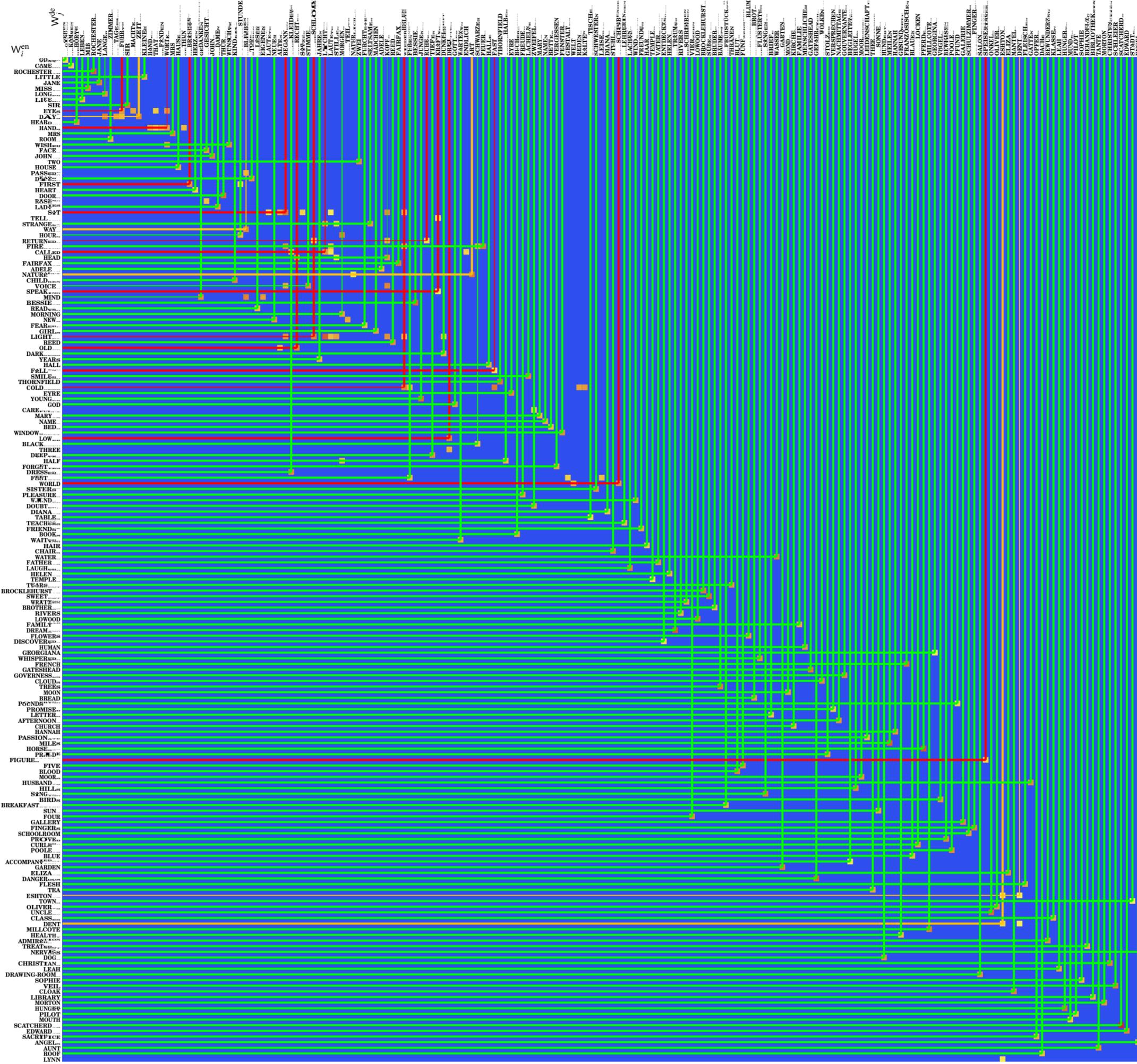


Fig. S5. Text mining in German. (Continued) (a') Statistically identified topics ( $n_{ii} \geq 20$ ) in a German version of *Jane Eyre*, with the same color encoding scheme as Fig. S3. (b') Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{de}})$  between selected topics in English and German versions of *Jane Eyre*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) rank of the similarity score for the optimal node. *Green* (resp. *red*) cross-hair indicates a correct (resp. incorrect) match. *Amber* cross-hair marks a link between distinct concepts that share the same hypernyms. (c') Results from control experiments with different choices of the  $\alpha$  parameter in hallmark screening criteria (1,12).

### 5.3 Modified Porter stemming algorithm for Dutch

**Definition 5.25** (Dutch stop words). If a word belongs to the following list<sup>64</sup>:

'k, 'n, 'ne, 'ns, 's, aan, achter, achteruit, af, al, allang, alle, allebei, allemaal, allen, alles, als, altijd, ander, andere, anders, behalve, beide, beiden, ben, beneden, bent, boven, bovendien, bij, bijna, d'n, daar, daardoor, daarin, daarmee, daarna, daarnaast, daarom, daaronder, daarop, daartoe, daarvan, dan, dat, datgene, datgenen, de, deden, deed, deedt, degene, degenen, den, der, dergelijk, dergelijke, dergelijks, des, desondanks, deze, dezelfde, dezen, dezer, dezes, die, diegene, diegenen, dien, diens, dier, dikwijs, dit, doch, doe, doen, doend, doet, door, dus, echt, echter, een, eenieder, eenmaal, eens, eigenlijk, elders, elk, elkaar, elkaars, elke, en, ene, enen, ener, enig, enige, enkel, enkele, enkels, er, erg, ergens, erin, ermee,ernaast, erom, eronder, erop, ertoer, ervan, even, eveneens, evengoed, evenmin, eventjes, eventueel, evenwel, evenzeer, ge, gedaan, geen, gehad, gekund, gemoeten, genoeg, geweest, geworden, gj, haar, had, hadden, hadt, hare, heb, hebbe, hebben, hebbend, hebt, heeft, heel, heen, hele, helemaal, hem, hen, het, hetgeen, hetwelk, hetzelfde, hier, hierin, hiermee, hiernaast, hierom, hieronder, hierop, hier toe, hiervan, hiervoor, hoe, hoelang, hoeveel, hoewel, hun, hunne, hunner, hij, ieder, iedere, iedereen, iemand, iets, ik, in, intussen, is, ja, je, jezelf, jou, jouw, jouwe, jouver, jullie, jij, jijzelf, k, kan, kon, konden, kondt, kunnen, kunnend, kunt, later, luttel, luttele, maar, mag, me, mee, meer, meest, meeste, men, menig, menige, met, meteen, minder, mindere, minderen, minst, minste, minstens, misschien, mocht, mochte, mochten, moest, moeste, moesten, moet, moetende, moetend, moge, mogen, mogend, moogt, mij, mijn, mijne, mijner, mijzelf, n, na, naar, naartoe, naast, nauwelijks, nee, neen, nergens, net, niemand, niet, niets, nietemin, noch, nog, nooit, nou, nu, of, om, omdat, onder, onmiddellijk, ons, onszelf, onze, onzer, ooit, ook, op, over, per, reeds, sindsdien, sint, sommige, sommigen, soms, straks, t, te, tegen, ten, tenminste, tenslotte, ter, terug, tevens, toch, toe, toen, tot, trouwens, u, uit, uw, uwe, uwer, uzelf, vaak, vaaks, vaakst, vake, vaker, van, vanaf, vandaan, vanwaar, veel, vele, voor, vooral, vooruit, vroeger, waar, waarheen, waarin, waarmee, waarnaast, waarom, waaronder, waarop, waartoe, waarvandaan, wanneer, want, waren, was, wat, we, weer, wees, weest, weinig, weinige, wel, welk, welke, wellicht, werd, werden, werdt, wezen, wie, wil, word, worden, wordend, wordt, wij, zal, ze, zeer, zelden, zelf, zich, zichzelf, zo, zoals, zojuist, zolang, zonder, zou, zoude, zouden, zoudt, zoveel, zulk, zulke, zulks, zulle, zullen, zullend, zult, zij, zijn, zijnd, zijne, zijner, zijt,

then we consider it a Dutch stop word. All the Dutch stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list.  $\square$

To accommodate to pre-1946 (Flanders)/pre-1947 (Netherlands) Dutch written standard, which employed a sophisticated case system as in modern German (see [https://en.wikipedia.org/wiki/Archaic\\_Dutch\\_declension](https://en.wikipedia.org/wiki/Archaic_Dutch_declension)), we need to add inflected forms of the definite and indefinite articles, as well as those of the demonstratives to the list of stop words. It should be noted that archaic Dutch declensions still survive in some stock phrases in modern Dutch, such as the genitive construction in *het Koninkrijk der Nederlanden* “the Kingdom of Netherlands”.

The 1946/1947 orthography reform changed the endings of certain words: *bosch* “forest”, *mensch* “human” and *visch* “fish” (singular noun) are now spelt *bos*, *mens* and *vis*, respectively [47, §2.1]. Meanwhile, the spellings of adjectival endings in *mathematisch* “mathematical” and *typisch* “typical” are not affected, despite having the same silent *ch* at word final positions. So long as a document consistently employs an orthographical standard, we do not need to worry about these changes in our clustering algorithm, and we will not forcibly change all word final *sch* to *s*.

#### 5.3.1 Effective spelling and essential root

Vowel lengths are indicated in Dutch spellings according to certain rules. For example, all these three words *boom* “tree”, *boomen* “trees (pre-1946/pre-1947)” and *bomen* “trees (post-1946/post-1947)” contain the same long vowel in their first syllable [47, §2.1]. The penultimate step in the following algorithm helps us detect the vowel length in a Dutch spelling heuristically.

**Algorithm 5.26** (Dutch effective spelling). Set  $\mathbf{V} = (a|A|e|\acute{e}|E|i|\acute{i}|o|\acute{o}|O|u|\acute{u}|\ddot{u}|U|y)$ ,  $\mathbf{V}_* = (a|e|o|u)$  and  $\mathbf{V}_0 = (a|e|i|o|u)$ . Set

$\mathbf{c} = (b|c|d|f|g|h|j|k|l|m|n|p|q|r|s|t|v|w|x|z)$  and  $\mathbf{C} = (B|C|D|F|G|H|J|K|L|M|N|P|Q|R|S|T|V|W|X|Z)$ .

For a Dutch word  $\hat{\sigma}$ , its effective spelling  $\text{EffSpell}(\hat{\sigma})$  is constructed in the following steps:

- (1) Convert to lowercase, and replace

<sup>64</sup>Our list of Dutch stop words is based on [snowball.tartarus.org/algorithms/dutch/stop.txt](https://snowball.tartarus.org/algorithms/dutch/stop.txt), with extensive additions that roughly match their counterparts in English. We have also included these words: *d'n*, *den*, *der*, *des*, *deze*, *dezen*, *dezer*, *dezes*, *dien*, *diens*, *dier*, *eens*, *eens*, *ene*, *enen*, *ener*, *'n*, *'ne*, *nen*, *'ns*, *'s*, which are involved in archaic Dutch declensions and records of colloquial speech.

(Ø ge)zegen~	aanwezig <b>X</b>	a <u>f</u> wezig <b>X</b>	bezet~	compl~	diner	eerb~	eerder~	ellend	forst~	gardiner
βλεστ̄n	πρεσεν̄t̄	αβσεν̄t̄	βζ̄et̄	κμλ̄	δινερ̄	ehr̄β̄	vroeḡ	μισρβ̄d̄	φορστ̄	φαρδινερ̄
gelo <b>X</b> ε(oʃv)~	genegen~	gezag	gezicht	jone~	kilo~	kind~	langza~	ma(Ø m)ma(Ø a)~	manier~	
βλλο <b>X</b>	nijgen	αυτρρḡ	φαστ̄	γοῦε	κιλο	κιλδ̄	σλωωza	moeder	μανίερ	
meester	negen~	o(or re)~	pappa(Ø a)~	ped~	plant	schook	schoonh~	tafel	triest~	vier~
μεεστ̄er̄	9e	εαρ̄	vader	peδ̄	plant̄	σχοκ̄	mooih̄	ταβλ̄l̄	τριεστ̄t̄	4er̄
wonder	(Ø ge)rend(Ø e en)		jane	ma(š š atje)(Ø s)		mamme(n tje)(Ø s)	oom(Ø s)	~mooie <b>X</b>		
wonðer	rennen		γιοῦα	moeder		moeder	οῦκλ̄	mooi		

(2) Replace<sup>65</sup>

(Ø ge)hoor	(Ø ge)zegd	X <sup>ε</sup> (bra(aʃv))~	X <sup>ε</sup> (do(d od))-	á	ä	bruik	docht~	é	è	ë	geëe~	geget	georg~
χοο̄r̄	zeg	p <b>X</b>	t <b>X</b>	a	a	βruik̄	dtocht̄	e	e	e	geĒ	et̄	jgeorḡ
gezond	hor̄	ī	janu~	koning	leraar	lie <b>X</b> ε(fv)(Ø st)~	lof	loof	los	lov	moe(Ø ie)~	ó	ö
jgezond	χor̄	ī	jnu	koning	leren	xli <b>X</b>	λof̄	λoof̄	xlos̄	λov̄	μμmoē	o	o
oē	oū	placht̄	su~	tien	tuin	ú	ü	voeld̄	ij̄	zicht̄	ziek~	zom~	zoon~
ó̄	û̄	pleeḡ	xsu	tzien	tuin̄	u	u	föl̄	ȳ	zien̄	xziek̄	ζom̄	xxzoon̄
beter(Ø e s)		<u>sir</u>				~(isme ist)(Ø en tje)(Ø s)				~dje(Ø s)		~ga(Ø an and at)	
gôd̄		ssir			e				d̄			ginḡ	
~ga(afʃ ʃve ven)			~kom(Ø t)			~kwa(amt m me men)				~la(agt g ge gen)		~la(ast s ze zen)	
geven			komen			komen				liggen		lezen	
~sta(Ø an and at)		~uk̄				~wist(Ø e en)				~zat(Ø e en)		~zei(Ø de den dt)	
stond		ukk̄				weten				zitten		zeḡ	

(3) Do *ei* → *ē*, *list* → *λist*, *ui* → *ü*;

(4) Do *ie* → *î*, *aan(Ø|ge)n* → *n*, *door(Ø|ge)r* → *r*, *dwars(Ø|ge)s* → *s*, *om(Ø|ge)m* → *m*, *op(Ø|ge)p* → *p*, *vort(Ø|ge)t* → *t*, *weg(Ø|ge)g* → *g*;

(5) Do *~end* → *e*, *~χ̄₁₄Vs* → *χ̄₁*, *χ̄₂(dom|vol)(e|l|r|s|t)ₘ₀* → *χ̄₂*, *χ̄₃(ery|vaardig)**X*** → *χ̄₃*;

(6) Do *χ̄₁v* → *χ̄₁f*, *χ̄₂z* → *χ̄₂s* before turning double letters into the capital form of the same letter (i.e. *ee* → *E* etc.);

(7) Do *~en* → *e*;

(8) Do *~V\_\*ce* → *V\_\*+c*, *~V₀Ce* → *V₀C\_-*, where *V\_\*+* is the upper case form of *V\_\** and *C\_-* is the lower case form of *C*;

(9) Do *lerares* → *lEr*, *vrE* → *βrE*, *~îr(Ø|es|s|tje)(Ø|s)* → *Ø*, *~An(Ø|en|tje)(Ø|s)* → *a*, *~sje(Ø|s)* → *s*, *~tje(Ø|s)* → *Ø*, *~Tje(Ø|s)* → *t*.

**Definition 5.27** (Dutch protected range). Set  $\mathbf{V}^* = (a|A|\hat{e}|E|i|i|o|\hat{o}|O|u|\hat{u}|i|i|U|y)$ , and write  $\bar{\mathbf{V}}_m$  (resp.  $\bar{\mathbf{V}}_{m_0}$ ) for one (resp. zero) or more repeats of any character other than  $\mathbf{V} = (a|A|\hat{e}|E|i|i|o|\hat{o}|O|u|\hat{u}|i|i|U|y)$ . Define  $\mathbf{C}_{m_0}^{**} = (\overline{a|A|\hat{e}|E|i|i|o|\hat{o}|O|r|u|\hat{u}|i|i|U})_{m_0}$ . Let  $\hat{\sigma}$  be the effective spelling of a Dutch word, its protected range  $\text{ProtRg}(\hat{\sigma}) = \max\{\lambda_1(\hat{\sigma}), \lambda_2(\hat{\sigma})\}$  is determined by two non-negative integers  $\lambda_1(\hat{\sigma})$  and  $\lambda_2(\hat{\sigma})$  specified through the following procedures:

- Look for the string pattern  $((Ø|a|A|be|e|E|ge|her|i|mede|na|O|\hat{o}|u|U|\hat{u}|i|i|ver)\bar{\mathbf{V}}_{m_0}\mathbf{V}^*|\bar{\mathbf{V}}_m(e|E))\mathbf{C}_{m_0}^{**}~$  in the string  $\hat{\sigma}$ ;
- If the string pattern above is found, the last position occupied by such a string defines  $\lambda_1(\hat{\sigma})$ ; otherwise, set  $\lambda_1(\hat{\sigma}) = 0$ ;
- Look for the pattern  $(a|A|o|\hat{o}|O|u|\hat{u}|U)$  in the string  $\hat{\sigma}$ ;
- If a letter in the pattern above is found, the last position occupied by such a letter defines  $\lambda_2(\hat{\sigma})$ ; otherwise, set  $\lambda_2(\hat{\sigma}) = 0$ . □

Unlike Danish, but similar to German, our Dutch essential root is free from common separable verb prefixes (see step (3) in the algorithm below).

<sup>65</sup>Here, the string pattern *ij* may refer to two Latin letters *i* and *j* sitting next to each other, or a single Unicode character LATIN SMALL LIGATURE IJ [U+0133].

**Algorithm 5.28** (Dutch essential root). Set  $\mathbf{V} = (a|A|e|\hat{e}|E|i|\hat{i}|o|\hat{o}|O|u|\hat{u}|U|y)$  and  $\mathbf{C} = \overline{(a|A|e|\hat{e}|E|i|\hat{i}|o|\hat{o}|O|u|\hat{u}|U|y)}$ . Let  $\hat{\sigma}$  be the effective spelling of a Dutch word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:

(1) Break down  $\hat{\sigma} = \hat{\sigma}_1\hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}^{[\text{ProtRg}(\hat{\sigma})]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma})$  is equal to the protected range of  $\hat{\sigma}$ .

(2) On  $\hat{\sigma}_2$ , perform the following substitutions in a sequel:<sup>66</sup>

(2.1) Do  $\sim(Er|s) \rightarrow \emptyset$ ;

(2.2) Do  $\sim e(d|e|m|n|r|st|t)_{m_0} \rightarrow \emptyset$ .

The result after these two steps of operations is called  $\hat{\sigma}'_2$ .

(3) Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ , before doing  $(\emptyset|af|An|by|dOr|dwars|mE|mede|na|om|op|t\hat{o}|vOrt|weg)(\emptyset|ge)\mathbf{X}\mathbf{V}\sim \rightarrow \mathbf{X}\mathbf{V}$ ;

(4) Do  $\sim\mathbf{C}t \rightarrow \mathbf{C}$ ;

(5) Do  $v \rightarrow f, z \rightarrow s$ ;

(6) Do  $\sim\hat{x} \rightarrow \hat{x}_-$  (i.e. reduce final letter to lowercase);

(7) Do  $ag \rightarrow Ag, ak \rightarrow Ak, al \rightarrow Al, at \rightarrow At, breng \rightarrow brach, denk \rightarrow dach, Ef \rightarrow ef, ep \rightarrow Ep, kOp \rightarrow koch, nam \rightarrow nAm, sôk \rightarrow soch$ .

### 5.3.2 Admissible mutation and approximate clustering

The Dutch vowel blotting mechanism is similar to the Danish version (Algorithm 5.5).

**Algorithm 5.29** (Dutch vowel blotting). Set  $\mathbf{V}_m = (a|A|e|\hat{e}|E|i|\hat{i}|o|\hat{o}|O|u|\hat{u}|\ddot{u}|U|y)_m$  and  $\mathbf{C}_{m_0} = \overline{(a|A|e|\hat{e}|E|i|\hat{i}|o|\hat{o}|O|u|\hat{u}|\ddot{u}|U|y)}_{m_0}$ . For a token string  $\hat{\sigma}$ , its blotted form  $\text{BlotV}_1(\hat{\sigma})$  is constructed as follows:

- If the string pattern  $(\emptyset|a|be|e|\hat{e}|i|o|u|\ddot{u})\mathbf{C}_{m_0}\mathbf{V}_m\sim$  can be found in the string  $\hat{\sigma}$ , then the last position occupied by such a pattern is replaced by the letter “a”.
- Otherwise, leave the string  $\hat{\sigma}$  intact.

Similar to what we did in §5.1.2 for the case of Danish, we will construct a bivariate Boolean-valued function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  on a “simple heredity test function” in Algorithm 5.30, and a set of “admissible suffix mismatch and vowel alternation” rules in Algorithm 5.31.

**Algorithm 5.30** (Simple heredity test). Set  $\mathbf{V} = (a|A|e|\hat{e}|E|i|\hat{i}|o|\hat{o}|O|u|\hat{u}|\ddot{u}|U|y)$ . The Boolean-valued function  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if  $\hat{\alpha}$  contains at least one instance of  $\mathbf{V}$ , **AND** at least one of the following six conditions holds:<sup>67</sup>

(i)  $\hat{\alpha} = \hat{\beta}$ ;

(ii)  $\hat{\beta} = \hat{\alpha}d$ ;

(iii)  $\hat{\beta} = \hat{\alpha}Er$ ;

(iv)  $\hat{\beta} = \hat{\alpha}s$ ;

(v)  $\hat{\beta} = \hat{\alpha}t$ ;

(vi)  $\ell(\hat{\beta}) > \ell(\hat{\alpha}) \geq \ell(\hat{\beta}) - 2$  **AND**  $\hat{\alpha} = \hat{\beta}^{[\ell(\hat{\alpha})]}$  **AND**  $\hat{\beta}^{\{\ell(\hat{\alpha})+1\}} = (e|n|s)$ . (See Definition 3.1 for the notations  $\hat{\beta}^{[n]}$  and  $\hat{\beta}^{\{n\}}$ .)

In what follows, we define  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta})$  and  $\text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta})$  as what is done in the Danish case (Algorithm 5.7).

<sup>66</sup>In other words, the core algorithm for essential root extraction runs as follows: keep the last “strong” vowel *a*, *i*, *o* or *u* in non-final position, plus one subsequent letter; delete final *a*; erase the final appearance of *e* and all the letters thereafter.

<sup>67</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical **OR**.

**Algorithm 5.31** (Admissible suffix mismatch and vowel alternation). *For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function*

$$\text{AdmMut}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta}))$$

*returns TRUE if*

$$\begin{aligned} & \text{SuffixNW}(\hat{\alpha}, \hat{\beta}) \\ = & [(\emptyset | (e|E|n|t)_m | ((acht|haf|mat)ig | (\emptyset|b)Ar|e(ling|lyk|nis|r)|h(e|\hat{e})d|i(n|ng|sch)|je|l(ing|yk)|nis|s(Am|e|t)|t(\hat{e}t|je)|yk)\mathbf{X})), \\ & (\emptyset | (e|E|n|t)_m | ((acht|haf|mat)ig | (\emptyset|b)Ar|e(ling|lyk|nis|r)|h(e|\hat{e})d|i(n|ng|sch)|je|l(ing|yk)|nis|s(Am|e|t)|t(\hat{e}t|je)|yk)\mathbf{X})) \end{aligned}$$

**AND** at least one of the following three conditions holds:<sup>68</sup>

- (i)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \emptyset$  **AND** the lowercase form of  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$  contains at least one instance of  $(a|e|\hat{e}|i|\hat{l}|o|\hat{o}|u)$ ;
- (ii)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = ([A, E] | [A, i] | [A, \hat{i}] | [A, \hat{o}] | [A, O] | [e, i] | [e, \hat{i}] | [e, o] | [E, i] | [E, \hat{i}] | [E, O] | [E, y] | [i, o] | [\hat{i}, O] | [O, \hat{i}])$ ;
- (iii)  $\text{NW}^*(\hat{\beta}, \hat{\alpha}) = ([A, E] | [A, i] | [A, \hat{i}] | [A, \hat{o}] | [A, O] | [e, i] | [e, \hat{i}] | [e, o] | [E, i] | [E, \hat{i}] | [E, O] | [E, y] | [i, o] | [\hat{i}, O] | [O, \hat{i}])$ .

Similarly, one can evaluate another Boolean-valued function

$$\text{AdmMut}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta}))$$

by trading all the occurrences of NW in the statements above with SW.

**Algorithm 5.32** (Heredity test function). *The structure of the Dutch heredity test function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  is identical to the German version (Algorithm 8.1.2), except that the functions  $\text{SimpHrdTest}$ ,  $\text{RootNW}$ ,  $\text{SuffixNW}$ ,  $\text{NW}^*$ ,  $\text{RootSW}$ ,  $\text{SuffixSW}$ ,  $\text{SW}^*$  must follow the Dutch rules stated above.*

**Algorithm 5.33** (Approximate clustering of Dutch words). *The algorithm is essentially the same as Algorithm 5.10, except that Dutch rules (instead of Danish rules) apply to all the tags (effective spelling, essential root, vowel blotting etc.).*

*Example 5.33.1.* As an illustration of our clustering algorithm, we pick the following families of Dutch words, where an “approximate translation” in English, enclosed in quotation marks, is appended to the end of each family.<sup>69</sup>

*best, beste, beter, betere, beters, goed, goede, goeden, goeder, goeds — “good”;  
daad, daden — “deed”;  
dooier, dooiers, dooiertje — “yolk”;  
ei, eieren, eitje, eitjes — “egg”;  
eidooier, eidooiers, eidooiertje, eierdooier, eierdooiers, eierdooiertje — “egg yolk”;  
eiwit, eiwitten, eiwittje — “egg white”;  
geleerd, leer, leerde, leerden, leert, lere, leren, lerend — “learn”;  
gelost, los, losse, lossen, lossend, lost, loste, losten — “dump”;  
geraasd, raas, raasde, raasden, raast, raze, razen, razend — “rage”;  
gewerkt, werk, werke, werken, werkend, werkt, werkte, werkten — “work”;  
klein, kleine, kleinen, kleiner, kleinere, kleiners, kleins, kleinst, kleinste — “small”;  
maan, manen, maantje — “moon”;  
man, manne, mannen, mans — “man”;  
vrouw, vrouwe, vrouwen — “woman”;  
wit, wits, witst, witste, witte, wittere, witters — “white”.*

Applying Algorithm 5.33 to this list of words, we obtain the following clustering results:

*{best, beste, beter, betere, beters, goed, goede, goeden, goeder, goeds}, {daad, daden}, {dooier, dooiers, dooiertje}, {ei, eieren, eitje, eitjes}, {eidooier, eidooiers, eidooiertje}, {eierdooier, eierdooiers, eierdooiertje}, {eiwit, eiwitten, eiwittje}, {geleerd, leer, leerde, leerden, leert, lere, leren, lerend}, {gelost, los, losse, lossen, lossend, lost, loste, losten}, {geraasd, raas, raasde, raasden, raast, raze, razen, razend}, {gewerkt, werk, werke, werken, werkend, werkt, werkte, werkten}, {klein, kleine, kleinen, kleiner, kleinere, kleiners, kleins, kleinst, kleinste}, {maan, maantje, man, manen, manne, mannen, mans}, {vrouw, vrouwe, vrouwen}, {wit, wits, witst, witste, witte, wittere, witters}.*

<sup>68</sup>We note that  $\text{NW}^*(\hat{\alpha}, \hat{\beta})$  and  $\text{NW}^*(\hat{\beta}, \hat{\alpha})$  differ only in the order of the two components in the bracket.

<sup>69</sup>This list includes archaic declension forms found in the paradigms on the following web-page: [en.wikipedia.org/wiki/Archaic\\_Dutch\\_declension](https://en.wikipedia.org/wiki/Archaic_Dutch_declension).

We note that two alternative forms of “egg yolk” *eidooier* and *eierdooier* are split into two different clusters in this method.

*Example 5.33.2.* To further test our algorithm against Dutch irregular verbs with various vowel alternation patterns across different tenses and moods, we throw the following list:

*bak, bakke, bakken, bakkend, bakt, bakte, bakten, gebakken* — “bake”;  
*ban, bande, banden, banne, bannen, bannend, bant, gebannen* — “expel”;  
*bederf, bederft, bederve, bederven, bedervend, bedierf, bedierft, bedorve, bedorven* — “spoil”;  
*bedrieg, bedriege, bedriegen, bedriegend, bedriegt, bedroge, bedrogen, bedroog, bedroogt* — “deceive”;  
*beet, bete, beten, bijt, bijte, bijten, bijtend, gebeten* — “bite”;  
*begin, beginne, beginnen, beginnend, begint, begon, begonne, begonnen, begont* — “begin”;  
*bied, biede, bieden, biedend, biedt, bode, boden, bood, boodt, geboden* — “bid”;  
*bind, binde, binden, bindend, bindt, bond, bonde, bonden, bondt, gebonden* — “tie”;  
*blaas, blaast, blaze, blazen, blazend, blijs, bliest, blieze, bliezen, geblazen* — “blow”;  
*bleek, bleekt, bleke, bleken, blyk, blyke, blyiken, blykend, blykt, gebleken* — “appear”;  
*boge, bogen, boog, boogt, buig, buige, buigen, buigend, buigt, gebogen* — “bend”;  
*braakt, brak, brake, braken, breek, breekt, breke, breken, brekend, gebroken* — “break”;  
*bracht, brachte, brachten, breng, brenge, brengen, brengend, brengt, gebracht* — “bring”;  
*dacht, dachte, dachten, denk, denke, denken, denkend, denkt, gedacht* — “think”;  
*draag, draagt, drage, dragen, dragend, droeg, droege, droegen, droegt, gedragen* — “carry”;  
*droop, droopt, drope, dropen, druip, druipe, druipen, druipend, druipt, gedropen* — “drip”;  
*gegouden, geld, gelde, gelden, geldend, geldt, gold, golde, golden, goldt* — “apply”;  
*geheven, hef, heffe, heffen, heffend, heft, hief, hieft, hieve, hieven* — “heave”;  
*geholpen, help, helpe, helpen, helpend, helpt, hielpt, hielpe, hielpen, hielpt* — “help”;  
*gekocht, kocht, kochte, kochten, koop, koopt, kope, kopen, kopend* — “buy”;  
*gelegen, laagt, lag, lage, lagen, lig, ligge, liggen, liggend, ligt* — “lie”;  
*genomen, naam, nam, name, namen, neem, neemt, neme, nemen, nemend* — “take”;  
*geschapen, schep, scheppen, scheppend, schept, schiep, schiepen, schiept* — “create”;  
*gescholden, scheld, schelde, schelden, scheldend, scheldt, schold, scholden, scholdt* — “scold”;  
*gevraagd, vraag, vraagd, vraagde, vraagden, vraagt, vrage, vragen, vragend, vroeg, vroege, vroegen, vroegt* — “ask”;  
*gewogen, weeg, weegt, wege, wegen, wegend, woge, wogen, woog, woogt* — “weigh”;  
*gezeten, zat, zate, zaten, zit, zitte, zitten, zittend* — “sit”;  
*gezocht, zocht, zachte, zochten, zoek, zoeke, zoeken, zoekend, zoekt* — “seek”

into our algorithm, and obtain

{*bak, bakke, bakken, bakkend, bakt, bakte, bakten, gebakken*}, {*ban, bande, banden, banne, bannen, bannend, bant, gebannen*}, {*bederf, bederft, bederve, bederven, bedervend, bedierf, bedierft, bedorve, bedorven*}, {*bedrieg, bedriege, bedriegen, bedriegend, bedriegt, bedroge, bedrogen, bedroog, bedroogt*}, {*beet, bete, beten, bijt, bijte, bijten, bijtend, gebeten*}, {*begin, beginne, beginnen, beginnend, begint, begon, begonne, begonnen, begont*}, {*bied, biede, bieden, biedend, biedt, bode, boden, bood, boodt, geboden*}, {*bind, binde, binden, bindend, bindt, bond, bonde, bonden, bondt, gebonden*}, {*blaas, blaast, blaze, blazen, blazend, blijs, bliest, blieze, bliezen, geblazen*}, {*bleek, bleekt, bleke, bleken, blyk, blyke, blyiken, blykend, blykt, gebleken*}, {*boge, bogen, boog, boogt, buig, buige, buigen, buigend, buigt, gebogen*}, {*braakt, brak, brake, braken, breek, breekt, breke, breken, brekend, gebroken*}, {*bracht, brachte, brachten, breng, brenge, brengen, brengend, brengt, gebracht*}, {*dacht, dachte, dachten, denk, denke, denken, denkend, denkt, gedacht*}, {*draag, draagt, drage, dragen, dragend, droeg, droege, droegen, droegt, gedragen*}, {*droop, droopt, drope, dropen, druip, druipe, druipen, druipend, druipt, gedropen*}, {*gegouden, geld, gelde, gelden, geldend, geldt, gold, golde, golden, goldt*}, {*geheven, hef, heffe, heffen, heffend, heft, hief, hieft, hieve, hieven*}, {*geholpen, help, helpe, helpen, helpend, helpt, hielpt, hielpe, hielpen, hielpt*}, {*gekocht, kocht, kochte, kochten, koop, koopt, kope, kopen, kopend*}, {*gelegen, laagt, lag, lage, lagen, lig, ligge, liggen, liggend, ligt*}, {*genomen, naam, nam, name, namen, neem, neemt, neme, nemen, nemend*}, {*geschapen, schep, scheppen, scheppend, schept, schiep, schiepen, schiept*}, {*gescholden, scheld, schelde, schelden, scheldend, scheldt, schold, scholden, scholdt*}, {*gevraagd, vraag, vraagd, vraagde, vraagden, vraagt, vrage, vragen, vragend, vroeg, vroege, vroegen, vroegt*}, {*gewogen, weeg, weegt, wege, wegen, wegend, woge, wogen, woog, woogt*}, {*gezeten, zat, zate, zaten, zit, zitte, zitten, zittend*}, {*gezocht, zocht, zachte, zochten, zoek, zoeke, zoeken, zoekend, zoekt*}.

### 5.3.3 Heuristic detection of compounds

The following algorithm for heuristic detection of Dutch compounds differs from the Danish version (Algorithm 5.12) only in some specific details. To make the context clear, we still state the algorithm in full. (In what follows, the string minus operation  $\hat{\beta} \ominus \hat{\alpha}$  is prescribed by Definition 5.11.)

**Algorithm 5.34** (Heuristic identification of Dutch binary compounds). *Let  $\Lambda^{\hat{\beta}} = \{\hat{\rho}_1, \dots, \hat{\rho}_Q\}$  be a list of distinct Dutch essential roots (without vowel blotting) that contain at least one instance of  $\mathbf{V} = (a|A|e|\hat{e}|E|i|\hat{i}|o|\hat{o}|O|u|\hat{u}|\hat{u}|U|y)$  and DO NOT match the following string patterns:*

(achtX|al|an|At|bar|beX|con|el|en|Er|ferX|fol|fOrX|ge|ging|grOt|kOm|le|man|me|o|of|on|ond|onge|per|u|üt|wAr|xlos).

The output of the function  $\text{CpdDet}(\Lambda^{\hat{\beta}})$  is obtained through the following procedures:

- (1) Construct a list  $\{(\hat{\rho}_1, \lambda_1^{\hat{\beta}}), \dots, (\hat{\rho}_Q, \lambda_Q^{\hat{\beta}})\}$  where  $\lambda_q^{\hat{\beta}} = \{\hat{\rho}_{(q,1)}, \dots, \hat{\rho}_{(q,n_q)}\}$  is a subset of  $\Lambda^{\hat{\beta}}$  whose members all match the string pattern  $\hat{\rho}_q \sim$ , for  $q \in \mathbb{Z} \cap [1, Q]$ .
- (2) Expand the aforementioned entry  $(\hat{\rho}_q, \lambda_q^{\hat{\beta}})$  into a list of triplets  $\{(\hat{\rho}_{(q,1)}, \hat{\rho}_q, \hat{\rho}_{(q,1)} \ominus \hat{\rho}_q), \dots, (\hat{\rho}_{(q,n_q)}, \hat{\rho}_q, \hat{\rho}_{(q,n_q)} \ominus \hat{\rho}_q)\}$  for every  $q \in \mathbb{Z} \cap [1, Q]$  such that  $\lambda_q^{\hat{\beta}} \neq \emptyset$ . Collect all these triplets as one runs through the list  $\{(\hat{\rho}_1, \lambda_1^{\hat{\beta}}), \dots, (\hat{\rho}_Q, \lambda_Q^{\hat{\beta}})\}$ . The list of these triplets  $\{(\hat{\rho}_{(1)}, \hat{\eta}_{(1)}, \hat{\rho}_{(1)} \ominus \hat{\eta}_{(1)}), \dots, (\hat{\rho}_{(Q')}, \hat{\eta}_{(Q')}, \hat{\rho}_{(Q')} \ominus \hat{\eta}_{(Q')})\}$  contains potentially valid decompositions of compounds.
- (3) Screen the aforementioned list of triplets as follows: for every  $q' \in \mathbb{Z} \cap [1, Q']$ , if  $(\hat{\rho}_{(q')}, \hat{\eta}_{(q')}, \hat{\tau}_{(q')} = \hat{\rho}_{(q')} \ominus \hat{\eta}_{(q')})$  satisfies

$$\ell(\hat{\rho}_{(q')} \ominus \hat{\eta}_{(q')}) \geq 2 \quad \text{AND} \quad \hat{\rho}_{(q')} \ominus \hat{\eta}_{(q')} = \mathbf{X}_1(a|A|e|\hat{e}|E|i|\hat{i}|o|\hat{o}|O|u|\hat{u}|\hat{u}|U|y)\mathbf{X}_2,$$

then construct  $\hat{\tau}_{(q')}^*$  by performing  $(e|n|s) \sim \rightarrow \emptyset$  on  $\hat{\tau}_{(q')}$  and  $\hat{\tau}_{(q')}^{**}$  by doing  $(en|er|es|ge) \sim \rightarrow \emptyset$  on  $\hat{\tau}_{(q')}$ , before generating a list  $\hat{\lambda}_{(q')}^*$  by members of  $\Lambda^{\hat{\beta}}$  that match the pattern  $(\hat{\tau}_{(q')} | \hat{\tau}_{(q')}^* | \hat{\tau}_{(q')}^{**})$ ; otherwise, set  $\hat{\lambda}_{(q')}^* = \emptyset$ .

- (4) Collect all the triplets  $(\hat{\rho}_{(q')}, \hat{\eta}_{(q')}, \hat{\lambda}_{(q')}^*)$  where  $\hat{\lambda}_{(q')}^*$  is non-void and  $\hat{\tau}_{(q')}$  DOES NOT match the following string patterns:

$$(h(e|\hat{e})d\mathbf{X}|ing\mathbf{X}|lyk\mathbf{X}|tet\mathbf{X}|um).$$

This list of triplets  $\text{CpdDet}(\Lambda^{\hat{\beta}})$  contains the heuristic decompositions of all the identified binary compounds.

**Algorithm 5.35** (Approximate clustering of Dutch words with heuristic detection of compounds). *The procedure runs essentially the same way as Algorithm 5.35, except that Dutch rules replace Danish rules.*

*Example 5.35.1.* Testing the algorithm above against the combined inputs from Examples 5.33.1 and 5.33.2, we obtain the following result:

{bak, bakke, bakken, bakkend, bakt, bakte, bakten, gebakken}, {ban, bande, banden, banne, bannen, bannend, bant, gebannen}, {bederf, bederft, bederve, bederven, bedervend, bedierf, bedierft, bedorve, bedorven}, {bedrieg, bedriege, bedriegen, bedriegend, bedriegt, bedroge, bedrogen, bedroog, bedroogt}, {beet, bete, beten, bijt, bijte, bijten, bijtend, gebeten}, {begin, beginne, beginnen, beginnend, begint, begon, begonne, begonnen, begont}, {best, beste, beter, betere, beters, goed, goede, goeden, goeder, goeds}, {bied, biede, bieden, biedend, biedt, bode, boden, bood, boodi, geboden}, {bind, binde, binden, bindend, bindt, bond, bonde, bonden, bondt, gebonden}, {blaas, blaast, blaze, blazen, blazend, blies, bliest, blieze, bliezen, geblazen}, {bleek, bleekt, bleke, bleken, blijk, blyke, blyiken, blykend, blykt, gebleken}, {boge, bogen, boog, boogt, buig, buige, buigen, buigend, buigt, gebogen}, {braakt, brak, brake, braken, breek, breekt, breke, breken, brekend, gebroken}, {bracht, brachte, brachten, breng, brenge, brengen, brengend, brengt, gebracht}, {daad, daden}, {dacht, dachte, dachten, denk, denke, denken, denkend, denkt, gedacht}, {dooier, dooiers, dooiertje, eidooier, eidooiers, eiderdooier, eiderdooiers, eiderdooiertje}, {draag, draagt, drage, dragen, dragend, droeg, droegen, droegt, gedragen}, {droop, droopt, drope, dropen, druip, druipe, druipen, druipend, druipt, gedropen}, {ei, eidooier, eidooiers, ei-dooiertje, eiderdooier, eiderdooiers, eiderdooiertje, eieren, eitje, eitjes, eiwit, eiwitten, eiwittje}, {eiwit, eiwitten, eiwittje, wit, wits, witst, witste, witte, witter, wittere, witters}, {gegolden, geld, gelde, gelden, geldend, geldt, gold, golde, golden, goldt}, {geheven, hef, heffe, heffen, heffend, heft, hief, hieft, hieve, hieven}, {geholpen, help, helpe, helpen, helpend, helpt, hielpe, hielpen, hielpt}, {gekocht, kocht, kochte, kochten, koop, koopt, kope, kopen, kopen}, {geleerd, leer, leerde, leert, lere, leren, lerend}, {gelegen, laagt, lag, lage, lagen, lig, ligge, liggen, liggend, ligt}, {gelost, los, losse, lossen, lossend, lost, loste, losten}, {genomen, naam, nam, name, namen, neem, neemt, neme, nemen, nemend}, {geraasd, raas, raasde, raasden, raast, raze, razen, razend}, {geschapen, schepp, scheppen, scheppend, schept, schiep, schiepe, schiepen, schiept}, {gescholden, scheld,

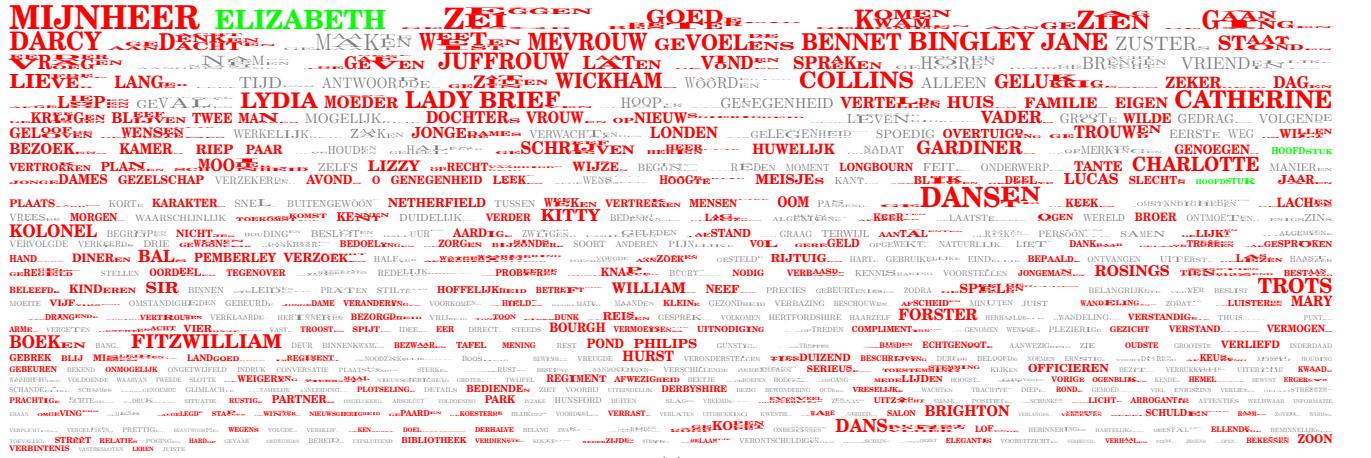
*schelde, schelden, scheldend, scheldt, schold, scholde, scholden, scholdt}, {gevraagd, vraag, vraagd, vraagde, vraagden, vraagt, vrage, vragen, vragend, vroeg, vroege, vroegen, vroegt}, {gewerkt, werk, werke, werken, werkend, werkt, werkte, werkten}, {gewogen, weeg, weegt, wege, wegen, wegend, woge, wogen, woog, woogt}, {gezeten, zat, zate, zaten, zit, zitte, zitten, zittend}, {gezocht, zocht, zochte, zochten, zoek, zoek, zoeken, zoekend, zoekt}, {klein, kleine, kleinen, kleiner, kleinere, kleiners, kleins, kleinst, kleinste}, {maan, maantje, man, manen, manne, mannen, mans}, {vrouw, vrouwe, vrouwen}.*

Here, the compounds *ei(Ø|er)dooier* “egg yolk” and *eiwit* “egg white” are correctly dissolved into their respective constituting components.

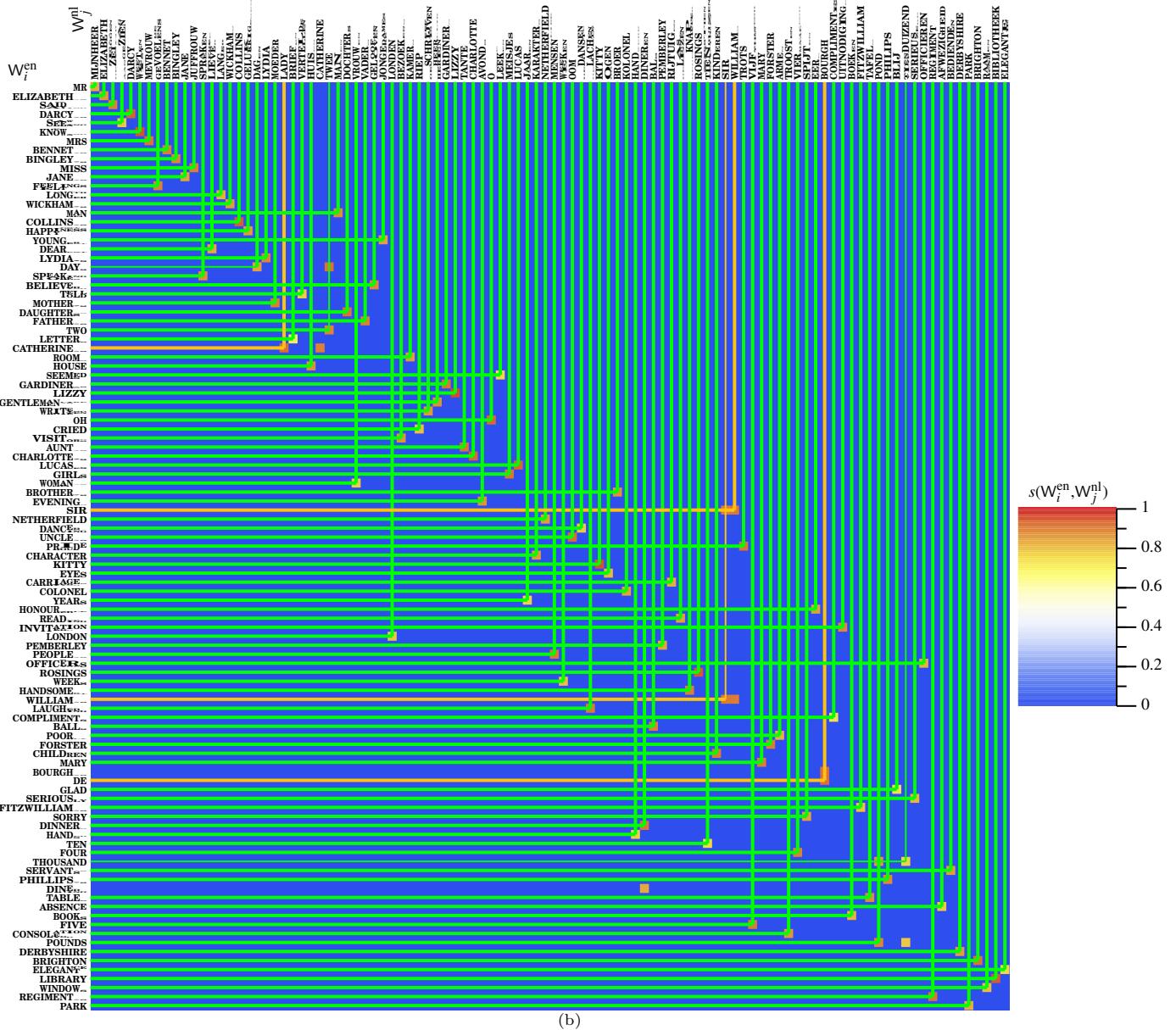
*Example 5.35.2.* In Fig. S6, we further apply the aforementioned word clustering algorithm to topic extraction and machine translation (see Table S1 for text source).

Note that *de* “the” is a Dutch stop word, so the English *de* cannot be translated exactly by our algorithm.

The co-occurring words in *Lady Catherine (de Bourgh)* and *Sir William (Lucas)* also make it difficult to resolve them exactly in our numerical experiments. These pairs of mismatches are indicated by amber cross-hairs in Fig. S6b.



(a)



(b)

Fig. S6. Text mining in Dutch. (a) Statistically identified topics ( $n_{ii} \geq 20$ ) in a Dutch version of *Pride and Prejudice*, with the same color encoding scheme as Fig. S3. (b) Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{nl}})$  between selected topics in English and Dutch versions of *Pride and Prejudice*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) rank of the similarity score for the optimal node. Green (resp. amber) cross-hair indicates an exact (resp. a close but non-exact) match.

## 6 Approximate word clustering in selected Romance languages

In this section, we present the approximate clustering algorithms for two representative modern Romance languages with a large number of speakers: Spanish and French. Their ancestral language, Latin, is also treated in this section.

Due to heavy loads of Latinate loanwords in English (borrowed either directly from Latin, or through Norman French), the word formation of English shares a lot of features with modern Romance languages. However, some critical differences must be noted:

- Unlike English, apostrophes and hyphens are used extensively in French to mark word boundaries.
- Unlike English, verbs in all these three languages conjugate in all persons for all tenses and aspects, thus making the verb morphology highly complicated.
- Unlike English, adjectives in all these three languages inflect in both gender and number.
- The classical language Latin additionally have declensions of nouns and adjectives in six cases.

The descriptions above (except the case declension of Latin) also extend to modern Italian, together with some Romance languages spoken in the Iberian Peninsula other than Spanish (namely, Catalan, Galician, Portuguese and Valencian), with minor modifications. Romanian, a Romance language spoken in Eastern Europe, has more conservative morphological structures (such as three genders and five cases), which are closer to the complexity of Latin (which has three genders and six cases) than all the other modern Romance languages.

### 6.1 Modified Porter stemming algorithm for Spanish

Before defining Spanish stop words, we need to define Spanish vowel extensions, and construct an algorithm to stress the penultimate vowel in a Spanish word, by heuristics.

**Definition 6.1** (Spanish Vowel Extensions). Hereafter in §6.1, the symbol  $\mathbf{V}^*$  stands for any member from the list  $\{a, e, i, \hat{i}, o, u\}$ , the so-called Spanish vowel extensions. In line with the multiplicity notations introduced in Definition 3.3, the symbol  $\mathbf{V}_m^*$  stands for a text string formed by consecutive appearance of one or more (not necessarily identical or distinct) members from the set of Spanish vowel extensions.

Dual to the notations above, the symbol  $\mathbf{C}^*$  stands for any character that does not belong to the list  $\{a, e, i, \hat{i}, o, u\}$ , and  $\mathbf{C}_{m_0}^*$  stands for a text string formed by consecutive appearance of zero or more (not necessarily identical or distinct) characters that do not belong to the same list.  $\square$

**Algorithm 6.2** (Penultimate Stress). *We define StrPenult( $\hat{\sigma}$ ) through the following operations on  $\hat{\sigma}$ :*

- (1) Remove diacritic marks from  $\hat{\sigma}$ , except that the Spanish letter  $\tilde{n}$  must be kept intact.
- (2) Do  $\sim(\mathbf{V}^*)^X \epsilon (\mathbf{C}_m^* \mathbf{V}_m^* \mathbf{C}_{m_0}^*) \rightarrow \mathbf{V}'^X$ , where one derives  $\mathbf{V}'$  from  $\mathbf{V}^*$  by adding stress marks to (a|e|i|o|u). The result so far is called  $\hat{\sigma}'$ .
- (3) If  $\Omega(\hat{\sigma}') = (d|r)$  or  $\hat{\sigma}' = \text{está}$ , then perform the operation in Step (1) on  $\hat{\sigma}'$  and output StrPenult( $\hat{\sigma}$ ). Otherwise, check whether  $\hat{\sigma}' = \mathbf{Xea}$ . If so, StrPenult( $\hat{\sigma}$ ) results from doing  $\sim ea \rightarrow \mathbf{é}a$  on  $\hat{\sigma}'$ ; if not, StrPenult( $\hat{\sigma}$ ) =  $\hat{\sigma}'$ .

**Definition 6.3** (Spanish stop words). If a word belongs to the following list<sup>70</sup>:

*a, abajo, acá, acaso, adentro, afuera, ahí, ahora, al, algo, alguien, algun, algúin, alguna, algunas, alguno, algúnos, allá, allí, alrededor, ambas, ambos, ante, antes, aquel, aquél, aquella, aquélla, aquellas, aquéllas, aquello, aquellos, aquéllos, aquí, arriba, así, atrás, aun, aúin, aunque, bajo, bien, cada, casi, cerca, como, cómo, con, connigo, connosco, consigo, contigo, contra, convusco, cual, cuál, cuales, cualesquier, cualesquiera, cualquier, cualquiera, cuan, cuán, cuando, cuándo, cuanta, cuánta, cuantas, cuántas, cuanto, cuánto, cuantos, cuántos, cu-ya, cuya, cuyas, cuyás, cuyo, cuyos, de, debajo, del, delante, dentro, desde, después, detrás, donde, dónde, durante, durea, e, el, él, ella, ellis, ello, ellos, empero, en, encima, enfrente, entonces, entre, esa, ésa, esas, ésas, ese, ése, éses, eso, esos, esta, ésta, estas, éstas, este, éste, éstes, esto, estos, frente, hacia, hasta, jamás, junto, la, las, le, lejos, les, lo, los, luego, mas, más, me, mediante, menos, mi, mí, mía, mías, mío, mios, mis, misma, mismas, mismísima, mismísimas, mismísimo, mismísimos, mismo, mismos, mucha, muchas, mucho, muchos, muy, nada, nadie, ni, ningún, ninguna, ningunas, ninguno, ningunos, no, nos, nosotras, nosotros, nuestra, nuestras, nuestro, nuestros, nunca, o, os, otra, otras, otro, otros, para, pero, poco, por, porque, pronto, pues, que, qué, quien, quién, quienes, quizá, segün, si, sí, siempre, sin, sino, siquiera, so, sobre, su, sus, suya, suyas, suyo, suyos,*

<sup>70</sup>Our list of Spanish stop words is based on [snowball.tartarus.org/algorithms/spanish/stop.txt](http://snowball.tartarus.org/algorithms/spanish/stop.txt), with extensive additions.

también, tampoco, tan, tanta, tantas, tanto, tantos, te, ti, toda, todas, todavía, todo, todos, tras, tu, tú, tus, tuyas, tuyos, u, un, una, uno, unos, usted, ustedes, vosotras, vosotros, vuestra, vuestras, vuestra, vuestras, vuestro, vuestros, y, ya, yo,

then we consider it a Spanish stop non-verb (notation: **StopNonVerbSpanish**). If a word belongs to the following list:

era, erais, éramos, eran, eras, eres, es, está, estaba, estabais, estaban, estabas, estad, estada, estadas, estado, estados, estáis, estamos, están, estando, estar, estará, estarán, estarás, estaré, estaréis, estaremos, estaría, estaríais, estaríamos, estarían, estarías, estás, esté, estéis, estemos, estén, estés, estoy, estuve, estuviera, estuvierais, estuvíeramos, estuvieran, estuvieras, estuviere, estuviereis, estuvíremos, estuvieren, estuvieres, estuvieron, estuviese, estuvieseis, estuvísemos, estuviesen, estuvieses, estuvimos, estuviste, estuvisteis, estuvo, fue, fuera, fuerais, fuéramos, fueran, fueras, fuere, fuereis, fuéremos, fueron, fueres, fueron, fuese, fueseis, fuésemos, fuesen, fuses, fui, fuimos, fuiste, fuisteis, ha, habe, habed, habéis, haber, había, habíais, habíamos, habían, habias, habida, habidas, habido, habidos, habiendo, habrá, habrán, habrás, habré, habréis, habremos, habría, habriais, habriamos, habrian, habrias, hace, hacé, haced, hacéis, hacemos, hacen, hacer, haces, hacestú, hacésvos, hacía, hacíais, hacíamos, hacían, hacías, haciendo, haga, hagáis, hagamos, hagan, hagas, hago, han, hará, harán, harás, haré, haréis, haremos, haría, haríais, haríamos, harían, harías, has, hay, haya, hayáis, hayamos, hayan, hayas, haz, he, hé, hecha, hechas, hecho, hechos, hemos, hice, hiciera, hiciera, hicíeras, hicíremos, hicieran, hicieras, hiciere, hiciereis, hicíremos, hicieren, hicieres, hicieron, hiciese, hicieseis, hicísemos, hiciesen, hicieses, hicimos, hiciste, hicisteis, hizo, hube, hubiera, hubiera, hubiéramos, hubieran, hubieras, hubiere, hubiereis, hubiéremos, hubieren, hubieres, hubieron, hubiese, hubieseis, hubísemos, hubiesen, hubieses, hubimos, hubiste, hubisteis, hubo, podáis, podamos, poded, podéis, podemos, poder, podia, podíais, podiamos, podían, podías, podida, podidas, podido, podidos, podrá, podrán, podrás, podré, podréis, podremos, podría, podríais, podríamos, podrían, podrías, pude, pudiendo, pudiera, pudiera, pudíramos, pudieran, pudieras, pudiere, pudiereis, pudíremos, pudieren, pudieres, pudieron, pudiese, pudieseis, pudísemos, pudiesen, pudieses, pudimos, pudiste, pudisteis, pudo, pueda, puedan, puedes, pueden, puedo, se, sé, sea, seáis, seamos, sean, seas, sed, ser, será, serán, serás, seré, seréis, seremos, sería, seríais, seríamos, serían, serías, sida, sidas, sido, sidos, siendo, sois, somos, son, sos, soy, ten, tendrá, tendrán, tendrás, tendré, tendréis, tendremos, tendría, tendríais, tendríamos, tendrían, tendrías, tené, tened, tenéis, tenemos, tener, tenés, tenga, tengáis, tengamos, tengan, tengas, tengo, tenía, teníais, teníamos, tenían, tenías, tenida, tenidas, tenido, tenidos, teniendo, tiene, tienen, tienes, tuve, tuviera, tuviera, tuvíramos, tuvieran, tuviera, tuviera, tuvíremos, tuvieren, tuviera, tuvieron, tuviese, tuvieseis, tuvísemos, tuviesen, tuvieses, tuvimos, tuviste, tuvisteis, tuvo,

then we consider it a Spanish stop verb (notation: **StopVerbSpanish**). A Spanish stop word matches the following string pattern:

(**StopNonVerbSpanish|StopVerbSpanish|StrPenult(StopVerbSpanish)SpanishVerbCombo**),

where the string pattern **SpanishVerbCombo** = (*la|las|le|les|lo|los|me|nos|os|se|te*) points to pronouns that may attach to the end of Spanish verbs to form “combined forms”. All the Spanish stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list.  $\square$

The following Spanish verbs have highly irregular conjugations:

*da, daba, dabais, dábamos, daban, dabas, dada, dadas, dado, dados, dais, damos, dan, dando, dar, dará, darán, darás, daré, daréis, daremos, daría, daríais, daríamos, darían, darías, das, dé, deis, demos, den, des, di, diera, dierais, diéramos, dieran, dieras, diere, diereis, diéremos, dieren, dieres, dieron, diese, dieseis, diésemos, diesen, dieses, dimos, dio, diste, disteis, doy — “give”;*

*fue, fuera, fuerais, fuéramos, fueran, fueras, fuere, fuereis, fuéremos, fueron, fueres, fuese, fueseis, fuésemos, fuesen, fuses, fui, fuimos, fuiste, fuisteis, iba, ibais, ibamos, iban, ibas, id, ida, idas, ido, idos, ir, irá, irán, irás, iré, iréis, iremos, iría, iríais, iríamos, irían, irías, va, vais, vamos, van, vas, vaya, vayáis, vayamos, vayan, vayas, ve, voy, yendo — “go”.*

We are not going to treat them as stop words, so as to be consistent with their English counterparts. Nevertheless, we will identify these string patterns as **giveSpanish** and **goSpanish**, to facilitate the clustering of content words.

### 6.1.1 Effective spelling and essential root

**Algorithm 6.4** (Spanish effective spelling). *For a Spanish word  $\hat{\sigma}$ , its effective spelling  $\text{EffSpell}(\hat{\sigma})$  is constructed in six sequential steps:*

## (1) Replace

$af \sim$	$ag \sim$	$amab(\emptyset i)l \sim$	$ami(g stad stos) \sim$	$año \sim$	$at \sim$	$biblio \sim$	$caus \sim$	$charlott$	$cinc \sim$	$com(u u)n \sim$
$\alpha f$	$\alpha g$	$\phi \rho e \tilde{\nu} \delta$	$\phi \rho e \tilde{\nu} \delta$	$\eta a \rho o$	$\alpha t$	$\beta i \beta l i \omega$	$\kappa a u \sigma$	$\check{s} a r l o t \tau$	$q i 5$	$\kappa \omega m u \tilde{v}$
$co \mu e n t$	$co \mu p o \mu$	$\kappa \tau \rho$	$\delta$	$\phi e \lambda i \zeta$	$\gamma a \rho \delta$	$\iota$	$\lambda i \sigma \tau$	$\lambda o \tilde{v} \gamma$	$\pi u \rho$	$\alpha a r \mu$
$p a r e (\emptyset z)c$	$p o s i b \sim$	$p o s i c \sim$	$p o s i t i v \sim$	$p r i m a v e r$	$p r i m i t$	$p u e r t a$	$v e n e r \sim$	$v i n o \sim$	$c a r t a (\emptyset s)$	$j u e v e s$
$\pi \varepsilon a \rho$	$\pi o \sigma i \beta$	$\pi \omega \sigma i \zeta$	$p w s i t i \beta$	$\sigma \pi \rho i \tilde{\nu} \gamma$	$\pi \rho i \mu i t$	$p y o r t r$	$\beta e \tilde{v} e \rho$	$\omega i \tilde{v} \epsilon$	$\lambda e t \tau \rho$	$j 4 u e 4$
										$m a r \eta$

## (2) Replace

$(\acute{a} a) n i m \sim$	$(d i   d i \acute{g} a) \text{SpanishVerbCombo} \sim$	$(d i a   d i \acute{a}) \sim$	$(i e   i \acute{e})$	$(n d r   n g)$	$(z c   z g)$	$\acute{a}$	$a b u e l o \sim$	$a c \sim$
$\alpha n i \mu$	$d e c i r$	$\delta i a$	$\hat{i}$	$n$	$z$	$a$	$a b u e l \omega$	$a c$
$dece \sim$	$d i f \sim$	$d i j \sim$	$d i r e m o s$	$d i r e \sim$	$e d a d \sim$	$e d u \sim$	$e n \sim$	$e s \sim$
$\delta e c e$	$\delta i f$	$d i c h$	$d e c i r$	$\delta i r e$	$\delta e a \delta$	$\delta e \bar{u}$	$\varepsilon \tilde{v}$	$\varepsilon s$
$ldr$	$marid \sim$		$mejor(\emptyset es) \sim$	$o c \sim$	$\acute{o} n$	$p a r (e ie) n t \sim$	$qu$	$s g$
$l$	$muarid$		$bueno$	$o c$	$o n e$	$\pi a p i e n t$	$k$	$s$
$t i o \sim$	$y e \sim$		<u>StrPenult(giveSpanish)SpanishVerbCombo</u>			<u>StrPenult(goSpanish)SpanishVerbCombo</u>		
$t i \omega$	$e$		$\sigma \gamma i b \sigma$			$\sigma g o \sigma$		
	<u><math>(m i s s   s e \tilde{n} o r i t a   s r t a ) (\emptyset s)</math></u>		<u><math>(m r   s e \tilde{n} o r) (\emptyset es)</math></u>		<u><math>(m r s   s e \tilde{n} o r a   s r a ) (\emptyset s)</math></u>	<u>giveSpanish</u>	<u>goSpanish</u>	
	$m d s e \tilde{n} o r i t a$		$m r s e \tilde{n} o r$		$f s e \tilde{n} o r a$	$\sigma \gamma i b \sigma$	$\sigma g o \sigma$	
	<u><math>m a m (a á i)(\emptyset s)</math></u>		<u><math>p a p (a á a c i t o   a i t o   i)(\emptyset s)</math></u>		$\sim \hat{x}(i i)s(m o   t a   t i c a   t i c o ) (\emptyset s)$	$\sim d o s$	$\sim t r i(c e s z)$	
	$m a d r e$		$p a d r e$		$\hat{x} e$	$das$	$tor$	

## (3) Replace

$(\emptyset h)(u ü)e$	$c(c t)$	$co^{\mathbf{X}\epsilon}(g j)$	$\acute{e} n$	$\acute{i}$
$o$	$z$	$r a c o \mathbf{X}$	$ena$	$i$

(4) Do  $\acute{e} \rightarrow e$ ,  $\tilde{n} \rightarrow n$ ,  $\acute{o} \rightarrow o$ ,  $(\acute{u}|ü) \rightarrow u$ .(5) Do  $c^{\mathbf{V}\epsilon}(a|o|u) \rightarrow k\mathbf{V}$ ,  $c^{\mathbf{V}\epsilon}(e|i|\tilde{i}) \rightarrow z\mathbf{V}$ ,  $g^{\mathbf{V}\epsilon}(e|i|\tilde{i}) \rightarrow j\mathbf{V}$ ,  $gu \rightarrow g$ .(6) Do  $\underline{ven} \rightarrow vis$ ,  $\sim(d|mo|mos|n|r)\text{SpanishVerbCombo} \rightarrow \emptyset$ ,  $\sim idad(\emptyset|es) \rightarrow \emptyset$ .**Definition 6.5** (Spanish protected range). Let  $\hat{\sigma}$  be a text string derived from a Spanish word, its protected range  $\text{ProtRg}(\hat{\sigma})$  is an integer determined as follows:

- Try to find the string pattern  $(\emptyset|al|des|es|ko|pre|pro|ra|re|tran|tras)\mathbf{C}_{m_0}^*\mathbf{V}^*(\emptyset|(b|d|g|m|n|r|s|t|v|z)_m) \sim$  in the string  $\hat{\sigma}$ ;
- If the string pattern above is found, the last position occupied by such a string defines  $\text{ProtRg}(\hat{\sigma})$ ; otherwise, set  $\text{ProtRg}(\hat{\sigma}) = 0$ .  $\square$

**Algorithm 6.6** (Spanish essential root). Let  $\hat{\sigma}$  be the effective spelling of a Spanish word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:(1) Do  $kom\mathbf{V}^* \sim \rightarrow kkom\mathbf{V}^*$  on  $\hat{\sigma}$  and call the result  $\hat{\sigma}_*$ .(2) Break down  $\hat{\sigma}_* = \hat{\sigma}_1 \hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}_*^{[\text{ProtRg}(\hat{\sigma}_*)]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma}_*)$  is equal to the protected range of  $\hat{\sigma}_*$ .(3) On  $\hat{\sigma}_2$ , perform the following substitutions in a sequel:(3.1) Do  $\sim(d|t)(a|o)(\emptyset|s)(\emptyset|\text{SpanishVerbCombo}) \rightarrow \emptyset$ ,  $\sim d(\emptyset|es) \rightarrow \emptyset$ .(3.2) Do  $(ba|(b(\emptyset|i)(\emptyset|l))|nt|(r\mathbf{V}^*|\sim r)|st|\sim n)(\emptyset|\text{SpanishVerbCombo}) \rightarrow \emptyset$ .(3.3) Do  $(d o r | g | i o n | k a | k o | l o g | m | s | t i v | u v | y | z) \rightarrow \emptyset$ .<sup>71</sup>(3.4) Do  $\sim\mathbf{V}_m^* \rightarrow \emptyset$ .<sup>71</sup>We note that longer matches take priority over shorter matches: if *log* is found, then delete these three letters altogether, instead of just a single letter *g*.

The result after these four steps of operations is called  $\hat{\sigma}'_2$ .

- (4) Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ .
- (5) Perform the following replacements on  $\hat{\sigma}_1\hat{\sigma}'_2$ :

$\mathbf{X}^{\infty}(\emptyset kon re)di$	$durm$	$o$	$\sim d(ez ich ig ij iz)$	$\sim duj$	$\sim jog$	$\sim k(e u)p$	$\sim kerr$	$\sim kis$	$\sim p(ost us)$
$\mathbf{X}dir$	$dorm$	$oxx$	$dir$	$duz$	$jug$	$kab$	$ker$	$kir$	$pon$
$\sim pud$	$\sim traj$		$\sim tuv$		$\sim v(e ed em er es im ind ir is ist)$			$\sim yaz$	
$pod$	$tra$		$ten$			$vis$			$yag$

### 6.1.2 Admissible mutation and approximate clustering

**Algorithm 6.7** (Spanish vowel blotting). Set  $\mathbf{V}_m^*$  and  $\mathbf{C}_{m_0}^*$  as in Definition 6.1. For a token string  $\hat{\sigma}$ , its blotted form  $\text{BlotV}_1(\hat{\sigma})$  is constructed as follows:

- If the string pattern  $(\emptyset|a|des|e|es|i|i|ko|o|pre|pro|ra|re|tran|tras|u)\mathbf{C}_{m_0}^*\mathbf{V}_m^*\sim$  can be found in the string  $\hat{\sigma}$ , then the last position occupied by such a pattern is replaced by the letter “a”.
- Otherwise, leave the string  $\hat{\sigma}$  intact.

Similar to what we did in §5 for three representative Germanic languages, we will construct a bivariate Boolean-valued function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  on a “simple heredity test function” in Algorithm 6.8, and a set of “admissible suffix mismatch and vowel alternation” rules in Algorithm 6.9.

**Algorithm 6.8** (Simple heredity test). The Boolean-valued function  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if the lowercase form of  $\hat{\alpha}$  contains at least one instance of  $\mathbf{V}^*$  (Definition 6.1) **AND** at least one of the following four conditions holds:<sup>72</sup>

- (i) After doing  $\sim s \rightarrow \emptyset$  on  $\hat{\beta}$ , we obtain  $\hat{\alpha}$ ;
- (ii) After doing  $\sim a(\emptyset|s) \rightarrow o$  on  $\hat{\alpha}$ , we obtain  $\hat{\beta}$ ;
- (iii)  $\hat{\beta} = \hat{\alpha}k$ ;
- (iv)  $\ell(\hat{\beta}) > \ell(\hat{\alpha}) \geq \frac{\ell(\hat{\beta})}{2}$  **AND**  $\hat{\alpha} = \hat{\beta}^{[\ell(\hat{\alpha})]}$  **AND**  $\hat{\beta}^{[\ell(\hat{\alpha}')+1]} = \mathbf{V}^*$ . (See Definition 3.1 for the notations  $\hat{\beta}^{[n]}$  and  $\hat{\beta}^{\{n\}}$ .)

In what follows, we define  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{NW}^*(\hat{\alpha}, \hat{\beta})$  and  $\text{SuffixSW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{RootSW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{SW}^*(\hat{\alpha}, \hat{\beta})$  as what is done in the Danish case (Algorithm 5.7), which was also applied to two other Germanic languages treated in §5.

**Algorithm 6.9** (Admissible suffix mismatch and vowel alternation). For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function

$$\text{AdmMut}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta}))$$

returns **TRUE** if

$$\text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = [(\emptyset|\mathbf{V}^*\mathbf{X}(d|t)|(r|s)\mathbf{X}|tiv\mathbf{X}), (\emptyset|\mathbf{V}^*\mathbf{X}(d|t)|(r|s)\mathbf{X}|tiv\mathbf{X})]$$

**AND** at least one of the following two conditions holds:

- (i)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \emptyset$  **AND**  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$  contains at least one instance of  $\mathbf{V}^*$ ;
- (ii)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = ([e, i]| [e, i]| [i, i]| [\hat{i}, i])$ .

Similarly, one can evaluate another Boolean-valued function

$$\text{AdmSM}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta}))$$

by trading all the occurrences of NW in the statements above with SW.

**Algorithm 6.10** (Heredity test function). The structure of the Spanish heredity test function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  is identical to the German version (Algorithm 8.1.2), except that the functions  $\text{SimpHrdTest}$ ,  $\text{RootNW}$ ,  $\text{SuffixNW}$ ,  $\text{NW}^*$ ,  $\text{RootSW}$ ,  $\text{SuffixSW}$ ,  $\text{SW}^*$  must follow the Spanish rules stated above.

<sup>72</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical OR.

**Algorithm 6.11** (Approximate clustering of Spanish words). *The algorithm is essentially the same as Algorithm 5.10, except that Spanish rules (instead of Danish rules) apply to all the tags (effective spelling, essential root, vowel blotting etc.).*

*Example 6.11.1.* We first select some “simple” Spanish word families with (more or less) regular inflections. For the selected verbs, we have not only included conjugated forms (in different moods, persons and tenses) but also derived forms (such as adjectives and nouns associated with a verb). An “approximate translation” in English, enclosed in quotation marks, is appended to the end of each family. The readers should be reminded that such an “approximate translation” only matches certain members with a word family.

*concluí, concluía, concluíais, concluíamos, concluían, concluías, concluid, concluida, concluidas, concluido, concluidos, concluimos, concluir, concluirá, concluirán, concluirás, concluiré, concluiréis, concluirémos, concluiría, concluiríais, concluiríamos, concluirían, concluirías, concluís, concluiste, concluisteis, conclusión, conclusivo, concluso, concluya, concluyáis, concluyamos, concluyan, concluyas, concluye, concluyen, concluyendo, concluyente, concluyera, concluyerais, concluyéramos, concluyeran, concluyeras, concluyere, concluyereis, concluyéremos, concluyeren, concluyeres, concluyeron, concluyes, concluyese, concluyeseis, concluyésemos, concluyesen, concluyeses, concluyo, concluyó — “conclude”;*

*continua, continua, continuaba, continuabais, continuaban, continuabas, continuación, continuad, continuada, continuadamente, continuadas, continuado, continuador, continuadora, continuados, continuás, continuamente, continuamiento, continuamos, continuán, continuando, continuar, continuara, continuará, continuara, continuais, continuáramos, continuaran, continuárán, continuaras, continuarás, continuare, continuará, continuareis, continuáreis, continuaremos, continuáremos, continuaren, continuares, continuaria, continuáis, continuáramos, continuarian, continuarias, continuaron, continuás, continuase, continuaseis, continuásemos, continuasen, continuases, continuaste, continuasteis, continuativa, continuativo, continué, continué, continué, continuéis, continuemos, continuén, continués, continuidad, continuo, continuó, continuó — “continue”;*

*rápida, rápidamente, rápidas, rápido, rápidos — “fast”;*

*rapideces, rapidez — “speed”.*

Applying Algorithm 6.11 to the list of words above, we arrive at the following result:

{*concluí, concluía, concluíais, concluíamos, concluían, concluías, concluid, concluida, concluidas, concluido, concluidos, concluimos, concluir, concluirá, concluirán, concluirás, concluiré, concluiréis, concluirémos, concluiría, concluiríais, concluiríamos, concluirían, concluirías, concluís, concluiste, concluisteis, conclusión, conclusivo, concluso, concluya, concluyáis, concluyamos, concluyan, concluyas, concluye, concluyen, concluyendo, concluyente, concluyera, concluyerais, concluyéramos, concluyeran, concluyeras, concluyere, concluyereis, concluyéremos, concluyeren, concluyeres, concluyeron, concluyes, concluyese, concluyeseis, concluyésemos, concluyesen, concluyeses, concluyo, concluyó*},

{*continua, continua, continuaba, continuabais, continuaban, continuabas, continuación, continuad, continuada, continuadamente, continuadas, continuado, continuador, continuadora, continuados, continuás, continuamente, continuamiento, continuamos, continuán, continuando, continuar, continuara, continuará, continuara, continuais, continuáramos, continuaran, continuárán, continuaras, continuarás, continuare, continuará, continuareis, continuáreis, continuaremos, continuáremos, continuaren, continuares, continuaria, continuáis, continuáramos, continuarian, continuarias, continuaron, continuás, continuase, continuaseis, continuásemos, continuasen, continuases, continuaste, continuasteis, continuativa, continuativo, continué, continué, continué, continuéis, continuemos, continuén, continués, continuidad, continuo, continuó, continuó*},

{*rápida, rápidamente, rápidas, rapideces, rapidez, rápido, rápidos*}.

*Example 6.11.2.* The representative verbs from the three regular conjugations of Spanish are given below.

*coma, comáis, comamos, coman, comas, comás, come, comé, comed, coméis, comemos, comen, comer, comerá, comerán, comerás, comeré, comeréis, comeremos, comería, comeríais, comeríamos, comerían, comerías, comes, comés, comí, comía, comíais, comíamos, comían, comías, comida, comidas, comido, comidos, comiendo, comiera, comierais, comiéramos, comieran, comieras, comiere, comiereis, comiéremos, comieren, comieres, comieron, comiese, comieseis, comiésemos, comiesen, comieses, comimos, comió, comiste, comisteis, como — “eat”;*

*habla, hablá, hablaba, hablabais, hablábamos, hablaban, hablabas, hablad, hablada, habladas, hablado, hablad, habláis, hablamos, hablan, hablando, hablar, hablara, hablará, hablarais, habláramos, hablaran, hablarán,*

*hablaras, hablarás, hablare, hablaré, hablaréis, hablaréis, hablaremos, habláremos, hablaren, hablares, habla-  
ría, habaríais, habaríamos, habarian, habarias, habaron, hablas, hablás, hablaseis, hablásemos,  
hablasen, hablaces, hablaste, hablasteis, hable, hablé, habléis, hablemos, hablen, hables, hablés, hablo, habló  
—“speak”;*

viva, viváis, vivamos, vivan, vivas, vivás, vive, viven, vives, viví, vivía, vivíais, vivíamos, vivían, vivías, vivid, vivida, vividas, vivido, vividos, viviendo, viviera, viviera, vivierais, viviéramos, vivieran, vivieras, viviere, viviereis, viviéremos, vivieren, vivieres, vivieron, viviese, vivieseis, viviésemos, viviesen, vivieses, vivimos, vivió, vivir, vivirá, vivirán, vivirás, viviré, viviréis, viviremos, viviría, viviríais, viviríamos, vivirían, vivirías, vivís, viviste, vivisteis, vivo — “live”.

Without exceptions, the conjugated forms of these regular verbs follow the pattern **root + ending**, where the root remains invariant, while the ending starts with a vowel, and consists of at most two syllables (discounting final ~*mos*). We shall refer to this pattern as the “canonical root-ending dichotomy” hereafter.

Applying the clustering algorithm to the conjugated forms of the three Spanish verbs above, we obtain

{*coma, comáis, comamos, coman, comas, comás, come, comé, comed, coméis, comemos, comen, comer, comerá, comerán, comerás, comeré, comeréis, comeremos, comería, comeríais, comeríamos, comerían, comerías, comes, comés, comí, comía, comíais, comíamos, comían, comías, comida, comidas, comido, comidos, comiendo, comiera, comieraís, comiéramos, comieran, comieras, comiere, comiereis, comiéremos, comieren, comieres, comieron, comiese, comieseis, comiésemos, comiesen, comieses, comimos, comió, comiste, comisteis, como*}

{habla, hablá, hablaba, hablabais, hablábamos, hablaban, hablabas, hablad, hablada, habladas, hablado, hablad-  
os, habláis, hablamos, hablan, hablando, hablar, hablara, hablará, hablaraís, habláramos, hablaran, hablarán,  
hablaras, hablarás, hablare, hablaré, hablareis, hablaréis, hablaremos, habláremos, hablaren, hablares, habla-  
ria, hablariais, hablariamos, hablarian, hablarias, hablaron, hablas, hablás, hablse, hablseis, hablásemos,  
hablasen, hablaces, hablaste, hablasteis, hable, hablé, habléis, hablemos, hablen, hables, hablés, hablo, habló},

{viva, viváis, vivamos, vivan, vivas, vivás, vive, viven, vives, viví, vivía, vivíais, vivíamos, vivían, vivías, vivid, vivida, vividas, vivido, vividos, viviendo, viviera, vivierais, viviéramos, vivieran, vivieras, viviere, viviereis, viviéremos, vivieren, vivieres, vivieron, viviese, vivieseis, viviésemos, viviesen, vivieses, vivimos, vivió, vivir, vivirá, vivirán, vivirás, viviré, viviréis, viviremos, viviría, viviríais, viviríamos, vivirían, vivirías, vivís, viviste, vivisteis, vivo}.

*Example 6.11.3.* Our typology for irregular Spanish verbs is largely based on Wiktionary ([https://en.wiktionary.org/wiki/Wiktionary:Spanish\\_verb\\_inflection-table\\_templates](https://en.wiktionary.org/wiki/Wiktionary:Spanish_verb_inflection-table_templates)). There are about one hundred different types of irregularities on this official template list. We are not going to exhaust all of them in our sample words.

In our clustering test in this example, we ignore some “trivially irregular verbs”, whose roots are at most off by a stress mark (which will be dropped in the effective spelling, anyway) in conjugated forms. We also ignore irregular types whose conjugation template follows the “canonical root-ending dichotomy”, even though the detailed endings may deviate from the three aforementioned regular types (for example, the conjugations of *distinguir* and *planir*). Some highly irregular verbs in common use are also excluded from our test.

The typology on the aforementioned Wiktionary page regards [*o*, *ue*]Xar, [*o*, *ue*]Xer ending and [*o*, *ue*]Xir ending as three separate cases. We pick only one of these three scenarios for our test below, because our algorithm treats vowel alternations and endings independently. The same principle applies to other types of vowel alternations and consonant changes. What remain in our sample pool are the irregular verbs exhibiting systematic vowel alternation and/or consonant changes, which account for the overwhelming majority of irregular factors in Spanish verbs.

Concretely speaking, we send the following verbs

*adquiera, adquieran, adquieras, adquiere, adquieren, adquieres, adquiero, adquiráis, adquiramos, adquirí, adquiría, adquiríais, adquiríamos, adquirían, adquirías, adquirid, adquirida, adquiridas, adquirido, adquiridos, adquiriendo, adquiriera, adquirierais, adquiriéramos, adquirieran, adquirieras, adquiriere, adquiriereis, adquiriéremos, adquirieren, adquirieres, adquirieron, adquiriese, adquirieseis, adquirísemos, adquiriesen, adquirieses, adquirimos, adquirió, adquirir, adquirirá, adquirirán, adquirirás, adquiriré, adquiriréis, adquiriremos, adquiriría, adquiriríais, adquiriríamos, adquirirían, adquirirías, adquirís, adquiriste, adquiristeis — “acquire”,*

*anda, andá, andaba, andabais, andábamos, andaban, andabas, andad, andada, andadas, andado, andados, andáis, andamos, andan, andando, andar, andará, andarán, andarás, andaré, andaréis, andaremos, andaría, andariais, andariámos, andarian, andarias, andas, andás, ande, andéis, andemos, anden, andes, ando, anduve, anduviera, anduvierais, anduvíeramos, anduvieran, anduvieras, anduviere, anduviereis, anduvíeremos, anduvieren, anduvieres, anduvieron, anduviese, anduvieseis, anduvísemos, anduviesen, anduvieses, anduvimos, anduviste, anduvisteis, anduvo — “walk”;*

*argüíi, argüía, argüíais, argüíamos, argüían, argüías, argüid, argüida, argüidas, argüido, argüidos, argüímos, argüir, argüirá, argüirán, argüirás, argüiré, argüiréis, argüiremos, argüiría, argüiríais, argüiríamos, argüiríán, argüirías, argüís, argüiste, argüisteis, arguya, arguyás, arguyamos, arguyan, arguyas, arguye, arguyen, arguendo, arguyera, arguyeras, arguyéramos, arguyeran, arguyeras, arguyere, arguyereis, arguyéremos, arguyeren, arguyeres, arguyeron, arguyes, arguyeseis, arguyésemos, arguyesen, arguyeses, arguyo, arguyó — “infer”;*

*ase, asen, ases, asga, asgáis, asgamos, asgan, asgas, asgo, así, asia, asíais, asíamos, asian, asías, asid, asida, asidas, asido, asidos, asiendo, asiera, asierais, asiéramos, asieran, asieras, asiere, asiereis, asiéremos, asieren, asieres, asieron, asiese, asieseis, asiésemos, asiesen, asieses, asimos, asió, asir, asirá, asirán, asirás, asiré, asi-réis, asiremos, asiría, asiríais, asiríamos, asiríán, asirías, asís, asiste, asisteis — “grab”;*

*avergoncé, avergoncéis, avergoncemos, avergonzaba, avergonzabais, avergonzábamos, avergonzaban, avergon-zabas, avergonzad, avergonzada, avergonzadas, avergonzado, avergonzados, avergonzáis, avergonzamos, aver-gonzando, avergonzar, avergonzara, avergonzará, avergonzarais, avergonzáramos, avergonzaran, avergonza-rán, avergonzaras, avergonzarás, avergonzare, avergonzaré, avergonzareis, avergonzaréis, avergonzaremos, avergonzáremos, avergonzaren, avergonzares, avergonzaría, avergonzariais, avergonzariámos, avergonzarian, avergonzarias, avergonzaron, avergonzase, avergonzaseis, avergonzásemos, avergonzasen, avergonzases, aver-gonzaste, avergonzasteis, avergonzó, avergüence, avergüencen, avergüences, avergüenza, avergüenzan, aver-güenzas, avergüenzo — “embarrass”;*

*cabe, cabed, cabéis, cabemos, caben, caber, cabes, cabía, cabíais, cabíamos, cabian, cabias, cabidas, ca-bido, cabidos, cabiendo, cabrá, cabrán, cabré, cabréis, cabremos, cabría, cabríais, cabríamos, cabrían, cabrías, cupe, cupiera, cupierais, cupiéramos, cupieran, cupieras, cupiere, cupiereis, cupiéremos, cupieren, cu-pieres, cupieron, cupiese, cupieseis, cupiésemos, cupiesen, cupieses, cupimos, cupiste, cupisteis, cupo, quepa, quepáis, quepamos, quepan, quepas, quepo — “fit”;*

*cae, caed, caéis, caemos, caen, caer, caerá, caerás, caeré, caeréis, caeremos, caeria, caeríais, caeria-mos, caerían, caerías, caes, caí, caía, caíais, caíamos, caían, caías, caída, caídas, caído, caídos, caiga, caigáis, caigamos, caigan, caigas, caigo, caímos, caíste, caísteis, cayendo, cayera, cayerais, cayéramos, cayeran, caye-ras, cayere, cayereis, cayéremos, cayeren, cayeres, cayeron, cayese, cayeseis, cayésemos, cayesen, cayeses, cayó — “fall”;*

*coge, coged, cogéis, cogemos, cogen, coger, cogerá, cogerán, cogerás, cogeréis, cogeremos, cogería, co-geríais, cogeríamos, cogerían, cogerías, cogen, cogí, cogía, cogíais, cogíamos, cogían, cogías, cogida, cogidas, cogido, cogidos, cogiendo, cogiera, cogieraís, cogíramos, cogieran, cogieras, cogiere, cogiereis, cogiéremos, cogieren, cogieres, cogieron, cogiese, cogieseis, cogiésemos, cogiesen, cogieses, cogimos, cogió, cogiste, cogisteis, coja, cojáis, cojamos, cojan, cojas, cojo — “catch”;*

*conoce, conocé, conoced, conocéis, conocemos, conocen, conocer, conocerá, conocerán, conocerás, conoce-ré, conoceréis, conoceremos, conocería, conoceríais, conoceríamos, conocerían, conocerías, conoce, conocés, conoci, conocía, conocíais, conocíamos, conocían, conocías, conocida, conocidas, conocido, conocidos, cono-ciendo, conociera, conocíerais, conocíeramos, conocieran, conocieras, conociere, conociereis, conocíremos, conocieren, conocieres, conocieron, conociese, conocieseis, conocísemos, conociesen, conocieses, conocimos, conoció, conociste, conocisteis, conozca, conozcáis, conozcamos, conozcan, conozcas, conozco — “know”;*

*cruce, crucé, crucéis, crucemos, crucen, cruces, crucés, crusa, cruzá, cruzaba, cruzabais, cruzábamos, cruzaban, cruzabas, cruzad, cruzada, cruzadas, cruzado, cruzados, cruzáis, cruzamos, cruzan, cruzando, cruzar, cruzara, cruzará, cruzarais, cruzáramos, cruzaran, cruzárán, cruzaras, cruzarás, cruzare, cruzaré, cruzareis, cruzaréis, cruzarem, cruzáremos, cruzaren, cruzares, cruzaría, cruzárais, cruzáramos, cruzárán, cruzarias, cruzaron, cruzas, cruzás, cruzase, cruzaseis, cruzasen, cruzases, cruzaste, cruzasteis, cruzo, cruzó — “cross”;*

*decía, decíais, decíamos, decían, decías, decid, decimos, decir, decís, di, dice, dicen, dices, dicha, dichas, dicho, dichos, diciendo, diga, digáis, digamos, digan, digas, digo, dije, dijera, dijerais, dijéramos, dijeron, dijeras, dijere, dijereis, dijéremos, dijeren, dijeres, dijeron, dijese, dijeseis, dijésemos, dijesen, dijeses, dijimos, dijiste, dijisteis, dijo, dirá, dirán, dirás, diré, diréis, diremos, diría, diríais, diríamos, dirían, dirías — “say”;*

*delinca, delincáis, delincamos, delincan, delincas, delinco, delinque, delinquen, delinques, delinquí, delinquiá, delinquiáis, delinquiámos, delinquián, delinquiás, delinquid, delinquida, delinquidas, delinquito, delinquidos, delinquiendo, delinquier, delinquierais, delinquiéramos, delinquieran, delinquieras, delinquierie, delinquierieis,*

*delinquiéremos, delinquieren, delinquieres, delinquieron, delinquiese, delinquiseis, delinquiésemos, delinquieten, delinquieres, delinquimos, delinquió, delinquir, delinquirá, delinquirán, delinquirás, delinquiré, delinquiréis, delinquiremos, delinquiría, delinquiríais, delinquiríamos, delinquirán, delinquirías, delinquis, delinquiste, delinquistis — “commit a crime”;*

*deshuesa, deshuesá, deshuesan, deshuesas, deshuese, deshuesen, deshueses, deshueso, desosabais, desosábamos, desosaban, desosabas, desosad, desosada, desosadas, desosado, desosados, desosáis, desosamos, desosando, desosar, desosara, desosará, desosarais, desosáramos, desosaran, desosarán, desosaras, desosarás, desosare, desosaré, desosareis, desosaréis, desosaremos, desosáremos, desosaren, desosares, desosaría, desosárias, desosariámos, desosarian, desosarias, desosaron, desosás, desosase, desosaseis, desosásemos, desosasen, desosases, desosaste, desosasteis, desosé, desoseis, desosó — “debone”;*

*discernáis, discernamos, discerní, discernía, discerníais, discerníamos, discernían, discernías, discernid, discernida, discernidas, discernido, discernidos, discerniendo, discerniera, discernierais, discerníeramos, discernieran, discernieras, discerniere, discerniereis, discerniéremos, discernieren, discernieres, discernieron, discerniese, discernieseis, discerniésemos, discerniesen, discernieses, discernimos, discernió, discernir, discernirá, discerníran, discernirás, discerniré, discernréis, discerniremos, discerniría, discerniríais, discerniríamos, discernirían, discernirías, discernís, discerniste, discernisteis, dicierna, diciernan, diciernas, discierne, diciernen, disciernes, dicierno — “discern”;*

*erra, erraba, errabais, errábamos, erraban, errabas, errad, errada, erradas, errado, errados, erráis, erramos, erran, errando, errar, errara, errará, errarais, erráramos, erraran, errarán, erraras, errarás, errare, erraré, errareis, errareis, errarem, errarens, errares, erraria, erráis, erráramos, errarian, errárias, erraron, erras, errase, erraseis, errásemos, errasen, errases, erraste, errasteis, erre, erré, erréis, erremos, ennen, erres, erro, erró, yerra, yerran, yerras, yerre, yerren, yerres, yerro — “miss”;*

*juega, juegan, juegas, juego, juegue, jueguen, juegues, jugá, jugaba, jugábamos, jugaban, jugabas, jugad, jugada, jugadas, jugado, jugados, jugáis, jugamos, jugando, jugar, jugara, jugará, jugarais, jugáramos, jugaran, jugarán, jugadoras, jugarás, jugaré, jugarais, jugaréis, jugaremos, jugáremos, jugaren, jugares, jugaría, jugaríais, jugaríamos, jugarian, jugarías, jugaron, jugás, jugase, jugaseis, jugásemos, jugasen, jugases, jugaste, jugasteis, jugó, jugué, juguéis, juguemos, jugués — “play”;*

*lea, leáis, leamos, lean, leas, lee, leé, leed, leéis, leemos, leen, leer, leerá, leerán, leerás, leeré, leeréis, leeremos, leería, leeríais, leeríamos, leeríam, leeríás, lees, leé, lei, leía, leíais, leíamos, leían, leías, leída, leídias, leído, leídios, leímos, leíste, leísteis, leo, leyendo, leyera, leyerais, leyéramos, leyieran, leyeras, leyere, leyereis, leyéremos, leyeren, leyeres, leyeron, leyese, leyeseis, leyésemos, leyesen, leyeses, leyó — “read”;*

*oí, oía, oíais, oíamos, oían, oías, oíd, oída, oídas, oído, oídos, oiga, oigáis, oigan, oigas, oigo, oímos, oír, oírá, oírán, oíras, oíré, oíréis, oíremos, oíria, oíriais, oíriamos, oírián, oíriás, oís, oíste, oísteis, oye, oyen, oyendo, oyera, oyerais, oyéramos, oyean, oyeras, oyere, oyereis, oyéremos, oyeren, oyeres, oyeron, oyes, oyese, oyeseis, oyésemos, oyesen, oyeses, oyó — “hear”;*

*pedí, pedía, pedíais, pedíamos, pedían, pedías, pedid, pedida, pedidas, pedido, pedidos, pedimos, pedir, pedirá, pedirán, pedirás, pediré, pediréis, pediremos, pediría, pediríais, pediríamos, pedirán, pedirías, pedís, pediste, pedisteis, pida, pidáis, pidamos, pidan, pidas, pide, pidén, pides, pidiendo, pidiera, pidiera, pidierais, pidierámos, pidieran, pidieras, pidiere, pidiereis, pidíremos, pidieren, pidieres, pidieron, pidiese, pidieseis, pidísemos, pidiesen, pidieses, pidió — “request”;*

*pervertí, pervertía, pervertíais, pervertíamos, pervertían, pervertías, pervertid, pervertida, pervertidas, pervertido, pervertidos, pervertimos, pervertir, pervertirá, pervertirán, pervertirás, pervertiré, pervertiréis, pervertíremos, pervertiría, pervertiríais, pervertiríamos, pervertirán, pervertirías, pervertís, pervertiste, pervertisteis, pervierta, perviertan, perviertas, pervierte, pervierten, perviertes, pervierto, pervirtáis, pervirtamos, pervirtás, pervirtiendo, pervirtiera, pervirtiera, pervirtierais, pervirtíramos, pervirtieran, pervirtieras, pervirtiere, pervirtiereis, pervirtíremos, pervirtieren, pervirtieres, pervirtieron, pervirtiese, pervirtieseis, pervirtísemos, pervirtiesen, pervirtieses, pervirtió — “corrupt”;*

*podáis, podamos, poded, podéis, podemos, poder, podía, podíais, podíamos, podían, podías, podidas, podido, podidos, podrá, podrán, podrás, podré, podréis, podremos, podría, podríais, podríamos, podrían, podrías, pude, pudiendo, pudiera, pudiera, pudiera, pudíramos, pudieran, pudieras, pudiere, pudiereis, pudíremos, pudieren, pudieres, pudieron, pudiese, pudieseis, pudísemos, pudiesen, pudieses, pudimos, pudiste, pudisteis, pudo, pueda, puedan, puedas, puede, pueden, puedes, puedo — “can”;*

*prevén, prevendrá, prevendrán, prevendrás, prevendré, prevendréis, prevendremos, prevendría, prevendríais, prevendriamos, prevendrían, prevendrías, prevenga, prevengáis, prevengamos, prevengan, prevengas, prevengo, prevenia, preveníais, preveníamos, prevenían, prevenías, previd, prevenida, previdas, previd, previdos, prevenimos, prevenir, prevenís, previene, previenan, previenes, previne, previniendo, previniera, previnierais, previniéramos, previnieran, previnieras, previniere, previniereis, previniéremos, previnieren, previnieres, previnieran, previniese, previnieseis, previniésemos, previniesen, previnieses, previnimos, previniste, previnisteis, previno — “prevent”;*

*produce, producen, produces, producí, producía, producíais, producíamos, producían, producías, producid, producida, producidas, producido, producidos, produciendo, producimos, producir, producirá, producirán, producirás, produciré, produciréis, produciremos, produciría, produciríais, produciríamos, producirían, producirías, producís, produje, produjera, produjerais, produjéramos, produjeron, produjeras, produjere, produjereis, produjéremos, produjeren, produjeres, produjeron, produjese, produjeseis, produjésemos, produjesen, produjeses, produjimos, produjiste, produjisteis, produjo, produzca, produzcais, produzcamos, produzcan, produzcas, produzcás, produzco — “produce”;*

*queráis, queramos, querás, queré, quered, queréis, queremos, querer, querés, quería, queríais, queríamos, querían, querías, querida, queridas, querido, queridos, queriendo, querrá, querrán, querrás, querré, querréis, querremos, querría, querriais, querriamos, querrián, querriás, quiera, quieran, quieras, quiere, quieren, quieres, quiero, quise, quisiera, quisieraís, quisiéramos, quisieran, quisieras, quisiere, quisiereis, quisiéremos, quisieren, quisieres, quisieron, quisiese, quisieseis, quisiésemos, quisiesen, quisieses, quisimos, quisiste, quisisteis, quiso — “desire”;*

*rehuí, rehuía, rehuíais, rehuíamos, rehuían, rehuías, rehuid, rehuida, rehuidas, rehuido, rehuidos, rehuimos, rehuir, rehuirá, rehuirán, rehuirás, rehuiré, rehuiréis, rehuiremos, rehuiría, rehuiríais, rehuiríamos, rehuirían, rehuirías, rehuís, rehuiste, rehuisteis, rehiya, rehuyás, rehuyamos, rehúyan, rehúyas, rehúye, rehúyen, rehuyendo, rehuyera, rehuyerais, rehuyéramos, rehuyeran, rehuyeras, rehuyere, rehuyereis, rehuyéremos, rehuyeren, rehuyeres, rehuyeron, rehúyes, rehuyese, rehuyeseis, rehuyésemos, rehuyesen, rehuyeses, rehuyó, rehúyo* — “avoid”;

*repón, responderá, responderán, responderás, responderé, responderéis, responderemos, respondería, responderíais, responderíamos, responderían, responderías, repone, reponed, reponéis, reponemos, reponen, reponer, repones, reponga, repongáis, repongamos, repongan, repongas, repongo, reponía, reponíais, reponíamos, reponían, reponías, reponiendo, repuesta, repuestas, repuesto, repuestos, repuse, repusiera, repusierais, repusiéramos, repusieran, repusieras, repusiere, repusiereis, repusiéremos, repusieren, repusieres, repusieron, repusiese, repusieseis, repusiésemos, repusiesen, repusieses, repusimos, repusiste, repusisteis, repuso — “replace”;*

*retén, retendrá, retendrán, retendrás, retendré, retendréis, retendremos, retendría, retendríais, retendríamos, retendrían, retendrías, retened, retenéis, retenemos, retener, retenga, retengáis, retengamos, retengan, retengas, retengo, retenía, reteníais, reteníamos, retenían, retenías, retenida, retenidas, retenido, retenidos, reteniendo, retenie, retenien, retienes, retuve, retuviera, retuvierais, retuviéramos, retuviieran, retuvieras, retuviere, retuviereis, retuviéremos, retuvieren, retuvieres, retuvieron, retuviese, retuvieseis, retuviésemos, retuviesen, retuvieses, retuvimos, retuviste, retuvisteis, retubo — “retain”;*

*sal, saldrá, saldrán, saldrás, saldré, saldréis, saldremos, saldría, saldríais, saldríamos, saldrían, saldrías, sale, salen, sales, salga, salgáis, salgamos, salgan, salgas, salgo, salí, salía, salíais, salíamos, salían, salías, salid, salida, salidas, salido, salidos, saliendo, saliera, salierais, saliéramos, salieran, salieras, saliere, saliereis, saliéremos, salieren, salieres, salieron, saliese, salieseis, saliésemos, saliesen, salieses, salimos, salió, salir, salís, saliste, salisteis* — “leave”;

*trae, traed, traéis, traemos, traen, traer, traerá, traerán, traerás, traeré, traeréis, traeremos, traería, traeríais, traeríamos, traerían, traerías, traes, traía, traíais, traíamos, traían, traías, traída, traidas, traído, traidos, traiga, traigáis, traigamos, traigan, traigas, traigo, traje, trajera, trajerais, trajéramos, trajeren, trajeras, trajere, trajereis, trajéremos, trajeren, trajeres, trajeron, trajese, trajeseis, trajésemos, trajesen, trajeses, trajimos, trajiste, trajisteis, trajo, trayendo — “bring”;*

*ve, vea, veáis, veamos, vean, veas, ved, veía, veíais, veíamos, veían, veías, veis, vemos, ven, veo, ver, verá, verán, verás, veré, veréis, veremos, vería, veríais, veríamos, verían, verías, ves, vi, viendo, viera, vieraís, viéramos, vieran, vieraís, viereis, viéremos, vieren, vieres, vieron, viese, vieseis, viésemos, viesen, vieses, vimos, vio, vista, vistas, viste, vistéis, visto, vistos* — “see”:

*yace, yaced, yacéis, yacemos, yacen, yacer, yacerá, yacerán, yacerás, yaceré, yaceréis, yaceremos, yacería, yaceríais, yaceríamos, yacerían, yacerías, yaces, yaci, yacia, yacíais, yacíamos, yacían, yacías, yacida, yacidas, yacido, yacidos, yaciendo, yaciera, yaciera, yaciéramos, yacieran, yacieras, yaciere, yaciereis, yaciéremos, yacieren, yacieres, yacieron, yacie, yacieis, yaciésemos, yacießen, yacießen, yacimos, yació, yaciste, yacisteis, yaga, yagáis, yagamos, yagan, yagas, yago, yaz, yazca, yazcái, yazcamos, yazcas, yazco, yazga, yazgáis, yazgamos, yazgan, yazgas, yazgo — “recline”*

as input, and receive

{*adquiera, adquieran, adquieras, adquiere, adquieren, adquieres, adquiero, adquiráis, adquiramos, adquirí, adquiría, adquiríais, adquiríamos, adquiríam, adquirías, adquirid, adquirida, adquiridas, adquirido, adquiridos, adquiriendo, adquiriera, adquiriera, adquiríeramos, adquirieran, adquireras, adquiriere, adquiriereis, adquiríremos, adquirieren, adquirieres, adquirieron, adquiriese, adquirieseis, adquirísemos, adquiriesen, adquirieses, adquirimos, adquirió, adquirir, adquirirá, adquirirán, adquirirás, adquiriré, adquiriréis, adquiriremos, adquiriría, adquiriríais, adquiriríamos, adquirirían, adquirís, adquiriste, adquiristeis*},

{*anda, andá, andaba, andabais, andábamos, andaban, andadas, andada, andado, andados, andáis, andamos, andan, andando, andar, andará, andarán, andarás, andaré, andaréis, andaremos, andaría, andariáis, andariamos, andarián, andariás, andas, andás, ande, andéis, andemos, anden, andes, ando, anduve, anduviera, anduviera, anduvíramos, anduvieran, anduvieras, anduviere, anduviereis, anduvíremos, anduvieren, anduvieres, anduvieron, anduviese, anduvieseis, anduvísemos, anduviesen, anduvieses, anduvimos, anduviste, anduvisteis, anduvo*},

{*argüí, argüía, argüíais, argüíamos, argüían, argüías, argüiid, argüida, argüidas, argüido, argüidos, argüímos, argüir, argüirá, argüirán, argüirás, argüiré, argüiréis, argüiremos, argüirla, argürial, argüíramos, argüiríán, argüirias, argüíis, argüiste, argüisteis, arguya, arguyáis, arguyamos, arguyan, arguyas, arguye, arguyen, arguyendo, arguyera, arguyerais, arguyéramos, arguyeran, arguyeras, arguyere, arguyereis, arguyéremos, arguyeren, arguyeres, arguyeron, arguyes, arguyese, arguyeseis, arguyésemos, arguyesen, arguyeses, arguyo, arguyó*},

{*ase, asen, ases, asga, asgáis, asgamos, asgan, asgas, asgo, así, ásia, ásiais, ásiams, ásián, ásías, asid, asida, asidas, asido, asidos, asiendo, asiera, asierais, ásiáramos, ásieran, ásieras, ásiere, ásiereis, ásiéremos, ásieren, ásieres, ásieron, ásiese, ásieseis, ásiésemos, ásiesen, ásieses, ásimos, ásió, asir, asirá, asirán, asirás, asiré, ásiréis, ásiremos, ásiría, ásiríais, ásiríamos, ásiran, ásirías, ásís, ásiste, ásisteis*},

{*avergoncé, avergoncéis, avergoncemos, avergonzaba, avergonzabais, avergonzábamos, avergonzaban, avergonzabas, avergonzad, avergonzada, avergonzadas, avergonzado, avergonzados, avergonzáis, avergonzamos, avergonzando, avergonzar, avergonzara, avergonzará, avergonzara, avergonzaramos, avergonzaran, avergonzarán, avergonzaras, avergonzarás, avergonzare, avergonzaré, avergonzareis, avergonzareis, avergonzaremos, avergonzáremos, avergonzaren, avergonzares, avergonzaría, avergonzaríais, avergonzariamos, avergonzarián, avergonzarias, avergonzaron, avergonzase, avergonzaseis, avergonzásemos, avergonzasen, avergonzases, avergonzaste, avergonzasteis, avergonzó, avergüence, avergüencen, avergüences, avergüenza, avergüenzan, avergüenzas, avergüenzzo*},

{*cabe, cabed, cabéis, cabemos, caben, caber, cabes, cabía, cabíais, cabíamos, cabían, cabías, cabida, cabidas, cabido, cabidos, cabiendo, cabrá, cabrán, cabrás, cabré, cabréis, cabremos, cabría, cabríais, cabríamos, cabrian, cabrias, cupe, cupiera, cupiera, cupíeramos, cupieran, cupieras, cupiere, cupiereis, cupiéremos, cupieren, cupieres, cupieron, cupiese, cupieseis, cupiésemos, cupiesen, cupieses, cupimos, cupiste, cupisteis, cupo, quepa, quepáis, quepamos, quepan, quepas, quepo*},

{*cae, caed, caéis, caemos, caen, caer, caerá, caerán, caerás, caeré, caeréis, caeremos, caería, caeríais, caeríamos, caerían, caerías, caes, caí, caía, caíais, caíamos, caían, caías, caída, caídas, caídos, caiga, caigáis, caigamos, caigan, caigas, caigo, caimos, caíste, caísteis, cayendo, cayera, cayerais, cayéramos, cayeran, cayeras, cayere, cayerais, cayéremos, cayeren, cayeres, cayeren, cayeron, cayese, cayeseis, cayésemos, cayesen, cayeses, cayó*},

{*coge, coged, cogéis, cogemos, cogen, coger, cogerá, cogerán, cogerás, cogeré, cogeréis, cogeremos, cogería, cogeríais, cogeríamos, cogerian, cogerías, cogen, cogí, cogía, cogíais, cogíamos, cogían, cogías, cogida, cogidas, cogido, cogidos, cogiendo, cogiera, cogiera, cogiera, cogíeramos, cogieran, cogieras, cogiere, cogiereis, cogíeremos, cogieren, cogieres, cogieron, cogiese, cogieseis, cogíesemos, cogiesen, cogieses, cogimos, cogió, cogiste, cogisteis, coja, cojáis, cojamos, cojan, cojas, cojo*},

{conoce, conocé, conoced, conocéis, conocemos, conocen, conocer, conocerá, conocerán, conocerás, conoce-ré, conoceréis, conoceremos, conocería, conoceríais, conoceríamos, conocerían, conocerías, conoce-s, conocés, conoci, conocía, conocíais, conocíamos, conocían, conocías, conocida, conocidas, conocido, conocidos, cono-ciendo, conociera, conocierais, conocíeramos, conocieran, conocieras, conociere, conociereis, conociéremos, conocieren, conocieres, conocieron, conociese, conocieseis, conociésemos, conociesen, conocieses, conocimos, conoció, conociste, conocisteis, conozca, conozcás, conozcamos, conozcan, conozcas, conozco},

{cruce, crucé, crucéis, crucemos, crucen, cruces, crucés, cruza, cruzá, cruzabais, cruzábamos, cru-zaban, cruzabas, cruzad, cruzada, cruzadas, cruzado, cruzados, cruzáis, cruzamos, cruzan, cruzando, cruzar, cruzara, cruzará, cruzarais, cruzáramos, cruzaran, cruzarán, cruzaras, cruzarás, cruzare, cruzaré, cruzareis, cruzaréis, cruzaremos, cruzáremos, cruzaren, cruzares, cruzaría, cruzáis, cruzáramos, cruzárian, cruzarias, cruzaron, cruzas, cruzás, cruzase, cruzaseis, cruzásemos, cruzasen, cruzases, cruzaste, cruzasteis, cruzo, cruzó},

{decía, decíais, decíamos, decían, decías, decid, decimos, decir, decís, dice, dicen, dices, dicha, dichas, dicho, dichos, diciendo, diga, digáis, digamos, digan, digas, digo, dije, dijera, dijerais, dijéramos, dijeron, dijeras, dijere, dijereis, dijéremos, dijeren, dijeres, dijeron, dijese, dijeseis, dijésemos, dijesen, dijeses, dijimos, dijiste, dijisteis, dijo, dirá, dirán, dirás, diré, diréis, diremos, diría, diráis, diríamos, dirían, dirás, discernáis},

{delinca, delincáis, delincamos, delinca, delinco, delinque, delinques, delinquí, delinquía, delinquíais, delinquíamos, delinquían, delinquías, delinquid, delinquida, delinquidas, delinquido, delinquidos, delinquiendo, delinquiera, delinquiera, delinquíeramos, delinquieran, delinquieras, delinquiere, delinquiereis, delinquíeremos, delinquieren, delinquieres, delinquieron, delinquiese, delinquieseis, delinquiésemos, delinquie-sen, delinquieses, delinquimos, delinquió, delinquir, delinquirá, delinquirán, delinquirás, delinquiré, delinqui-reis, delinquiremos, delinquiría, delinquiria, delinquiríamos, delinquirian, delinquirías, delinquís, delinquiste, delinquisteis},

{deshuesa, deshuesá, deshuesan, deshuesas, deshuese, deshuesen, deshueses, deshueso, desosaba, desosabais, desosábamos, desosaban, desosabas, desosad, desosada, desosadas, desosado, desosados, desosáis, desosamos, desosando, desosar, desosara, desosará, desosarais, desosáramos, desosaran, desosárán, desosaras, desosarás, desosare, desosaré, desosareis, desosáreis, desosaremos, desosáremos, desosaren, desosares, desosaría, desosárias, desosáriamos, desosárian, desosariás, desosaron, desosás, desosase, desosaseis, desosásemos, desosasen, desosases, desosaste, desosasteis, desosé, desoséis, desosemos, desosó},

{di},

{discernamos, discerní, discernía, discerníais, discerníamos, discernían, discernías, discernid, discernida, dis-cernidas, discernido, discernidos, discerniendo, discerniera, discerniera, discerníeramos, discernieran, discernieras, discerniere, discerniereis, discerníeremos, discernieren, discernieres, discernieron, discerniese, discernieseis, discerníe-semos, discerniesen, discernieses, discernimos, discernió, discernir, discernirá, discernirán, discernirás, discerniré, discernréis, discerniremos, discerniría, discerníerai, discerníeramos, discerníeran, discernírias, discernís, discerniste, discernisteis, discierna, disciernan, disciernas, discierne, disciernen, disciernes, discierno},

{erra, erraba, errabais, errábamos, erraban, errabas, errad, errada, erradas, errado, errados, erráis, erramos, erran, errando, errar, errara, errará, errara, errara, erráramos, erraran, errará, erraras, errará, erraré, errareis, errareis, erráremos, erráremos, erraren, errares, erraria, erráis, erráramos, errarian, errárias, erra-ron, erras, errase, erraseis, errásemos, errasen, errases, erraste, errasteis, erre, erré, erréis, erremos, erren, erres, erro, erró, yerra, yerran, yerras, yerre, yerren, yerres, yerro},

{juega, juegan, juegas, juego, juegue, jueguen, juegues, jugá, jugaba, jugabais, jugaban, jugabas, jugad, jugada, jugadas, jugado, jugados, jugáis, jugamos, jugando, jugar, jugara, jugará, jugarais, jugáramos, jugaran, jugarán, jugaras, jugarás, jugaré, jugarais, jugaréis, jugaréis, jugarémos, jugarémos, jugaren, jugares, jugaría, jugaríais, jugaríamos, jugarían, jugarías, jugaron, jugás, jugase, jugaseis, jugásemos, jugasen, jugases, jugaste, jugasteis, jugó, jugué, juguéis, juguemos, jugués},

{lea, leáis, leamos, lean, leas, lee, leé, leed, leéis, leemos, leen, leer, leerá, leerán, leerás, leeré, leeréis, leeremos, leería, leeríais, leeríamos, leerían, leerías, lees, leés, leí, leía, leíais, leíamos, leían, leías, leída, leídas, leído, leí-dos, leímos, leíste, leísteis, leo, leyendo, leyera, leyerais, leyéramos, leyeron, leyeras, leyere, leyereis, leyéremos, leyeren, leyeres, leyeron, leyese, leyeseis, leyésemos, leyesen, leyeses, leyó},

{oí, oía, oíais, oíamos, oían, oítas, oíd, oída, oídas, oído, oídos, oiga, oigáis, oigamos, oigan, oigas, oigo, oímos, oír, oirá, oirán, oirás, oiré, oiréis, oiremos, oiría, oiríais, oiríamos, oirían, oirías, oís, oíste, oísteis, oye, oyen, oyendo, oyera, oyerais, oyéramos, oyeran, oyeras, oyere, oyereis, oyéremos, oyeren, oyeres, oyeron, oyes, oyese, oyeseis, oyésemos, oyesen, oyeses, oyó},

{pedí, pedia, pedíais, pedíamos, pedian, pedias, pedid, pedida, pedidas, pedido, pedidos, pedimos, pedir, pedirá, pedirán, pedirás, pediré, pediréis, pediremos, pediría, pediríais, pediríamos, pedirán, pedirías, pedís, pediste, pedisteis, pida, pidáis, pidamos, pidan, pidas, pide, pidan, pides, pidiendo, pidiera, pidierais, pidíramos, pidieran, pidieras, pidiere, pidiereis, pidíremos, pidieren, pidieres, pidieron, pidiese, pidieseis, pidísemos, pidiesen, pidieses, pidío, pido},

{pervertí, pervertía, pervertíais, pervertíamos, pervertían, pervertías, pervertid, pervertida, pervertidas, pervertido, pervertidos, pervertimos, pervertir, pervertirá, pervertirán, pervertirás, pervertiré, pervertiréis, pervertíremos, pervertiría, pervertírais, pervertíramos, pervertirán, pervertirías, pervertís, pervertiste, pervertisteis, pervierta, perviertan, perviertas, pervierte, pervierten, perviertes, pervierto, pervirtáis, pervirtamos, pervirtás, pervirtiendo, pervirtiera, pervirtiera, pervirtíramos, pervirtieran, pervirtieras, pervirtiere, pervirtiereis, pervirtiéremos, pervirtieren, pervirtieres, pervirtieron, pervirtiese, pervirtieseis, pervirtísemos, pervirtiesen, pervirtieses, pervirtio},

{podáis, podamos, poded, podéis, podemos, poder, podía, podíais, podíamos, podían, podida, podidas, podido, podidos, podrá, podrán, podrás, podré, podréis, podremos, podría, podríais, podríamos, podrían, podrías, pude, pudiendo, pudiera, pudiera, pudíramos, pudieran, pudieras, pudiere, pudiereis, pudíremos, pudieren, pudieres, pudieron, pudiese, pudieseis, pudísemos, pudiesen, pudieses, pudimos, pudiste, pudisteis, pudo, pueda, puedan, puedas, puede, pueden, puedes, puedo},

{prevén, prevendrá, prevendrán, prevendrás, prevendré, prevendréis, prevendremos, prevendría, prevendríais, prevendríamos, prevendrían, prevendrías, prevenga, prevengáis, prevengamos, prevengan, prevengas, prevengo, prevenía, preveníais, preveníamos, prevenían, prevenías, prevenid, prevenida, prevenidas, prevenido, prevenidos, prevenimos, prevenir, prevenís, previene, previenen, previenes, previne, previniendo, previniera, previnie, previníramos, previnieran, previnieras, previnire, previnireis, previníremos, previnieren, previnieres, previnieron, previniese, previnieseis, previnísemos, previniesen, previnieses, previnimos, previniste, previnisteis, previno},

{produce, producen, produces, producí, producía, producíais, produciamos, producían, producías, producid, producida, producidas, producido, producidos, produciendo, producimos, producir, producirá, producirán, producirás, produciré, produciréis, produciremos, produciría, produciríais, produciríamos, producirían, producirías, producis, produje, produjera, produjerais, produjéramos, produjeren, produjeras, produjere, produjereis, produjéremos, produjeren, produjeres, produjeron, produjeres, produjeron, produjese, produjeseis, produjésemos, produjesen, produjeses, produjimos, produjiste, produjisteis, produjo, produzca, produzcais, produzcamos, produzcan, produzcas, produzcás, produzco},

{queráis, queramos, querás, queré, quered, queréis, queremos, querer, querés, quería, queríais, queríamos, querían, querías, querida, queridas, querido, queridos, queriendo, querrá, querrán, querrás, querré, querréis, querremos, querría, queríais, queríamos, querían, querías, quiera, quieran, quieras, quiere, quieren, quieres, quiero, quise, quisiera, quisiera, quisíramos, quisieran, quisieras, quisiere, quisiereis, quisíremos, quisieren, quisieres, quisieron, quisiese, quisieseis, quisísemos, quisiesen, quisieses, quisimos, quisiste, quisisteis, quiso},

{rehuí, rehuía, rehuíais, rehuíamos, rehuían, rehuías, rehuid, rehuida, rehuidas, rehuido, rehuidos, rehuimos, reuir, reuirá, reuirán, reuirás, reuiré, reuiréis, reuiramos, reuiría, reuiríais, reuiríamos, reuirían, reuirías, rehuís, rehuiste, rehuisteis, rehuya, rehuyás, rehuyamos, rehúyan, rehuyás, rehúye, rehúyen, rehuyendo, rehuyera, rehuyerais, rehuyéramos, rehuyeren, rehuyeras, rehuyere, rehuyereis, rehuyéremos, rehuyeren, rehuyeres, rehuyeron, rehúyes, rehuyeseis, rehuyésemos, rehuyesen, rehuyeses, rehuyó, rehúyo},

{repón, respondrá, respondrán, respondrás, respondré, respondréis, respondremos, respondría, respondríais, respondriamos, respondrían, respondrías, repone, reponed, reponéis, reponemos, reponen, reponer, repones, reponga, repongáis, repongamos, repongan, repongas, repongo, reponía, reponíais, reponíamos, reponían, reponías, reponiendo, repuesta, repuestas, repuesto, repuestos, repuse, repusiera, repusiera, repusíramos, repusieran, repusieras, repusiere, repusiereis, repusíremos, repusieren, repusieres, repusieron, repusiese, repusieseis, repusísemos, repusiesen, repusieses, repusimos, repusiste, repusisteis, repuso},

{retén, retendrá, retendrán, retendrás, retendré, retendréis, retendremos, retendría, retendríais, retendriamos, retendrían, retendrías, retened, retenéis, retenemos, retener, retenga, retengáis, retengamos, retengan, retengas, retengo, retenía, reteníais, reteníamos, retenían, retenías, retenida, retenidas, retenido, retenidos, reteniendo, reteniene, retienes, retienen, retuve, retuviera, retuvierais, retuvíeramos, retuvieran, retuvieras, retuviere, retuviereis, retuvíeremos, retuvieren, retuvieres, retuvieron, retuviese, retuvieseis, retuvísemos, retuviesen, retuvieses, retuvimos, retuviste, retuvisteis, retuvo},

{sal, saldrá, saldrán, saldrás, saldré, saldréis, saldremos, saldría, saldríais, saldríamos, saldrían, saldrías, sale, salen, sales, salga, salgáis, salgamos, salgan, salgas, salgo, salí, salía, salíais, salíamos, salían, salías, salid, salida, salidas, salido, salidos, saliendo, saliera, salieraís, salíeramos, salieran, salieras, saliere, saliereis, salíremos, salieren, salieres, salieron, saliese, salieseis, saliésemos, saliesen, salieses, salimos, salió, salir, salís, saliste, salisteis},

{trae, traed, traéis, traemos, traen, traer, traerá, traerán, traerás, traeré, traeréis, traeremos, traería, traeríais, traeríamos, traerían, traerías, traes, traia, traíais, traíamos, traian, traías, traída, traídas, traído, traídos, traiga, traigáis, traigamos, traigan, traigas, traigo, traje, trajera, trajerais, trajéramos, trajeren, trajeras, trajere, trajereis, trajéremos, trajeren, trajeres, trajeron, trajese, trajeseis, trajésemos, trajesen, trajeses, trajimos, trajiste, trajisteis, trajo, trayendo},

{ve},

{vea, veáis, veamos, vean, veas, ved, veía, veíais, veíamos, veían, veías, veis, vemos, ven, veo, ver, verá, verán, verás, veré, veréis, veremos, vería, veríais, veríamos, verían, verías, ves, vi, viendo, viera, vieraís, viéramos, vieran, vieraís, viere, viereis, viéremos, vieran, vieres, vieron, viese, vieseis, viésemos, viesen, vieses, vimos, vio, vista, vistas, viste, visteis, visto, vistos},

{yace, yaced, yacéis, yacemos, yacen, yacer, yacerá, yacerán, yacerás, yaceré, yaceréis, yaceremos, yacería, yaceríais, yaceríamos, yacerian, yacerías, yaces, yaci, yacia, yaciais, yacíamos, yacían, yacias, yacida, yacidas, yacido, yacidos, yaciendo, yaciera, yacieraís, yaciéramos, yacieran, yacieras, yaciere, yaciereis, yaciéremos, yacieren, yacieres, yacieron, yacie, yacieis, yaciésemos, yacieen, yaciees, yacimos, yació, yaciste, yacisteis, yaga, yagáis, yagamos, yagan, yagas, yago, yaz, yazca, yazcái, yazcamos, yazcas, yazco, yazga, yazgáis, yazgamos, yazgan, yazgas, yazgo}

as output. Here, the ambiguous words *di* and *ve* become orphans in our algorithm: *di* is a conjugated form of both *dar* “give” and *decir* “say”; *ve* is a conjugated form of both *ir* “go” and *ver* “see”.

*Example 6.11.4.* In Fig. S7, we further apply the aforementioned word clustering algorithm to topic extraction and machine translation (see Table S1 for text source).

Note that the machine translation task is particularly challenging in the example shown, for at least two reasons: (1) The translator has merged the name *Elizabeth* and the nickname *Lizzy* of the protagonist into one; (2) The Spanish noun *casa* “house” is spelt the same way as a conjugated form of the Spanish verb *casar* “marry”. Despite these disadvantages, the topics that are unrelated to these challenges are correctly translated.

In Fig. S7b, we consider the Spanish word *capital* “capital (city)” an exact match to the English word *town*, because the latter refers to London in *Pride and Prejudice*.

In Spanish, both *bella* and *hermosa* mean “beautiful” and modify feminine nouns. It thus comes to no surprise that the row for “beautiful” in Fig. S7b contains two hot spots that correspond to these two Spanish adjectives.

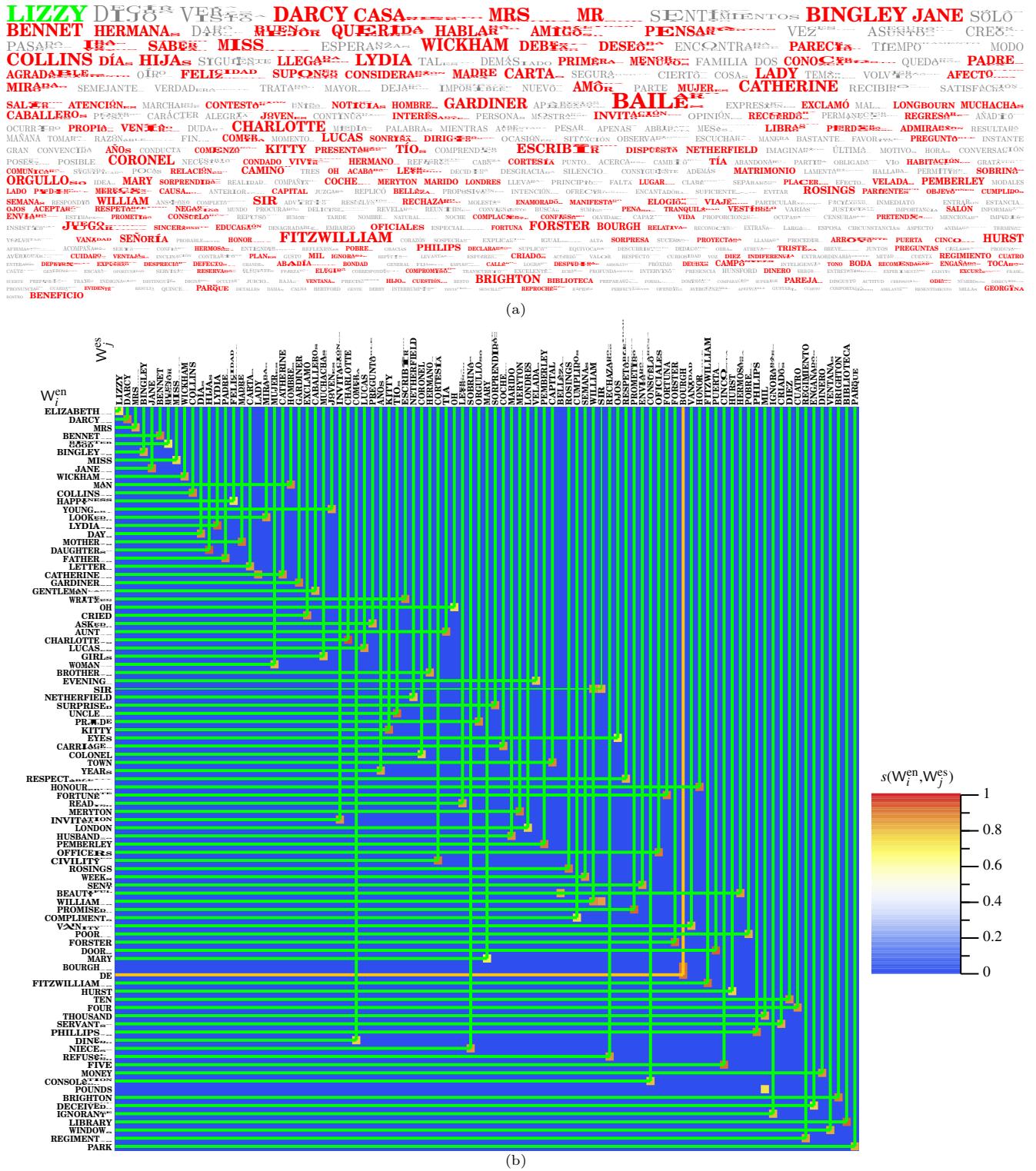


Fig. S7. Text mining in Spanish. (a) Statistically identified topics ( $n_{ii} \geq 20$ ) in a Spanish version of *Pride and Prejudice*, with the same color encoding scheme as Fig. S3. (b) Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{es}})$  between selected topics in English and Spanish versions of *Pride and Prejudice*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) ranking of the similarity score for the optimal node. Green (resp. amber) cross-hair indicates an exact (resp. a close but non-exact) match.

## 6.2 Modified Porter stemming algorithm for French

The French verbs for “to be” and “to have” are highly irregular. These verbs are also used as auxiliaries in various tenses and moods of arbitrary verbs, similar to their counterparts in English. All their conjugated forms (including the infinitive forms) are treated as stop words.

*a, ai, aie, aient, aies, ait, as, aura, aurai, auraient, aurais, aurait, auras, aurez, auriez, aurions, aurons, auront, avaient, avais, avait, avez, aviez, avions, avoir, avons, ayant, ayez, ayons, eu, eûmes, eurent, eus, eusse, eussent, eusses, eussiez, eussions, eut, eût, eûtes, ont — “have”;*

*es, est, étaient, étais, était, étant, été, êtes, étiez, étions, être, fûmes, furent, fus, fusse, fussent, fusses, fussiez, fussions, fut, fût, fûtes, sera, serai, seraient, serais, serait, seras,erez, seriez, serions, serons, seront, soient, sois, soit, sommes, sont, soyez, soyons, suis — “be”.*

In addition, we treat the following verbs as stop words.

*faire, fais, faisaient, faisais, faisait, faisant, faisiez, faisions, fait, faites, fasse, fassent, fasses, fassiez, fassions, fera, ferai, feraient, serais, ferait, feras, ferez, ferions, ferons, feront, fimes, furent, fis, fissse, fisssent, fissses, fissiez, fisssons, fit, fit, fites, font — “do”;*

*fallait, falloir, fallu, fallut, faudra, faudrait, faut* — “be necessary”;

*peut, peuvent, peux, pourra, pourrai, pourraient, pourrais, pourrait, pourras, pourrez, pourriez, pourrions, pourrons, pourront, pouvaient, pouvais, pouvait, pouvant, pouvez, pouviez, pouvions, pouvoir, pouvons, pu, puisse, puissent, puisses, puissiez, puissions, pûmes, parent, pus, pusse, pussent, pusses, pussiez, pussions, put, pût, pûtes — “can”*

**Definition 6.12** (French stop words). If a word belongs to the following list<sup>73</sup>:

a, à, afin, ai, aie, aient, aies, ainsi, ait, alors, après, arrière, as, assez, au, aucun, aucune, aucunes, aucuns, aura, aurai, auraient, aurais, aurait, auras, aurez, auriez, aurions, aurons, auront, aussi, autant, autour, autre, autrefois, autres, aux, avaient, avais, avait, avant, avec, avez, aviez, avions, avoir, avons, ayant, ayez, ayons, beaucoup, bien, bientôt, c, ça, ça, car, ce, ceci, cela, celà, celle, celles, celui, cependant, ces, cet, cette, ceux, chacun, chacune, chaque, chez, ci, combien, comme, comment, contre, d, dans, de, dedans, dehors, déjà, depuis, dernier, dernière, dernières, derniers, derrière, des, dès, dessous, dessus, devaient, devais, devait, devant, devez, deviez, devions, devoir, devons, devra, devrai, devraient, devrais, devras, devrez, devriez, devrions, devrons, devant, dois, doit, doive, doivent, doives, donc, dont, du, dû, dûmes, durant, durent, dus, dusse, dussent, dusses, dussiez, dussions, dut, dût, dûtes, elle, elles, en, encore, envers, es, est, et, étaient, étaient, était, étant, été, êtes, étiez, étions, être, eu, eue, eues, eûmes, eurent, eus, eusse, eussent, eusses, eussiez, eussions, eut, eût, eûtes, eux, faire, fais, faisaient, faisais, faisait, faisiez, faisions, fait, faite, faites, faits, fallait, falloir, fallu, fallut, fasse, fassent, fasses, fassiez, fassions, faudra, faudrait, faut, fera, ferai, feraient, ferais, ferait, feras, ferez, feriez, ferions, ferons, feront, fîmes, firent, fis, fissee, fissent, fisses, fissiez, fissions, fit, fît, fites, fois, font, fûmes, furent, fus, fusse, fussent, fusses, fussiez, fussions, fut, fût, fûtes, hors, ici, il, ils, j, jamais, je, jusqu, jusque, juste, justement, l, la, là, laquelle, le, lequel, les, lesquelles, lesquels, leur, leurs, lors, lorsqu, lorsque, lui, m, ma, maintenant, mais, malgré, me, même, mêmes, mes, mien, mienne, miennes, miens, moi, moindre, moins, mon, n, ne, ni, non, nos, notre, nôtre, nôtres, nous, nul, nulle, nulles, nuls, on, ont, ou, où, oui, par, parce, parfois, parmi, partout, pas, pendant, personne, peu, peut, peuvent, peux, plupart, plus, plusieurs, plutôt, pour, pourquoi, pourra, pourrai, pourraient, pourrais, pourrait, pourras, pourrez, pourriez, pourrions, pourrons, pourront, pourtant, pouvaient, pouvais, pouvait, pouvant, pouvez, pouviez, pouvions, pouvoir, pouvons, presqu, presque, prochain, prochaine, prochaines, prochains, propos, pu, puis, puisque, puisse, puissent, puisses, puissiez, puissions, pûmes, purent, pus, pusse, pussent, pusses, pussiez, pussions, put, pût, pûtes, qu, quand, quant, que, quel, quelconque, quelle, quelles, quelqu, quelque, quelquefois, quelques, quels, qui, quiconque, quoi, quoique, rien, s, sa, sans, sauf, se, selon, sera, serai, seraient, serais, serait, seras, serez, seriez, serions, serons, seront, ses, si, sien, sienne, siennes, siens, sitôt, soi, soient, sois, soit, sommes, son, sont, sous, souvent, soyez, soyons, suis, sur, surtout, t, ta, tandis, tant, tantôt, te, tel, telle, tellement, telles, tels, tes, tien, tienne, tiennes, tiens, toi, ton, toujours, tous, tout, toute, toutefois, toutes, travers, très, trop, tu, un, une, unes, uns, vers, voici, voilà, vos, votre, vôtre, vôtres, vous, y,

then we consider it a French stop word. All the French stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list.  $\square$

<sup>73</sup>Our list of French stop words is based on [snowball.tartarus.org/algorithms/french/stop.txt](http://snowball.tartarus.org/algorithms/french/stop.txt), with extensive modifications that roughly match their counterparts in English.

The French verbs for “to come” and “to go” also have highly irregular conjugations:<sup>74</sup>

*aille, aillent, ailles, alla, allai, allaient, allais, allait, allâmes, allant, allas, allasse, allassent, allasses, allassiez, allassions, allât, allâtes, allé, allés, aller, allèrent, allez, alliez, allions, allons, ira, irai, iraient, irais, irait, iras, irez, iriez, irions, irons, iront, va, vais, vas, vont — “go”;*

*venaient, venais, venait, venant, venez, veniez, venions, venir, venons, venu, venus, viendra, viendrai, viendrait, viendras, viendrait, viendras, viendrez, viendriez, viendrions, viendrons, viendront, vienne, viennent, viennes, viens, vient, vîmes, vinrent, vins, vînse, vînsent, vînsses, vînssiez, vînssions, vînt, vînt, vîntes — “come”.*

We are not going to treat them as stop words in French, so as to be consistent with their counterparts in other languages. However, to prepare for clustering of content words, we will refer to these two groups of strings as **goFrench** and **comeFrench**, respectively.

### 6.2.1 Effective spelling and essential root

**Definition 6.13** (French Vowel Extensions). Hereafter in §6.2, the symbol  $\mathbf{V}^*$  stands for any member from the list  $\{a, \hat{a}, e, \hat{e}, i, \hat{i}, o, \hat{o}, \hat{u}\}$ , the so-called French vowel extensions. In line with the multiplicity notations introduced in Definition 3.3, the symbol  $\mathbf{V}_m^*$  stands for a text string formed by consecutive appearance of one or more (not necessarily identical or distinct) members from the set of French vowel extensions.

Dual to the notations above, the symbol  $\mathbf{C}^*$  stands for any character that does not belong to the list  $\{a, \hat{a}, e, \hat{e}, i, \hat{i}, o, \hat{o}, u, \hat{u}\}$ , and  $\mathbf{C}_{m_0}^*$  stands for a text string formed by consecutive appearance of zero or more (not necessarily identical or distinct) characters that do not belong to the same list.  $\square$

**Algorithm 6.14** (French effective spelling). *For a French word  $\hat{\sigma}$ , its effective spelling  $\text{EffSpell}(\hat{\sigma})$  is constructed in six sequential steps:*

(1) *Replace*<sup>75</sup>

<i>abC*~</i>	<i>amabilité~</i>	<i>chance~</i>	<i>charlotte</i>	<i>défX~</i>	<i>demi~</i>	<i>déTECTX~</i>	<i>detei~</i>
<i>bαaC*</i>	<i>ami</i>	<i>βance</i>	<i>chapλωτ</i>	<i>δφ</i>	<i>dēmi</i>	<i>τδεκτ</i>	<i>τβει</i>
<i>détend~</i>	<i>détestX~</i>	<i>diab(Ø o)l~</i>	<i>discour~</i>	<i>foyer~</i>	<i>frém~</i>	<i>guet~</i>	<i>imposs~</i>
<i>ρλaξ</i>	<i>ηατε</i>	<i>δiaβλ</i>	<i>δσκουρ</i>	<i>φοκαρε</i>	<i>φρμ</i>	<i>γυετ</i>	<i>γuid</i>
<i>list~</i>	<i>litan~</i>	<i>litt~</i>	<i>livrX~</i>	<i>louis~</i>	<i>message~</i>	<i>muet~</i>	<i>muf~</i>
<i>λιστ</i>	<i>λιταν</i>	<i>λιττ</i>	<i>λuiσ</i>	<i>μεβαγε</i>	<i>μιωττ</i>	<i>mf</i>	<i>μιστι</i>
<i>présent(c i)~</i>	<i>quotidienX</i>	<i>récemment</i>	<i>régim~</i>	<i>région~</i>	<i>remerc~</i>	<i>rou(sse x)~</i>	<i>salad~</i>
<i>prseñt</i>	<i>jour</i>	<i>récent</i>	<i>ρεγμ</i>	<i>εργῆ</i>	<i>ρμερκ</i>	<i>rouge</i>	<i>saλδ</i>
<i>sommet~</i>	<i>tableau</i>	<i>taill</i>	<i>témér~</i>	<i>term~</i>	<i>triste~</i>	<i>trois~</i>	<i>volt~</i>
<i>σομμετ</i>	<i>ταβλεαυ</i>	<i>tail</i>	<i>τεμερ</i>	<i>terμ</i>	<i>trιστε</i>	<i>3trois</i>	<i>βολτ</i>
<i>(Ø ma mes)dame(Ø s)</i>	<i>(Ø ma mes)demoiselle(Ø s)</i>		<i>bessie</i>	<i>hui</i>	<i>joyau(Ø x)</i>	<i>mlle(Ø s)</i>	<i>mme(Ø s)</i>
<i>lady</i>		<i>miβ</i>	<i>βεσθιε</i>	<i>αujourd</i>	<i>joaillerie</i>	<i>miβ</i>	<i>lady</i>
<i>nu(Ø e)(Ø s)</i>	<i>rite(Ø s)</i>		<i>yie(Ø s)</i>		<i>~X<sup>ε</sup>(V<sup>*</sup>C<sub>m_0</sub><sup>*</sup>)ement(Ø s)</i>		<i>~an(ce i)</i>
<i>νιδε</i>	<i>ριτ</i>		<i>vivre</i>		<i>Xu</i>		<i>ants</i>

(2) *Replace*

<i>(maman matern mère)X~</i>	<i>(meilleur mieux)~</i>	<i>C<sup>*</sup>y</i>
<i>μαμα</i>	<i>μειλλευ</i>	<i>C<sup>*</sup>η</i>
<i>X<sup>ε</sup>(aper con per)ceptX~</i>	<i>a<sup>X<sup>ε</sup>(cc ss tt)~</sup></i>	<i>accrûX~</i>
<i>Xcev</i>	<i>xX<sup>[1]</sup></i>	<i>accr<sup>υ</sup>issons</i>
<i>bon(Ø ne)(Ø s)~</i>	<i>cess~</i>	<i>concl(u ū)~</i>
<i>μeilleur</i>	<i>ξεστ</i>	<i>courage~</i>
		<i>ction</i>
		<i>dame~</i>
		<i>dat~</i>
		<i>déceptX~</i>
		<i>déte~</i>
		<i>déτ</i>
		<i>det</i>

<sup>74</sup>We have chosen not to include the feminine forms of past participles (*allée, allées, venue, venues*) in the lists below, because their functions as nouns deviate from the original verbs. Similarly, when the present participles *allant* and *venant* are used as adjectives, their meanings no longer align perfectly with the respective verbs. Therefore, we do not include the feminine and plural forms of these adjectives (such as *allante* and *allants*) in our lists of verb forms.

<sup>75</sup>To avoid confusion of Greek nu with Latin vee, we write the former as ν in substitution rules.

<i>differend~</i>	<i>difficulté~</i>	<i>distant~</i>	<i>dort~</i>	<i>éclor</i>	<i>eill</i>	<i>excl(u û)~</i>	<i>faço~</i>	<i>fraternX~</i>	<i>front~</i>	<i>gén~</i>
<i>different</i>	<i>difficile</i>	<i>distance</i>	<i>dorm</i>	<i>eclos</i>	<i>eil</i>	<i>qwexclu</i>	<i>φaço</i>	<i>frère</i>	<i>front</i>	<i>gzan</i>
<i>hiv~</i>	<i>hommage~</i>	<i>incl(u û)~</i>	<i>ivité</i>	<i>kC*</i>	<i>matin~</i>	<i>matrimonX~</i>	<i>mer~</i>	<i>mesu~</i>		
<i>hib</i>	<i>ηommaγe</i>	<i>winclu</i>	<i>ive</i>	<i>κC*</i>	<i>μatiū</i>	<i>marier</i>	<i>mer</i>	<i>meσw</i>		
<i>meur(ent)s</i>	<i>misér~</i>	<i>mois</i>	<i>moment~</i>	<i>muV*~</i>	<i>muet~</i>	<i>mul~</i>	<i>mur~</i>	<i>occ~</i>	<i>op~</i>	<i>parent~</i>
<i>mort</i>	<i>μiσer</i>	<i>μoiσ</i>	<i>μoμeñt</i>	<i>μwu</i>	<i>μuet</i>	<i>μuλ</i>	<i>myρ</i>	<i>cc</i>	<i>zp</i>	<i>pareñt</i>
<i>petit~</i>	<i>peur~</i>	<i>pré</i>	<i>prés~</i>	<i>qua(d)t)r</i>	<i>querell~</i>	<i>quest~</i>	<i>rd</i>	<i>rdin</i>	<i>reC*~</i>	
<i>petit</i>	<i>πeυρ</i>	<i>prae</i>	<i>rπeσ</i>	4	<i>qreλ</i>	<i>qqeστ</i>	<i>ρδ</i>	<i>ρρδiū</i>	<i>rC*</i>	
<i>réceptX·</i>	<i>rey~</i>	<i>rn</i>	<i>rst</i>	<i>saisi~</i>	<i>salu~</i>	<i>sc</i>	<i>sit~</i>	<i>sororX·</i>	<i>spoir</i>	<i>squ</i>
<i>rcefer</i>	<i>pey</i>	<i>r̄v</i>	<i>ρστ</i>	<i>saiσi</i>	<i>σalu</i>	<i>σc</i>	<i>sit</i>	<i>sæur</i>	<i>sper</i>	<i>σque</i>
								<i>tot</i>	<i>trava~</i>	<i>vers~</i>
								<i>trava</i>	<i>βerσ</i>	<i>zyrt</i>
									<i>φvite</i>	<i>vitρ</i>
<i>voiture</i>	<i>voyag</i>	<i>vulg~</i>			<i>(messieurs monsieur mr)</i>				<i>(mistress mme mrs)</i>	
<i>voiture</i>	<i>βoηay</i>	<i>βuły</i>			<i>umur</i>				<i>μadame</i>	
<i>Xε(Ø con de re)comeFrench</i>	<i>ans</i>	<i>bal(Ø s)</i>		<i>courent</i>	<i>dors</i>	<i>écri(V* è é)X</i>		<i>filleul(Ø s)</i>	<i>gens</i>	
<i>XfcomeφXi</i>	<i>an</i>	<i>βal</i>		<i>courez</i>	<i>dorm</i>	<i>fjeλ</i>		<i>fils</i>	<i>γenoσ</i>	
<i>goFrench</i>	<i>heure(Ø s)</i>	<i>mari(Ø s)</i>	<i>maria</i>	<i>mot(Ø s)</i>	<i>né(Ø e)(Ø s)</i>	<i>nez</i>	<i>papa(Ø s)</i>	<i>parc(Ø s)</i>	<i>vieux</i>	
<i>fgoφ</i>	<i>zheure</i>	<i>μari</i>	<i>μaria</i>	<i>μot</i>	<i>naître</i>	<i>ñεζ</i>	<i>père</i>	<i>pparc</i>	<i>vieil</i>	
	<i>voit</i>	<i>yeux</i>			<i>~χ̄(e)aux</i>			<i>~(imes îtes)</i>		
	<i>verra</i>	<i>œil</i>			<i>χ̄al</i>			<i>i</i>		
<i>~Xε(V* Cm0*)ab(ilité le)(Ø s)</i>	<i>~eau(Ø té tés x)</i>	<i>~ial(Ø e ité)(Ø s)</i>	<i>~iaux</i>	<i>~k</i>	<i>~m(f i)t</i>	<i>~ss</i>	<i>~trice(Ø s)</i>	<i>~urent</i>	<i>~y</i>	
<i>Xu</i>	<i>eler</i>		<i>e</i>	<i>e</i>	<i>κ</i>	<i>mis</i>	<i>β</i>	<i>teur</i>	<i>u</i>	<i>η</i>

(3) *Do* (é|è|ë) → *e*, (i|y) → *ii*, (*k|q*) → *κ*, *Xε(a|o|u)r* → *Xρ*, *ç* → *ce*, *disC\*~* → *dssC\**, *gn* → *nad*, *jourñée~* → *joup*, *naîtrX~* → *naquit*, *ô* → *os*, *sXε(in|ui)* → *βX*, *sûr* → *βur*, *ü* → *u*, *û* → *uû*, *mypent* → *mouvez*, *~au(Ø|x)* → *al*, *~c* → *k*, *~eux* → *eu*, *iXε(me|~t)* → *iX*, *~in(s|t)* → *in*;

(4) Replace

<i>n(d t)r</i>	<i>χv</i>	<i>pt</i>	<i>C*nee</i>	<i>cXε(a o u C*)</i>	<i>dsss~</i>	<i>entier~</i>	<i>entreν~</i>	<i>gXε(e i)</i>	<i>qu</i>	<i>ui</i>	<i>vieilX</i>	<i>~(ez nt)</i>	<i>~eup(Ø s)</i>
<i>n</i>	<i>χf</i>	<i>pt</i>	<i>C*</i>	<i>kX</i>	<i>diss</i>	<i>εv̄tier</i>	<i>ēv̄tp̄ev</i>	<i>jX</i>	<i>k</i>	<i>uy</i>	<i>oλδ</i>	<i>i</i>	<i>euse</i>

(5) *Do ell* → *el*, *enn* → *en*, *ett* → *et*, *Vε(a|e|o)i* → *Vy*, *Vε(a|e|o)u* → *Vw*.

(6) *Do Xε(m|p|v)ew~* → *Xow*, *kow~* → *kkow*, *kru~* → *kroy*, *mu~* → *mowf*, *v(i|i|rey|i|l|u)~* → *voyre*, *bow(s|t)* → *bowyllons*, *vawt* → *vawdrons*.

(7) *Do Xε(trafa|konse)y(l|s)~* → *Xyl*, *mowfX~* → *mowf*, *Xε(p|s)u(rey|s|i)* → *Xuss*, *mow(s|t)* → *mowf*, *~ewse(Ø|s)* → *e*.

**Definition 6.15** (French protected range). Let  $\hat{\sigma}$  be a text string derived from a French word, its protected range  $\text{ProtRg}(\hat{\sigma})$  is an integer determined as follows:

- Try to find the string pattern  $(Ø|a|d(a|e|ia|if|ε)|e|f|i|ko|mo|ni|ob|par|pr(ae|o)|r(a|e)|siū|sou|su(b|p|r)|ter|tran|x(Ø|k)) \mathbf{C}_{m_0}^* \mathbf{V}^* \mathbf{C}_{m_0}^*$  in the string  $\hat{\sigma}$ ;
- If the string pattern above is found, the last position occupied by such a string defines  $\text{ProtRg}(\hat{\sigma})$ ; otherwise, set  $\text{ProtRg}(\hat{\sigma}) = 0$ .  $\square$

**Algorithm 6.16** (French essential root). Let  $\hat{\sigma}$  be the effective spelling of a French word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:

- Break down*  $\hat{\sigma} = \hat{\sigma}_1\hat{\sigma}_2$  *into the concatenation of two strings*  $\hat{\sigma}_1 = \hat{\sigma}^{[\text{ProtRg}(\hat{\sigma})]}$  *(see the notation in Definition 3.1) and*  $\hat{\sigma}_2$ , *where the length of the first string*  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma})$  *is equal to the protected range of*  $\hat{\sigma}$ .
- On*  $\hat{\sigma}_2$ , *perform the following substitutions in a sequel:*
  - Do*  $\sim(ons|t) \rightarrow Ø$ .
  - Do*  $(b(Ø|i)(Ø|l)|nt|~(ā|i|r|ū)X|st|~n) \rightarrow Ø$ .

(3.3)  $\text{Do } (\text{ion}|\text{ke}|\text{log}|\text{s}|t\text{if}|\text{y}) \rightarrow \emptyset$ .<sup>76</sup>

(3.4)  $\text{Do } \sim \mathbf{V}_m^* \rightarrow \emptyset$ .

The result after these four steps of operations is called  $\hat{\sigma}'_2$ .

(3) Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ .

(4) Perform the following replacements on  $\hat{\sigma}_1\hat{\sigma}'_2$ :

$mo(\rho\tau w\rho)\mathbf{X}\sim$	$mow(dr l)\mathbf{X}\sim$	$dekh(err u(\emptyset s ss t))$	$di(\emptyset r s ss t)$	$ekrif$	$kkows$
$mow\rho$	$mowd$	$dekhoyr$	$dict$	$ekris$	$kkowd$
$lu(\emptyset r s ss t)$	$p(ow(\emptyset f rr t) u(\emptyset ys))$	$ri$	$s(a f kh w\rho ys yt) u)$	$v(is(\emptyset s))(\emptyset err)$	
$lis$	$puss$	$rir$	$suss$	$voyr$	
$va(wdr yll)$	$vow(\emptyset dr t yll)$		$xs(e(\emptyset y) i(\emptyset r s ss t) oy(\emptyset r s t))$		
$val$	$vowl$		$xsied$		
$\sim ce(o w(\emptyset r s ss t))$	$\sim enferr$	$\sim na(k y(\emptyset s t))$	$\sim pr(i(\emptyset r s ss t) en)$	$\sim sol(\emptyset f)$	$\sim ti(\emptyset en)$
$cef$	$ensoy$	$nayss$	$prendr$	$sowd$	$tindr$
				$vendr$	$vif$

### 6.2.2 Admissible mutation and approximate clustering

Unlike the three representative Germanic languages and Spanish, we do not need vowel blotting for French, because vowel alternation is not found systematically in French verb conjugations.

In what follows, we will construct a bivariate Boolean-valued function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  on a “simple heredity test function” in Algorithm 6.17, and a set of “admissible suffix mismatch” rules in Algorithm 6.18.

**Algorithm 6.17** (Simple heredity test). *Let  $\hat{\alpha}'$  be the result from doing  $\sim(\emptyset|e)s \rightarrow \emptyset$  on  $\hat{\alpha}$ , and define  $\hat{\beta}'$  similarly. The Boolean-valued function  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if the lowercase form of  $\hat{\alpha}$  contains at least one instance of  $\mathbf{V}^*$  (Definition 6.13) **AND** at least one of the following three conditions holds:<sup>77</sup>*

- (i)  $\hat{\alpha}' = \hat{\beta}'$ ;
- (ii) After doing  $\sim(d|f|fr|t) \rightarrow \emptyset$  on  $\hat{\beta}$ , we obtain  $\hat{\alpha}$ ;
- (iii)  $\ell(\hat{\beta}) > \ell(\hat{\alpha}) \geq \frac{\ell(\hat{\beta})}{2}$  **AND**  $\hat{\alpha} = \hat{\beta}^{[\ell(\hat{\alpha})]} \text{ AND } \hat{\beta}^{[\ell(\hat{\alpha}')+1]} = \mathbf{V}^*$ . (See Definition 3.1 for the notations  $\hat{\beta}^{[n]}$  and  $\hat{\beta}^{[n]}$ .)

In what follows, we define  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{NW}^*(\hat{\alpha}, \hat{\beta})$  and  $\text{SuffixSW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{RootSW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{SW}^*(\hat{\alpha}, \hat{\beta})$  as what is done in the Danish case (Algorithm 5.7), which was also applied to two other Germanic languages treated in §5.

**Algorithm 6.18** (Admissible suffix mismatch). *For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function*

$$\text{AdmSM}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta}))$$

returns **TRUE** if

$$\text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = [(\emptyset|\mathbf{V}^*\mathbf{X}r|(m|r|s|t|y)\mathbf{X}|t\text{if}\mathbf{X}), (\emptyset|\mathbf{V}^*\mathbf{X}r|(m|r|s|t|y)\mathbf{X}|t\text{if}\mathbf{X})]$$

**AND**  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \emptyset$  **AND**  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$  contains at least one instance of  $\mathbf{V}^*$ .

Similarly, one can evaluate another Boolean-valued function

$$\text{AdmSM}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta}))$$

by trading all the occurrences of NW in the statements above with SW.

**Algorithm 6.19** (Heredity test function). *The structure of the French heredity test function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  is identical to the German version (Algorithm 8.1.2), except that the functions  $\text{SimpHrdTest}$ ,  $\text{RootNW}$ ,  $\text{SuffixNW}$ ,  $\text{NW}^*$ ,  $\text{RootSW}$ ,  $\text{SuffixSW}$ ,  $\text{SW}^*$  must follow the French rules stated above.*

**Algorithm 6.20** (Approximate clustering of French words). *The algorithm is essentially the same as Algorithm 5.10, except that French rules (instead of Danish rules) apply to all the tags (effective spelling, essential root etc.), and vowel blotting is not used.*

Concretely speaking, the approximate clustering of a list of French words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  is completed in two stages:

<sup>76</sup>We note that longer matches take priority over shorter matches: if *log* is found, then delete these three letters altogether, instead of just a single letter *g*.

<sup>77</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical OR.

- (1) We sort the list  $\{(\hat{\alpha}_1, \text{EffSpell}(\hat{\alpha}_1)), \dots, (\hat{\alpha}_N, \text{EffSpell}(\hat{\alpha}_N))\}$  alphabetically according to the second component. If two consecutive neighbors in the alphabetized list  $\{(\hat{\alpha}_{(1)}, \text{EffSpell}(\hat{\alpha}_{(1)})), \dots, (\hat{\alpha}_{(N)}, \text{EffSpell}(\hat{\alpha}_{(N)}))\}$  satisfy  $\text{HrdTest}(\text{EffSpell}(\hat{\alpha}_{(n)}), \text{EffSpell}(\hat{\alpha}_{(n+1)})) = \text{FALSE}$ , where  $n \in \mathbb{Z} \cap [1, N]$ , then add a demarcation line between these two entries. In this way, the alphabetized list is divided into separate groups of words tagged with their effective spellings:  $\{g_1 = \{(\hat{\alpha}_{(1,1)}, \text{EffSpell}(\hat{\alpha}_{(1,1)})), \dots, (\hat{\alpha}_{(1,n_1)}, *)\}, \dots, g_M = \{(\hat{\alpha}_{(M,1)}, *), \dots, (\hat{\alpha}_{(M,n_M)}, *)\}\}$ ,<sup>78</sup> where each sublist also preserves alphabetic order. In particular, we have  $\hat{\alpha}_{(1,1)} = \hat{\alpha}_{(1)}$  and  $\hat{\alpha}_{(M,n_M)} = \hat{\alpha}_{(N)}$ .
  - (2) For each group of words  $g_m = \{(\hat{\alpha}_{(m,1)}, *), \dots, (\hat{\alpha}_{(m,n_m)}, *)\}$  where  $m \in \mathbb{Z} \cap [1, M]$ , we augment it into a tagged entry  $G_m = (g_m, \text{EffSpell}(\hat{\alpha}_{(m,1)}), \text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})))$ . The list  $\{G_1, \dots, G_M\}$  is sorted alphabetically, with respect to  $\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))$  (with higher priority) and  $\text{EffSpell}(\hat{\alpha}_{(m,1)})$  (with lower priority). If two consecutive neighbors in the alphabetized list  $\{G_{(1)} = (g_{(1)}, \hat{\gamma}'_{(1)}, \hat{\gamma}''_{(1)}), \dots, G_{(M)} = (g_{(M)}, \hat{\gamma}'_{(M)}, \hat{\gamma}''_{(M)})\}$  satisfy

$$\text{HrdTest}(\hat{\gamma}_{(m)}'', \hat{\gamma}_{(m+1)}'') = \text{FALSE}$$

AND

**HrdTest**( $\hat{\gamma}'_{(m)}, \hat{\gamma}'_{(m+1)}$ ) = FALSE

where  $m \in \mathbb{Z} \cap [1, M]$ , then add a demarcation line between these two entries. In this manner, the alphabetized list  $\{G_{(1)}, \dots, G_{(M)}\}$  is divided into separate groups of tagged clusters  $\{\Gamma_1 = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, \Gamma_K = \{G_{(K,1)}, \dots, G_{(K,m_K)}\}\}$ . Finally, the output list of word clusters  $\{\check{\Gamma}_1, \dots, \check{\Gamma}_K\}$  is constructed by discarding all the tags (effective spellings, essential roots, vowel blotted forms) from each  $\Gamma_k$ , where  $k \in \mathbb{Z} \cap [1, K]$ .

*Example 6.20.1.* French verbs are divided into three conjugation groups. The first two groups are “regular”, while the verbs in the third group are all irregular. Unlike Spanish regular verbs, the root of a verb in the first two conjugations in French sometimes may undergo cosmetic changes for phonological reasons — their counterparts in Spanish would be classified as “irregular”, but we adhere to the traditional French grammar (as in [https://en.wiktionary.org/wiki/Appendix:French\\_verbs](https://en.wiktionary.org/wiki/Appendix:French_verbs)) below, and include them as exceptional cases of the first two “regular” conjugations in French.

We test our algorithm on the following regular French verbs belonging to the first and second conjugation groups:

*aima, aimai, aimaient, aimais, aimait, aimâmes, aimant, aimas, aimasse, aimassent, aimasses, aimassiez, aimassions, aimât, aimâtes, aime, aimé, aiment, aimer, aimera, aimeraï, aimeraient, aimerais, aimerait, aimeras, aimèrent, aimerez, aimeriez, aimerions, aimerons, aimeront, aimes, aimez, aimiez, aimions, aimons — “love”;*

*appela, appelai, appelaient, appelais, appelait, appelâmes, appelant, appelas, appelassee, appelassent, appelasses, appelassiez, appelassions, appelât, appelâtes, appelé, appeler, appellèrent, appelez, appeliez, appelions, appelle, appellent, appellera, appelleraï, appelleraient, appellerais, appelleraït, appelleras, appelleriez, appelleriez, appellerions, appellerons, appelleront, appelles, appelons — “call”;*

*créa, créai, créaient, créais, créait, créâmes, créant, créas, créasse, créassent, créasses, créassiez, créassions, créât, créâtes, crée, créé, créent, créer, créera, crérai, créraient, crérais, crérait, créeras, créèrent, créerez, créeriez, créerions, créerons, créeront, crées, créez, créiez, créions, créons — “create”;*

*fini, finîmes, finir, finira, finirai, finiraient, finirais, finirait, finiras, finirent, finirez, finiriez, finirions, finirons, finiront, finis, finissaient, finissais, finissait, finissant, finisse, finissent, finissez, finissiez, finissions, finissons, finit, finît, finîtes* — “finish”:

*hai, haïmes, haïr, haïra, haïrai, haïraient, haïrais, haïrait, haïras, haïrent, haïrez, haïriez, haïrions, haïrons, haïront, hais, haïs, haïssaien, haïssais, haïssait, haïssant, haïsse, haïssent, haïsses, haïsez, haïssiez, haïssions, haïssons, hait, haït, haïtes — “hate”.*

*mange, mangé, mangea, mangeai, mangeaient, mangeais, mangeait, mangeâmes, mangeant, mangeas, mangeasse, mangeassent, mangeasses, mangeassiez, mangeassions, mangeât, mangeâtes, mangent, mangeons, manger, mangera, mangeraï, mangeraient, mangerais, mangerait, mangeras, mangèrent, mangerez, mangeriez, mangerions, mangerons, mangeront, manges, mangez, mangiez, mangions* — “eat”;

*noie, noient, noiera, noierai, noieraient, noieraies, noierait, noieras, noierez, noieriez, noierions, noierons, noieront, noies, noya, noyai, noyaient, noyais, noyait, noyâmes, noyant, noyas, noyasse, noyassent, noyasses, noyassez, novassions, novât, novâtes, nové, nover, novèrent, novez, noviez, novions, novons* — “drown”;

<sup>78</sup> Here, we use an asterisk (\*) to abbreviate components that are clear from context.

*paie, paient, paiera, paierai, paieraient, paierais, paierait, paieras, paieriez, paierions, paierons, paieront, paies, paya, payai, payaient, payais, payait, payâmes, payant, payas, payasse, payassent, payasses, payassez, payassions, payât, payâtes, paye, payé, payent, payer, payera, payeraient, payerais, payerait, payeras, payèrent, payerez, payeriez, payerions, payeron, payes, payez, payiez, payions, payons — “pay”;*

*plaça, plaçai, plaçaient, plaçais, plaçait, plaçâmes, plaçant, plaças, plaçasse, plaçassent, plaçasses, plaçassiez, plaçassions, plaçât, plaçâtes, place, placé, placent, placer, placera, placeraient, placerais, placerait, placeras, placèrent, placerez, placeriez, placerions, placerons, placeront, places, placez, placiez, placion, placons — “place”.*

and obtain the following results:

{aima, aimai, aimaien, aimais, aimait, aimâmes, aimant, aimas, aimasse, aimassent, aimasses, aimassiez, aimassions, aimât, aimâtes, aime, aimé, aiment, aimer, aimera, aimeraient, aimerais, aimerait, aimeras, aimèrent, aimerez, aimeriez, aimerions, aimeront, aimes, aimez, aimiez, aimions, aimons},

{appela, appelai, appelaient, appelaïs, appelaït, appelaïmes, appelaient, appelas, appelaïsse, appelaïssent, appelaïsses, appelaïssiez, appelaïssions, appelaït, appelaïtes, appelaïé, appelaïer, appelaïèrent, appelaïez, appelaïiez, appelaïons, appelle, appellen, appellerai, appelleraien, appelleraias, appelleraiat, appelleraias, appelleraiiez, appelleraiions, appelleraior, appelleraiel, appelleraiel},

{créa, créai, créaien, créais, créait, créâmes, créant, créas, créasse, créassent, créasses, créassiez, créassions, créât, créâtes, crée, créé, créent, créer, créera, créeraient, créerais, créerait, créeras, créèrent, créerez, créeriez, créerions, créeront, crées, créez, créiez, créions, créons},

{fini, finîmes, finir, finirai, finiraient, finirais, finirait, finiras, finirent, finirez, finiriez, finirions, finirons, finiront, finis, finissaient, finissais, finissait, finissant, finisse, finissent, finisses, finissez, finissiez, finissions, finissons, finit, finît, finîtes},

{hai, haïmes, haïr, haïra, haïrai, haïraient, haïrais, haïrait, haïras, haïrent, haïrez, haïriez, haïrions, haïrons, haïront, hais, haïs, haïssaien, haïssais, haïssait, haïssant, haïsse, haïssent, haïsses, haïsez, haïssiez, haïssions, haïssons, hait, haït, haïtes},

{mange, mangé, mangea, mangeai, mangeaient, mangeais, mangeait, mangeâmes, mangeant, mangeas, mangeasse, mangeassent, mangeasses, mangeassiez, mangeassions, mangeât, mangeâtes, mangent, mangeons, mangera, mangerai, mangeraient, mangerais, mangerait, mangeras, mangèrent, mangerez, mangeriez, mangerions, mangeron, mangeron, manges, mangez, mangiez, mangions},

{noie, noient, noiera, noierai, noieraient, noierais, noierait, noieras, noierez, noieriez, noierions, noierons, noieront, notes, noya, noyai, noyaient, noyais, noyait, noyâmes, noyant, noyas, noyasse, noyassent, noyasses, noyassiez, noyassions, noyât, noyâtes, noyé, noyer, noyèrent, noyez, noyiez, noyions, noyons},

{paie, paient, paiera, paierai, paieraient, paierais, paierait, paieras, paieriez, paierions, paierons, paieront, paies, paya, payai, payaient, payais, payait, payâmes, payant, payas, payasse, payassent, payasses, payassez, payassions, payât, payâtes, paye, payé, payent, payer, payera, payeraient, payerais, payerait, payeras, payèrent, payerez, payeriez, payerions, payeron, payes, payez, payiez, payions, payons},

{plaça, plaçai, plaçaient, plaçais, plaçait, plaçâmes, plaçant, plaças, plaçasse, plaçassent, plaçasses, plaçassiez, plaçassions, plaçât, plaçâtes, place, placé, placent, placer, placera, placeraient, placerais, placerait, placeras, placèrent, placerez, placeriez, placerions, placerons, placeront, places, placez, placiez, placion, placons, placons}.

*Example 6.20.2.* Unlike Spanish, French verb conjugations do not exhibit systematic vowel alternations. In fact, the vowels in almost all French verb stems are invariant. There are a few dozen commonly used French verbs that are irregular in their own unique ways. They are treated by dedicated exceptions to rules in our clustering algorithm. Similar to the sampling of Spanish verbs, we ignore cases where the French verb root is different by at most a diacritic mark on the letter e (as both é and è will be converted to e in the effective spelling) in conjugated forms.

We use the following representative irregular verbs in French as our input:

accrois, accroissaient, accroissais, accroissait, accroissant, accroisse, accroissent, accroisses, accroissez, accroissiez, accroissions, accroissons, accroît, accroîtra, accroîtrai, accroîtraien, accroîtrais, accroîtrait, accroîtras, accroître, accroîtrez, accroîtriez, accroîtrions, accroîtrons, accroîtront, accrû, accrûmes, accrûrent, accrûs, accrûsse, accrûssent, accrûsses, accrûssiez, accrûssions, accrût, accrûtes — “increase”;

*acquéraient, acquérais, acquérait, acquérant, acquérez, acquéries, acquérons, acquéris, acquerra, acquerrai, acquerraient, acquerrais, acquerrait, acquerras, acquerez, acquerriez, acquerrions, acquerrons, acquerront, acquière, acquièrent, acquières, acquiers, acquiert, acquîmes, acquirent, acquis, acquise, acquissent, acquisses, acquissiez, acquissions, acquit, acquît, acquîtes — “acquire”;*

*assaillaient, assaillais, assaillait, assaillant, assaille, assaillé, assaillent, assailles, assailez, assailliez, assaillîmes, assaillions, assaillir, assaillira, assaillrai, assaillraient, assaillirais, assaillirait, assailliras, assaillirent, assaillirez, assailliriez, assaillirions, assaillirons, assailliront, assailis, assaillisse, assaillissent, assaillisses, assaillissiez, assaillissions, assaillit, assaillîtes, assaillons — “assail”;*

*asseoir, asseyaien, asseyais, asseyait, asseyant, asseye, asseyent, asseyes, asseyez, asseyiez, asseyions, assied, assieds, assiéra, assiérai, assiéraient, assiérais, assiérait, assiéras, assiérez, assiéries, assiérons, assieront, assîmes, assirent, assis, assisse, assisent, assises, assissiez, assissions, assit, assít, assîtes, assoie, assoient, assoies, assoira, assoirai, assoiraient, assoirais, assoirait, assoiras, assoirez, assoiriez, assoirions, assoirons, assoiront, assois, assoit, assoyaient, assoyais, assoyait, assoyant, assoyez, assoyiez, assoyions, assoyons — “sit”;*

*bat, bats, battaient, battais, battait, battant, batte, battent, bates, battez, battiez, battîmes, battions, battirent, battis, battisse, battissent, battisses, battissiez, battissions, battit, battît, battîtes, battons, battra, batrai, batraient, battrais, battrait, battras, battre, battrez, battriez, battrians, battrons, battron, battu — “beat”;*

*bouillaient, bouillais, bouillait, bouillant, bouille, bouillent, bouilles, bouillez, bouilli, bouilliez, bouillîmes, bouillions, bouillira, bouillirai, bouillraient, bouillirais, bouillirait, bouilliras, bouillirent, bouillirez, bouilliriez, bouillirions, bouillirons, bouilliront, bouillis, bouillisse, bouillissent, bouillisses, bouillissiez, bouillissions, bouillit, bouillît, bouillîtes, bouillons, bous, bout — “boil”;*

*braie, braient, braies, braira, brairai, brairaient, brairais, brairait, brairas, braire, brairez, brairiez, brairions, brairons, brairont, brais, brait, braya, brayai, brayaient, brayas, brayait, brayâmes, brayant, brayas, brayasse, brayassent, brayasses, brayassiez, brayassions, brayât, brayâtes, brayèrent, bravez, brayiez, brayions, brayons — “bray”;*

*conclu, concluaient, concluas, concluait, conluant, concluent, conclue, conclues, concluez, concluez, concluions, conclûmes, concluons, conclura, conclurai, concluraien, conclurais, conclurait, concluras, conclude, concludeut, conclurez, concluriez, conclurions, conclurons, concluron, conclus, conclusse, conclusent, conclusses, conclusiez, conclussons, conclut, conclût, conclûtes — “conclude”;*

*coud, coudra, coudrai, coudraient, coudrais, coudrait, coudras, coudre, coudrez, coudriez, coudrions, coudrons, coudront, couds, cousaient, cousaïs, cousat, cousat, cose, cousent, courses, cosez, cousiez, cousîmes, cousions, cousirent, cousis, cousisse, cousissent, cousisses, cousissiez, cousissions, cousít, cousîtes, cousons, cousu — “sew”;*

*couraient, courais, courait, courant, courre, courent, courres, couriez, couriez, courions, courir, courons, courra, courrai, courraient, courrais, courrait, courras, courrez, courriez, courrions, courrons, courront, cours, court, couru, courûmes, coururent, courus, courusse, courusent, courusses, courussiez, courussions, courut, courût, courûtes — “run”;*

*craignaient, craignais, craignait, craignant, craigne, craignent, craignes, craignez, craigniez, craignîmes, craignions, craignirent, craignis, craignisse, craignissent, craignisses, craignissiez, craignissions, craignit, craignît, craignîtes, craignons, craindra, craindrai, craindraient, craindrais, craindrait, craindras, craindre, craindrez, craindriez, craindrions, craindrons, craindront, crains, crant — “fear”;*

*croie, croient, croies, croira, croirai, croiraient, croirais, croirait, croiras, croire, croirez, croiriez, croirions, croirons, croiront, crois, croit, croyaient, croyais, croyait, croyant, croyez, croyiez, croyions, croyons, cru, crûmes, crurent, crus, crusse, crusent, crusses, crussiez, crussions, crut, crût, crûtes — “believe”;*

*cueillaient, cueillais, cueillait, cueillant, cueille, cueillent, cueillera, cueillerai, cueilleraien, cueillerais, cueillerait, cueilleras, cueillerez, cueilleriez, cueillerions, cueillerons, cueilleront, cueilles, cueillez, cueilli, cueilliez, cueillîmes, cueillions, cueillir, cueillirent, cueillis, cueillisse, cueillissent, cueillisses, cueillissiez, cueillissions, cueillit, cueillît, cueillîtes, cueillons — “gather”;*

décherra, décherrai, décherraien, décherrais, décherrait, décherras, décherrez, décherriez, décherriions, décherrons, décherront, déchet, déchoie, déchoient, déchoies, déchoir, déchoira, déchoirai, déchoiraient, déchoirais, déchoirait, déchoiras, déchoirez, déchoiriez, déchoirions, déchoirons, déchoiront, déchois, déchoit, déchoyant, déchoyez, déchoyez, déchoyions, déchoyons, déchu, déchumes, déchurent, déchus, déchusse, déchussent, déchusses, déchussiez, déchussions, déchut, déchût, déchûtes — “wane”;

*dîmes, dira, dirai, diraient, dirais, dirait, diras, dire, dirent, direz, diriez, dirions, dirons, diront, dis, disaient, disais, disait, disant, dise, disent, dises, disiez, disions, disons, disse, dissent, disses, dissiez, dissions, dit, dît, dites, dîtes* — “say”;

*dissolûmes, dissolurent, dissolus, dissolut, dissolûtes, dissolvaient, dissolvais, dissolvait, dissolvant, dissolve, dissolvent, dissolves, dissolvez, dissolviez, dissolvions, dissolvons, dissoudra, dissoudrai, dissoudraient, dissoudrais, dissoudrait, dissoudras, dissoudre, dissoudrez, dissoudriez, dissoudrions, dissoudrons, dissoudront, dissous, dissout — “dissolve”;*

*dormaient, dormais, dormait, dormant, dorme, dorment, dormes, dormez, dormi, dormiez, dormîmes, dormions, dormir, dormira, dormirai, dormiraient, dormirais, dormirait, dormiras, dormirent, dormirez, dormiriez, dormirions, dormirons, dormiront, dormis, dormisse, dormissent, dormisses, dormissiez, dormissions, dormit, dormît, dormîtes, dormons, dors, dort — “sleep”;*

*éclora, éclorai, écloraient, éclorais, éclorait, écloras, éclore, éclorez, écloriez, éclorions, éclorons, écloront, éclos, éclosant, éclose, éclosent, écloses, éclosez, éclosiez, éclosions, éclosons, éclôt — “hatch”;*

*écrira, écrirai, écriraient, écrirais, écrirait, écriras, écrire, écrirez, écririez, écririons, écrirons, écriront, écris, écrit, écrivaient, écrivais, écrivait, écrivant, écrive, écrivent, écrives, écrivez, écriviez, écrivîmes, écrivions, écrivirent, écrivis, écrivisse, écrivissent, écrivisses, écrivissiez, écrivissions, écrivit, écrivît, écrivîtes, écrivons — “write”;*

*enverra, enverrai, enverraient, enverrais, enverrait, enverras, enverrez, enverriez, enverrions, enverrons, enverront, envoie, envoient, envies, envoyas, envoyai, envoyaien, envoyais, envoyait, envoyâmes, envoyant, envoyas, envoyasse, envoyassent, envoyasses, envoyassiez, envoyassions, envoyât, envoyâtes, envoyé, envoyer, envoyèrent, envoyez, envoviez, envovions, envovons — “send”;*

*fui, fuiε, fuient, fuiεs, fuiμes, fuir, fira, firaι, firaient, firaιs, fuirait, firas, furent, firez, firiez, firions, furons, firon, fuis, fuisse, fuisse, fuisse, fuisse, fuisse, fuit, fuit, futes, fuyaient, fuyaιs, fuyait, fuyant, fuez, fuyiez, fuyions, fuyons — “escape”;*

*incluaient, incluais, incluait, incluant, inclue, incluent, inclus, incluez, incliez, incluions, inclûmes, incluons, inclura, inclurai, incluraient, inclurais, inclurait, incluras, inclure, inclurent, inclurez, incluriez, inclurions, inclurons, incluront, inclus, inclusse, inclussent, inclusses, inclussiez, inclussions, inclut, inclût, inclûtes — “include” ;*

*joignaient, joignais, joignait, joignant, joigne, joignent, joignes, joignez, joigniez, joignîmes, joignions, joignirent, joignis, joignisse, joignissent, joignisses, joignissiez, joignissions, joignit, joignît, joignîtes, joignons, joindra, joindrai, joindraient, joindrais, joindrait, joindras, joindre, joindrez, joindriez, joindrions, joindrons, joindront, joins. joint — “join”;*

*lira, lirai, liraient, lirais, lirait, liras, lire, lirez, liriez, lirions, lirons, liront, lis, lisaien, lisais, lisait, lisant, lise, lisent, lises, lisez, lisiez, lisions, lisons, lit, lu, lûmes, lurent, lus, lusse, lussent, lusses, lussiez, lussions, lut, lût, lûtes* — “read”;

*meure, meurent, meures, meurs, meurt, mort, mouraient, mourais, mourait, mourant, mourez, mouriez, mourions, mourir, mourons, mourra, mourrai, mourraient, mourrais, mourrait, mourras, mourrez, mourriez, mourrions, mourrons, mourront, mourûmes, moururent, mourus, mourusse, mourussent, mourusses, mourussiez, mourusions, mourut, mourût, mourûtes* — “die”:

*meus, meut, meuve, meuvent, meuves, mouvaient, mouvais, mouvait, mouvant, mouvez, mouviez, mouvions, mouvoir, mouvons, mouvra, mouvrai, mouvraient, mouvrais, mouvrait, mouvras, mouvrez, mouvriez, mouvrions, mouvrons, mouvront, mû, mûmes, murent, mus, musse, mussent, musses, mussiez, mussions, mut, mût, mûtes*  
— “move”.

*moud, moudra, moudrai, moudraient, moudrais, moudrait, moudras, moudre, moudrez, moudriez, moudrions, moudrons, moudront, mouds, moulaien, moulais, moulait, moulant, moule, moulen, moules, moulez, mouliez, moulions, moulons, moulu, moulûmes, moulurent, moulus, moulusse, moulussent, moulusses, moulussiez, moulussions, moulut, moulût, moulûtes* — “grind”;

*nais, naissaient, naissais, naissait, naissant, naisse, naissent, naisses, naissez, naissiez, naissions, naissons, naît, naîtra, naîtrai, naîtraient, naîtrais, naîtrait, naîtras, naître, naîtrez, naîtriez, naîtrions, naîtrons, naîtront, naquîmes, naquirent, naquis, naquisse, naquissent, naquisses, naquissiez, naquissions, naquit, naquît, naquîtes, né* — “be born”;

*nui, nuira, nuirai, nuiraien, nuirais, nuirait, nuiras, nuire, nuirez, nuiriez, nuirions, nuirons, nuiront, nuis, nui-saient, nuisais, nuisait, nuisant, nuise, nuisent, nuses, nuisez, nuisiez, nûsîmes, nûsions, nûsirent, nûsis, nûsisse, nûsissent, nûsisses, nûsissiez, nûsissions, nûosit, nûsîtes, nûsions, nuit* — “spoil”;

*offert, offraient, offrais, offrait, offrant, offre, offrent, offres, offrez, offrîmes, offrions, offrir, offrira, offrîrai, offriraient, offrîrais, offrirait, offriras, offrîrent, offrire, offririez, offririons, offrîrons, offriront, offris, offrisse, offrissent, offrisses, offrissiez, offrissions, offrit, offrît, offrîtes, offrons* — “offer”;

*peignaient, peignais, peignait, peignant, peigne, peignent, peignes, peignez, peigniez, peignîmes, peignions, peignirent, peignis, peignisse, peignissent, peignisses, peignissiez, peignissions, peignit, peignît, peignîtes, peignons, peindra, peindrai, peindraient, peindrais, peindrait, peindras, peindre, peindrez, peindriez, peindrions, peindrons, peindront, peins, peint* — “paint”;

*plaira, plairai, plairaient, plairais, plairait, plairas, plaire, plairez, plairiez, plairions, plairons, plairont, plais, plaisaient, plaisais, plaisait, plaisant, plaise, plaisent, plaises, plaisez, plaisiez, plaisions, plait, plaît, plu, plûmes, plurent, plus, plusse, plussent, plusses, plussiez, plussions, plut, plût, plûtes* — “please”;

*prenaient, prenais, prenait, prenant, prend, prendra, prendrai, prendraient, prendrais, prendrait, prendras, prendre, prendrez, prendriez, prendrions, prendrons, prendront, prends, prenez, preniez, prenions, prenne, prennent, prennes, prenons, prîmes, prirent, pris, prissee, prissent, prisses, prissiez, prissions, prit, prît, prîtes* — “take”;

*recevaient, recevais, recevait, recevant, receivez, receviez, recevions, recevoir, recevons, recevra, recevrai, recevraient, recevrais, recevrait, recevras, recevez, recevriez, recevrions, recevrons, recevront, reçois, reçoit, reçoive, reçoivent, reçoives, reçu, reçûmes, reçurent, reçus, reçusse, reçussent, reçusses, reçussiez, reçussions, reçut, reçût, reçûtes* — “receive”;

*riaient, riais, riait, riant, rie, rient, ries, riez, riiez, riions, rîmes, rions, rira, rirai, riraient, rirais, rirait, riras, rire, rirent, rirez, ririez, ririons, riront, ris, risse, rissent, risses, rissiez, rissions, rit, rît, rîtes* — “laugh”;

*sachant, sache, sachent, saches, sachez, sachiez, sachions, sachons, sais, sait, saura, saurai, sauraient, saurais, saurait, sauras, saurez, sauriez, saurions, saurons, sauront, savaien, savais, savent, savez, saviez, savions, savoir, savons, su, sûmes, surent, sus, susse, sussent, susses, sussiez, sussions, sut, sût, sûtes* — “know”;

*suffi, suffîmes, suffîra, suffîrai, suffîraient, suffîrais, suffîrait, suffîras, suffire, suffirent, suffîrez, suffîriez, suffîrions, suffîrons, suffîront, suffîs, suffisaient, suffisais, suffisait, suffisant, suffise, suffisent, suffises, suffisez, suffîiez, suffisions, suffisons, suffîsse, suffisent, suffisses, suffissiez, suffissions, suffit, suffît, suffîtes* — “suffice”;

*traie, traient, traies, traïra, traïrai, traïraient, traïrais, traïrait, traïras, traire, traïrez, traïriez, traïrions, traïrons, traïront, traïs, trait, trayâ, trayai, trayaien, trayais, trayait, trayâmes, trayant, trayas, trayasse, trayassent, trayasses, trayassiez, trayassions, trayât, trayâtes, trayèrent, trayez, trayiez, trayions, trayons* — “milk”;

*vaille, vaillent, vailles, valaient, valais, valait, valant, valent, valez, valiez, valions, valoir, valons, valu, valûmes, valurent, valus, valusse, valusent, valusses, valussiez, valussions, valut, valût, valûtes, vaudra, vaudrai, vaudraient, vaudrais, vaudrait, vaudras, vaudrez, vaudriez, vaudrions, vaudrons, vaudront, vaut, vaux* — “earn”;

*vainc, vaincra, vaincrai, vaincraient, vaincrais, vaincrait, vaincras, vaincre, vaincrez, vaincriez, vaincrions, vaincrons, vaincront, vaincs, vaincu, vainquaient, vainquais, vainquait, vainquant, vainque, vainquent, vainques, vainquez, vainquiez, vainquîmes, vainquions, vainquient, vainquis, vainquise, vainquissent, vainquisses, vainquissiez, vainquissions, vainquit, vainquît, vainquîtes, vainquons* — “win”;

*vécu, vécûmes, vécurent, vécus, vécusse, vécussent, vécusses, vécussiez, vécussions, vécut, vécût, vécûtes, vis, vit, vivaien, vivais, vivait, vivant, vive, vivent, vives, vivez, viviez, vivions, vivons, vivra, vivrai, vivraient, vivrais, vivrait, vivras, vivre, vivrez, vivriez, vivrions, vivrons, vivront* — “live”;

vêt, vêtaient, vêtais, vêtait, vêtant, vête, vêtent, vêtes, vêtez, vêtiez, vêtîmes, vêtions, vêtir, vêtira, vêtirai, vêtiraient, vêtirais, vêtirait, vêtiras, vêtirent, vêtirez, vêtiriez, vêtirions, vêtirons, vêtiront, vêtis, vêtisse, vêtissent, vêtisses, vêtissiez, vêtissions, vêtit, vêtû, vêtîtes, vêtions, vêts, vêtu — “dress”;

veuille, veuillent, veuilles, veulent, veut, veux, voudra, voudrai, voudraient, voudrais, voudrait, voudras, voudrez, voudriez, voudrions, voudrons, voudront, voulaient, voulais, voulait, voulant, voulez, vouliez, voulions, vouloir, voulons, voulu, voulûmes, voulurent, voulus, voulusse, voulussent, voulusses, voulussiez, voulussions, voulut, voulût, voulûtes — “wish”.

and obtain the following output:

{accrois, accroissaient, accroissais, accroissait, accroissant, accroisse, accroissent, accrois, accroissez, accroissiez, accroissions, accroissons, accroît, accroîtra, accroîtrai, accroîtraien, accroîtrais, accroîtrait, accroîtras, accroître, accroîtrez, accroîtriez, accroîtrions, accroîtrons, accroîtront, accrû, accrûmes, accrûrent, accrûs, accrûsse, accrûsset, accrûsses, accrûssiez, accrûssions, accrût, accrûtes},

{acquéraient, acquérais, acquérait, acquérant, acquérez, acquéries, acquérions, acquérir, acquérons, acquerra, acquerrai, acquerraient, acquerraits, acquerrait, acquerras, acquerez, acquerriez, acquerriens, acquerrons, acquerront, acquière, acquières, acquiers, acquiert, acquîmes, acquirent, acquis, acquisse, acquissent, acquisses, acquissiez, acquisitions, acquit, acquît, acquîtes},

{assaillaient, assaillais, assaillait, assaillant, assaille, assaillé, assaillent, assailles, assaillez, assailliez, assaillîmes, assaillions, assaillir, assaillira, assaillirai, assailliraient, assaillirais, assaillirait, assailliras, assaillirent, assaillirez, assailliriez, assaillirions, assaillirons, assailliront, assaillis, assaillisse, assaillissent, assaillisses, assaillissiez, assaillissions, assaillit, assaillîtes, assaillons},

{asseoir, asseyaien, asseyais, asseyait, asseyant, asseye, asseyent, asseyes, asseyez, asseyions, asseyons, assied, assieds, assiéra, assiérai, assiéraient, assiérais, assiérait, assiéras, assiérez, assiériez, assiériens, assiérons, assiéront, assîmes, assirent, assis, assisse, assisent, assisses, assissiez, assissions, assit, assít, assîtes, assoie, assoient, assoies, assoira, assoirai, assoiraient, assoirais, assoirait, assoiras, assoirez, assoiriez, assoirions, assoirons, assoiront, assois, assoit, assoyaient, assoyais, assoyait, assoyant, assoyez, assoyiez, assoyions, assoyons},

{bat, bats, battaient, battais, battait, battant, batte, battent, battez, battiez, battîmes, battions, battirent, battis, battise, battissent, battises, battissiez, battissions, battit, battît, battîtes, battons, battra, battraï, battraient, battrais, battrait, battras, battre, battrez, battriez, battrians, battrons, battron, battu},

{bouillaient, bouillais, bouillait, bouillant, bouille, bouillent, bouilles, bouillez, bouilli, bouilliez, bouillîmes, bouillions, bouillira, bouillirai, bouilliraient, bouillirais, bouillirait, bouilliras, bouillirent, bouillirez, bouilliriez, bouillirions, bouillirons, bouilliront, bouillis, bouillisse, bouillissent, bouillisses, bouillissiez, bouillissions, bouillit, bouillît, bouillîtes, bouillons, bout, bout},

{braie, braient, braies, braira, brairai, brairaien, brairais, brairait, brairas, braire, brairez, brairiez, brairions, brairons, brairont, brais, brait, braya, brayai, brayaient, brayas, brayait, brayâmes, brayant, brayas, brayasse, brayassent, brayasses, brayassiez, brayassions, brayât, brayâtes, brayèrent, brayez, brayiez, brayions, brayons},

{conclu, concluaient, concluas, concluait, concluant, conclue, conluent, conclues, concluez, concluiez, concluions, conclûmes, concluons, conclura, conclurai, concluraien, conclurais, conclurait, concluras, conclure, conclurent, conclurez, concluriez, conclurions, conclurons, concluron, conclus, conclusse, conclusent, conclusses, conclusiez, conclusions, conclut, conclût, conclûtes},

{coud, coudra, coudrai, coudraient, coudrais, coudrait, coudras, coudre, coudrez, coudriez, coudriens, coudrons, coudront, couds, cousaient, couais, couasit, couant, cose, couent, courses, cosez, couiez, couîmes, couisons, couisirent, couisis, couisse, couissent, couisses, couissiez, couissions, couisit, couîtes, cousons, cousu},

{couraint, courais, courait, courant, courre, courent, courres, couriez, couriez, courions, courir, courons, courra, courrai, courraient, courrais, courrait, courras, courrez, courriez, courriens, courrons, courront, cours, couru, courûmes, coururent, courus, courusse, courussen, courusses, courussiez, courussions, courut, courût, courûtes},

{court},

{craignaient, craignais, craignait, craignant, craigne, craignent, craignes, craignez, craigniez, craignîmes, craignions, craignirent, craignis, craignisse, craignissent, craignisses, craignissiez, craignissions, craignit, craignît, craignîtes, craignîtes, craignons, craindra, craindrai, craindraient, craindrais, craindrait, craindras, craindre, craindrez, craindriez, craindrions, craindrons, craindront, crains, craint},

{croie, croient, croies, croira, croirai, croiraient, croirais, croirait, croiras, croire, croirez, croiriez, croirions, croirons, croiront, crois, croit, croyaient, croyais, croyait, croyant, croyez, croyiez, croyions, croyons, cru, crûmes, crurent, crus, crusse, crussent, crusses, crussiez, crussions, crut, crût, crûtes},

{cueillaient, cueillais, cueillait, cueillant, cueille, cueillent, cueillera, cueillerai, cueilleraien, cueillerait, cueilleras, cueillerez, cueilleriez, cueillerions, cueillerons, cueilleront, cueilles, cueillez, cueilli, cueilliez, cueillîmes, cueillions, cueillir, cueillirent, cueillis, cueillisse, cueillissent, cueillisses, cueillissiez, cueillissions, cueillit, cueillît, cueillîtes, cueillons},

{décherra, décherrai, décherraient, décherrais, décherrait, décherras, décherrez, décherriez, décherriions, décherrons, décherront, déchet, déchoie, déchoient, déchoies, déchoir, déchoira, déchoirai, déchoiraient, déchoirais, déchoirait, déchoiras, déchoirez, déchoiriez, déchoirions, déchoirons, déchoiront, déchois, déchoit, déchoyant, déchoyez, déchoyiez, déchoyions, déchoyons, déchu, déchûmes, déchurent, déchus, déchusse, déchussent, déchusses, déchussiez, déchussions, déchut, déchût, déchûtes},

{dîmes, dira, dirai, diraient, dirais, dirait, diras, dire, dirent, direz, diriez, dirions, dirons, diront, dis, disaient, disais, disait, disant, dise, disent, dises, disiez, disions, disons, disse, dissent, disses, dissiez, dissions, dit, dît, dites, dîtes},

{*dissolûmes, dissolurent, dissolus, dissolut, dissolûtes, dissolvaient, dissolvais, dissolvait, dissolvant, dissolve, dissolvent, dissolves, dissolvez, dissolviez, dissolvions, dissolvons, dissoudra, dissoudrai, dissoudraient, dissoudrais, dissoudrait, dissoudras, dissoudre, dissoudrez, dissoudriez, dissoudrions, dissoudrons, dissoudront, dissous, dissout*},

{dormaient, dormais, dormait, dormant, dorme, dorment, dormes, dormez, dormi, dormiez, dormîmes, dormions, dormir, dormira, dormirai, dormiraient, dormirais, dormirait, dormiras, dormirent, dormirez, dormiriez, dormirions, dormirons, dormiront, dormis, dormisse, dormissent, dormisses, dormissiez, dormissions, dormit, dormît, dormîtes, dormons, dors, dort},

{éclora, éclorai, écloraient, éclorais, éclorait, écloras, éclore, éclorez, écloriez, éclorions, éclorons, écloront, éclos, éclosant, éclose, éclosent, écloses, éclosez, éclosiez, éclosions, éclosons, éclôt},

{écrira, écrirai, écriraient, écrirais, écrirait, écriras, écrire, écrirez, écririez, écririons, écriront, écris, écrit, écrivaient, écrivais, écrivait, écrivant, écrive, écrivent, écrives, écrivez, écriviez, écrivîmes, écrivions, écrivirent, écrivis, écrivisse, écrivissent, écrivisses, écrivissiez, écrivissions, écrivit, écrivît, écrivîtes, écrivons},

{*enverra, enverrai, enverraient, enverrais, enverrait, enverras, enverrez, enverriez, enverrions, enverrons, enverront, envoie, envoient, envoies, envoya, envoyai, envoyaien, envoyais, envoyait, envoyâmes, envoyant, envoyas, envoyasse, envoyassent, envoyasses, envoyassiez, envoyassions, envoyât, envoyâtes, envoyé, envoyer, envoyèrent, envoyez, envoyiez, envoyions, envoyovons*},

{*fui, fuie, fument, fumes, fuir, fura, furaient, furais, fuirait, furas, furent, firez, firiez, fuirions, furons, firont, fuis, fuisse, fissent, fusses, fuissiez, fuissions, fuit, fût, fûtes, fuyaient, fuyais, fuyait, fuyant, fuyez, fuyiez, fuyions, fuvons*}},

{*inclusaient, inclusaais, inclusaait, inclusaant, inclusaie, inclusaient, inclusaas, inclusaiez, inclusaions, inclusaumes, inclusaons, inclusra, inclusrai, inclusraient, inclusraais, inclusraait, inclusras, inclusre, inclusrent, inclusreuz, inclusriez, inclusrions, inclusrons, inclusront, inclus, inclusse, inclussant, inclusses, inclussiez, inclussions, inclusut, inclusit, inclusutes*}},

{*joignaient, joignais, joignait, joignant, joigne, joignent, joignes, joignez, joigniez, joignîmes, joignions, joignirent, joignis, joignisse, joignissent, joignisses, joignissiez, joignissions, joignit, joignît, joignîtes, joignons, joindra, joindrai, joindraient, joindrais, joindrait, joindras, joindre, joindrez, joindriez, joindrions, joindrons, joindront, joins, joint*}},

{*lira, lirai, liraient, lirais, lirait, liras, lire, lirez, liriez, lirions, lirons, liront, lis, lisaien, lisaien, lisais, lisait, lisant, lise, lisent, lises, lisez, lisiez, lisions, lisons, lit, lu, lûmes, lurent, lus, lusse, lussent, lusses, lussiez, lussions, lut, lût, lûtes*}.

{meure, meurent, meures, meurs, meurt, mort, mouraient, mourais, mourait, mourant, mourez, mouriez, mourions, mourir, mourons, mourra, mourrai, mourraient, mourrais, mourrait, mourras, mourrez, mourriez, mourrions, mourrons, mourront, mourûmes, moururent, mourus, mourusse, mourussent, mourusses, mourussiez, mourus-  
sions, mourut, mourût, mourûtes},

{*meus, meut, meuve, meuvent, meuves, mouvaient, mouvais, mouvait, mouvant, mouvez, mouviez, mouvions, mouvoir, mouvons, mouvra, mouvrai, mouvrainent, mouvrais, mouvrait, mouvras, mouvrez, mouvriez, mouvrions, mouvrons, mouvront, mû, mûmes, murent, mus, musse, mussent, musses, mussiez, mussions, mut, mût, mûtes*} ,

{*moud, moudra, moudrai, moudraient, moudrais, moudrait, moudras, moudre, moudrez, moudriez, moudrions, moudrons, moudront, mouds, moulaient, moulais, moulait, moulant, moule, mourent, moules, moulez, mouliez, moulions, moulons, moulu, moulûmes, moulurent, moulus, moulusse, moulussent, moulusses, moulussiez, moulussions, moulut, moulût, moulûtes*}},

{naïs, naissaient, naissais, naissait, naissant, naisse, naissent, naisses, naissez, naissiez, naissions, naissons, naït, naïtra, naïtrai, naïtraien, naïtrais, naïtrait, naïtras, naïtre, naïtrez, naïtriez, naïtrions, naïtrons, naïtront, naïquîmes, naïquirent, naïquis, naïquisse, naïquissent, naïquisses, naïquissiez, naïquissions, naïquit, naïquît, naïquîtes, né}.

{*nui, nuira, nuirai, nuiraien, nuirais, nuirait, nuiras, nuire, nuirez, nuiriez, nuirions, nuirons, nuiront, nuis, nui-saient, nuisais, nuisait, nuisant, naise, nisen, nises, nisez, nisiez, nisimes, nisions, nisirent, nisis, nui-sisse, nisissent, nisisse, nisissez, nisissons, nisit, nisit, nisites, nisons, nuit*},

{offert, offraient, offrais, offrait, offrant, offre, offrent, offres, offrez, offriez, offrimes, offrions, offrir, offrira, offrirai, offriraient, offrirais, offrirait, offriras, offrrent, offrire, offririez, offritions, offrions, offriront, offris, offrisse, offrissent, offrisses, offrissiez, offrissions, offrit, offrit, offrites, offrons}.

{*peignaient, peignais, peignait, peignant, peigne, peignent, peignes, peignez, peigniez, peignîmes, peignions, peignirent, pegnis, pegnisse, pegnissent, pegnisses, pegnissiez, pegnissions, pegnit, pegnît, pegnîtes, pegnons, peindra, peindrai, peindraient, peindrais, peindrait, peindras, peindre, peindrez, peindriez, peindrions, peindrons, peindront, peins, peint*}.

{*plaira, plairai, plairaient, plairais, plairait, plairas, plaire, plairez, plairiez, plairions, plairons, plairont, plais, plaisaient, plaisais, plaisait, plaisant, plaise, plaisent, plaises, plaisez, plaisiez, plaisions, plaisons, plait, plait*}},

*{plu, plûmes, plurent, plus, plusse, plussent, plusses, plussiez, plussions, plut, plût, plûtes}*,

{*prenaient, prenais, prenait, prenant, prend, prendra, prendrai, prendraient, prendrais, prendrait, prendras, prendre, prendrez, prendriez, prendrions, prendrons, prendront, prends, prenez, preniez, prenions, prenne, prennent, prennes, prenons, prîmes, prirent, pris, prissee, prissent, prisses, prissiez, prissions, prit, prît, prîtes*}.

{*recevaient, recevais, recevait, recevant, recevez, receviez, recevions, recevoir, recevons, recevra, recevrai, recevraient, recevrais, recevrait, recevas, recevrez, recevriez, recevrions, recevrons, recevront, reçois, reçoit, reçoive, reçoivent, reçoives, reçu, reçumes, reçurent, reçus, reçusse, reçussent, reçusses, reçussiez, reçussions, reçut, reçût, reçûtes*}},

{riaient, riais, riait, riant, rie, rient, ries, riez, riiez, riions, rîmes, rions, rira, rirai, riraient, rirais, rirait, riras, rire, rirent, rirez, ririez, ririons, rirons, riront, ris, risse, rissent, risses, rissiez, rissions, rit, rît, rîtes},

{*sachant, sache, sachent, saches, sachez, sachiez, sachions, sachons, sais, sait, saura, saurai, sauraient, saurais, saurait, sauras, saurez, sauriez, saurions, saurons, sauront, savaient, savais, savait, savent, savez, saviez, savions, savoir, savons, su, sûmes, surent, sus, susse, sussent, susses, sussiez, sussions, sut, sût, sûtes*}},

{*suffi, suffimes, suffira, suffirai, suffiraient, suffirais, suffirait, suffiras, suffire, suffirent, suffirez, suffiriez, suffirions, suffirons, suffiront, suffis, suffisaient, suffisais, suffisait, suffisant, suffise, suffisent, suffises, suffisez, suffiez, suffisions, suffisons, suffisse, suffissent, suffisses, suffissiez, suffissions, suffit, suffit, suffites*},

{*traie, traient, traies, traira, trairai, traираient, traираis, traирait, trairas, traire, trairez, trairiez, traирions, traироnсs, traирoнт, traис, trait, traya, trayai, trayаient, trayаis, trayait, trayâmes, trayant, trayas, trayasse, trayassent, trayasses, trayassiez, trayassions, trayât, trayâtes, trayèrent, trayez, trayiez, trayions, trayons*},

{vaille, vaillent, vailles, valaient, valais, valait, valant, valent, valez, valiez, valions, valoir, valons, valu, valûmes, valurent, valus, valusse, valussent, valusses, valussiez, valussions, valut, valût, valûtes, vaudra, vaudrai, vaudraient, vaudrais, vaudrait, vaudras, vaudrez, vaudriez, vaudrions, vaudrons, vaudront, vaut, vaux},

{vainc, vaincra, vaincrai, vaincraient, vaincrais, vaincrait, vaincras, vaincre, vaincrez, vaincriez, vaincrions, vaincrons, vaincront, vaincs, vaincu, vainquaient, vainquais, vainquait, vainquant, vainque, vainquent, vainques, vainquez, vainquiez, vainquîmes, vainquions, vainquirent, vainquis, vainquise, vainquissent, vainquisses, vainquissiez, vainquissions, vainquit, vainquîtes, vainquons},

{vêcu, vécûmes, vécurent, vécus, vécusse, vécussent, vécusses, vécussiez, vécussions, vécut, vécût, vécûtes, vi-vaient, vivais, vivait, vivant, vive, vivent, vives, vivez, viviez, vivions, vivons, vivra, vivrai, vivraient, vivrais, vi-vrait, vivras, vivre, vivrez, vivriez, vivrions, vivrons, vivront},

{vêt, vêtaient, vêtais, vêtait, vêtant, vête, vêtent, vêtes, vêtez, vêtîmes, vêtions, vêtir, vêtira, vêtirai, vêtiraient, vêtirais, vêtirait, vêtiras, vêtirent, vêtirez, vêtiriez, vêtirions, vêtiront, vêtis, vêtisse, vêtissent, vêtisses, vêtissiez, vêtissions, vêtit, vêtît, vêtîtes, vêtôns, vêts, vêtu},

{veuille, veuillent, veuilles, veulent, veut, veux, voudra, voudrai, voudraient, voudrais, voudrait, voudras, voudrez, voudriez, voudrions, voudrons, voudront, voulaient, voulais, voulait, voulant, voulez, vouliez, voulions, vouloir, voulons, voulu, voulûmes, voulurent, voulus, voulusse, voulussent, voulusses, voulussiez, voulussions, voulut, voulût, voulûtes},

{vis, vit},

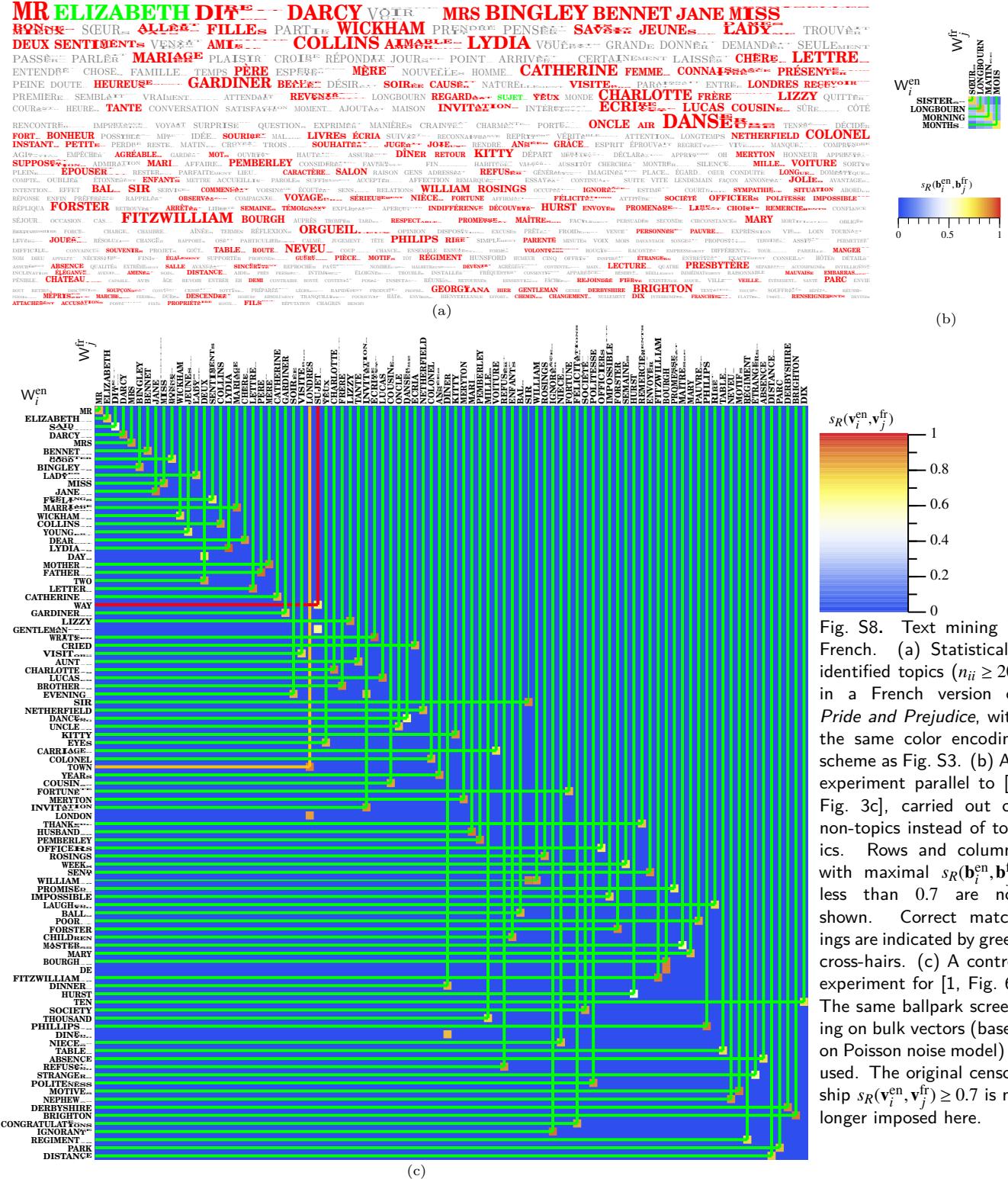
Here, some verb conjugations are not clustered as expected, because we want to accommodate to inherent ambiguities of certain French verb forms. For example, both *vis* and *vit* may be conjugated forms of *vivre* “live” or *voir* “see”, and the infinitive form of *plu* may be either *plaire* “please” or *pleuvoir* “rain”.

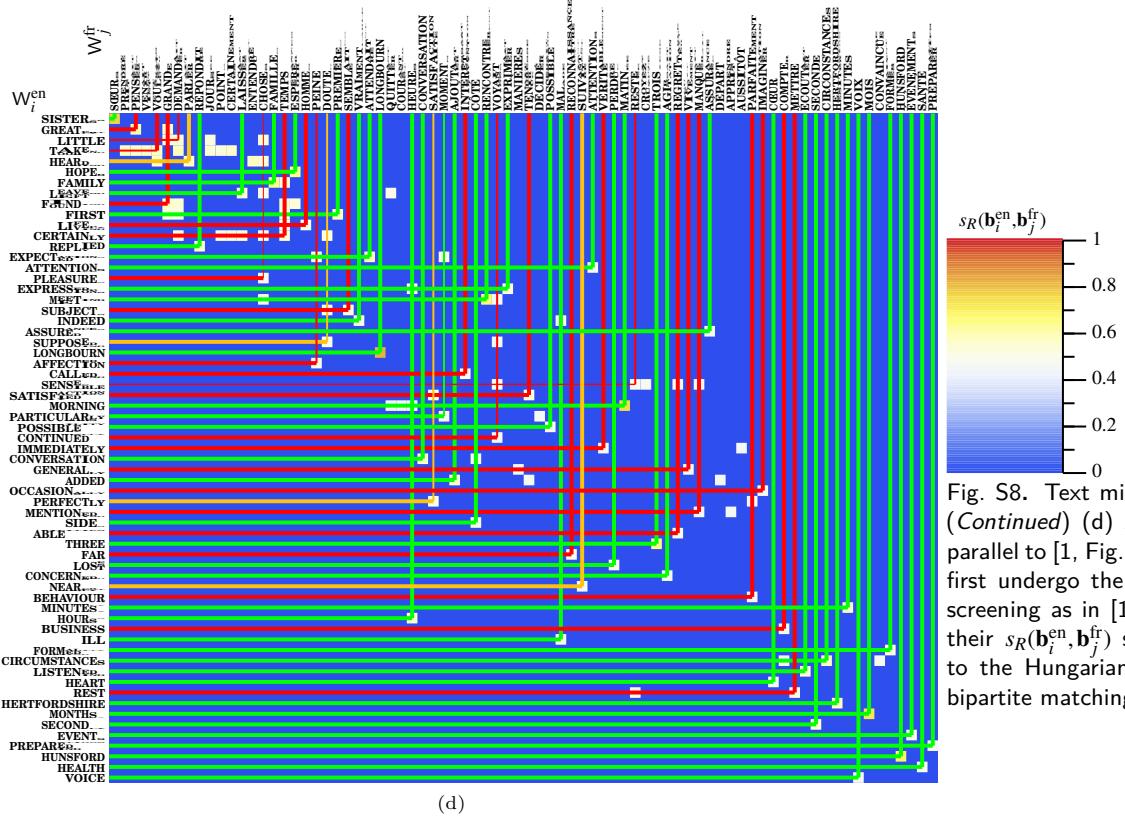
*Example 6.20.3.* In [1, Figs. 3c, 6] and Fig. S8, we further apply the aforementioned word clustering algorithm to topic extraction and machine translation (see Table S1 for text sources).

It should be noted that the performance of machine translation is affected by polysemy of certain commonly occurring French words. For example, *livre* “book, pound” and *livres* “books, pounds” can also be conjugated forms of the verb *livrer* “deliver”. Furthermore, a word in English may correspond to a phrase in French, such as *jeune fille* “girl” (in contrast with *fille* “daughter”). This explains why our algorithm does not identify French translations for such common words like “girl” and “daughter”.

Topicality (non-Poissonian behavior) plays an essential rôle in our algorithms. If a word pattern generates a significantly non-Poissonian trajectory, then it will have distinguishing features in its time structure ([1, Figs. 3c, 6] and Fig. S8c), which allow better semantic resolution than nearly Poissonian trajectories (Fig. S8b and Fig. S8d). The recurrence kinetics of non-topics are very close to single exponential decays, so the vector embedding of recurrence eigenvalues ([1, Figs. 3c, 6] and Fig. S8c) will not help here.

In French, “bat” (a kind of flying mammal) is called *chauve-souris* (literally “bald mouse”), so we consider *bat* an exact match to *chauve* “bald” in Fig. S8b”. (Also note that the row for *bat* in Fig. S8b” indeed contains two hot spots, one for *chauve* and one for *souris*.) In Darwin’s *Origin of Species*, the adjective *fresh* is mostly used in the context of *fresh water*, which is *eau douce* in French. Therefore, we consider *fresh* an exact match to *douce* (meaning “sweet” in general) in Fig. S8b”.





(d)

Fig. S8. Text mining in French. (Continued) (d) An experiment parallel to [1, Fig. 6]. Non-topics first undergo the same ballpark screening as in [1, Fig. 6], then their  $s_R(\mathbf{b}_i^{\text{en}}, \mathbf{b}_j^{\text{fr}})$  scores are sent to the Hungarian algorithm for bipartite matching.

DIT MR N SIEUR VOIR ALLES DEMANDA ROCHESTER VOULE REGARDAR MME VENIR AIMER PENSE BONS SEMBLANT PREPARE PARLER SAVOIR PORTE  
MÈLE MOISELLE JANE TROUVA ENTENDRE GRANDE JOUR SEULE PETITE CHOSE PASSER YEUX CHAMBRE RÉPONDRE DEUX TEMPS PARTIE DONNER LEVÉE BEAUTE FEMME LONGTEMPS HEURE MAIN LIT ARRIVÉE MOMENT NOUVEAU JOHN DÉSIRER ENFANTS MORT VIE PLACE CONTINUER COEUR PRÈS JEUNES MAÎTRESSE LAISSE OUVRIR MAISON SENTIR CHERCHER CROIS QUITTER FILLES HOMME FAIRFAX SAINT TENIR VOYAGE  
MURÀCROS SORTIE MARIE FROIDE FEU VOIX NUIT FIGURE ÉCRIA FORCE MOTS TRISTE APERCUS ANS ESPRIT SERVIR ADÈLE PLAISIR ENTRE TÊTE NATURE FORTUNE BESSIE ÉLÈVES ASSIS CÔTÉ OH MONDE MIS CERTAINEMENT HEUREUSE  
SOIR BRILLAIS PAROISSE ENFIN SILENCE CRAINS SENTIMENT EH TRANQUILLE REED ESPRÉSSE MALADE LIVRES FERME DEMEURER ATTENDRE COUP DIEU CONNAISSANCE SUIVRE APPELLE BAS ARRÊTA PEINE POINT TROIS AIR POSSIBLE TERRE PREMIERE CROYAS RAPPELÉ APPROCHÉ JOYE TOMBE GARDE MANIÈRE CALME BESOIN CHÂTEAU VIVRE BRUIT DOUTE THORNFIELD SOURIRE MÈRE POSITION NOIRS SOUFFRANCE RESTE DESCENDRE COMPRENDRE DOULSURE DEVANT BRAS COMMENCE OUBLIÉ APPRÉHENDER TRAITS MATIN INGRAM RECUE CESSÉ FENÊTRE SALLE ÉTRANGE REVENU VIEILLER REPRIT ÉLOIGNE PAYS RESTER SOEURS VRAI TABLE CHAISE LUMIÈRE AGITÉE FRAPPÉ AJOUTA MILIEU RETIRE RETOURNA EYRE MASON SOMBRE CONDUITE CAUSE CIEL BOZ PRÉSENCE FIN ROBE SÛRE LIEU NOM GOÛTEHANTE EXAMINER PAROLE INSTANT RECONNAISSANCE COUCHER MILLES TIERS DOMESTIQUES PÈRE ÉCLAIRÉ PAUVRE LOIN DIANA ÉCOUTER OCCUPER JOUES DURE AMIS ÉTRANGER PÂLE AVANCÉ PIEDS MARCHE LIBREMENT ÉTUDE PROFOND ORDRE MANGER HELENÉ CHEVEUX ATTENTIVEMENT DOUCE SENS VÉRITÉ MINUTE  
LUNE CRATE TRAVAIL CHEMIN TOURNÉ REMARQUE MOUVEMENT ENTRA HAUT SATISFAIRE VÊTEMENTS PIERRE EAU DROIT ÉCOLE DOUX ENVOYÉ TEMPLE CACHÉ SÉPARÉ BRISÉ SIMPLEMENT RAPIDE FATIGUE EFFORTS ÉTENDRE PLEINE FRÈRE SOLEIL ROUTE ROUGE PROPRE ÂME EXCITER CHARMÉ  
RICHE CREATURE VOITURE LOWOOD GÉNÉRAL CHANGER PRÉPARÉ SEMAINE CONVERSATION RIVERS IMPORTÉ CARACTÈRE PERDRE PLUIE VOLONTE SOIN FAIBLE RAISON MOIS ABOUD PROMÉTHÉE ÉCRIT VOYAGE JETÉ MARIAGE FAMILLE HUXLEY VENT ETAT PERMETTRE  
INTÉRÊT DÉSESPOIR COURSE SAUVAGE RÊVÉ COMPAGNIE IDÉE ÉPROUVE DECLARA HABITAGE VILLE SALON RIDEAUX RENDRE QUESTION RESSEMBLAIS AIDER PRÊTER ENTRER ESCALIER RIRE ÉPOUSER BUT FAUTE EXISTENCE ENTHÈRE BOIS BONHEUR VAISSE HABILLE BLANCHE FLEUR  
COUVERT EXPLIQUE ALIBI ORIELLES SUPPORTER OR METTRE CHEVAL PAIN FACILE DEMI ENTOURAPPORTES CHER MONTAGNES REPÈRE MALHEUR REUSSIR COULEUR PARENTS NUAGE LETTRE GEORGIANA RECOUVERTE SOUPER PUR UNION FOLLE COUC TACHE CHAMPS OFFREZ INDIQUEZ EMBRASSÉE VENUE JETAZ ADDRESSE  
SOLITAIRE MAUVAIS BONTE ARBRES SANG ATTACHE PENCHÉ NOUVELLE ACCOMPLIR ENVIE EXPRESSION MOTIF SUFISSE BRULEZ LEGÈREMENT EGALLEMENT REPOS SOUVENIR LARMES ACORDÉON ABANDONNE SONNÉZ EXPIMAS TROMPEZ DECIDEZ MINISTRE PEUR ESSAYER THOULE RACONTE ENSEMBLE COMPLÈTEZ NÉCESSAIRE LUTTE AUSSIUT RATONZ IMPOSSIBLE DIX ARGENT TRAVESSIE VERBÉ DÉCOUVRIR APPUYEZ REMPLIR DESTINÉ CHARGE RESTA DIFFICULTÉ CONTENTE VOILE RUE ROSES ORDINAIRE SUITE QUATRE CING TOUCHER VISITEZ ETONNEZ ENFERMER OBÉIR ORGUE NOBLE FLAME FOYER PRÉFÉREZ PARDONNEZ EFFORÇEZ SUPPOSEZ AFFAIRE  
THÈME FRANÇAIS HABITUDE OBSCURITÉ SOUTIEN RESPECTÉ EVILLE DESSINÉ CONFIANCE ENDORMIS REPAS HISTOIRE VINGT INTENTION ÉTAT AFFECTION ÉGLISE CONTRAIRE AVERTIS ALLUME TENTEZ LINTEZ SÉVERE MANOIS EMPARE ANDIE CHAUSSE CHANTE DIGNE CHAPEAU SUET CORSET CÔTEZ APPARTENÉS ATTENDRE ARBREZ GAIEZ APARTEZ OUD  
GATESHEAD DEBOUT CORRIDOR JUGÉ INQUIETÉ ETEINQUÉ FOND ALLEZ AVERTIS ASSEZ MANTEAU COURAGE CONSCIENCE BIBLIOTHÈQUE DINER REFUSE PLEUREZ RENCONTREZ DISPOSÉ DIRECTRESS ANNONSEZ ADMIRABLY CHEVÉ PUISSEZ FERAS PARTAISEZ SECRET OMBRE EPARS CAS TAILLE INUTILE EPAILL COIN GENS ENVELOPPE CONSEIL ACTIVER JOUISSANCE ENVIRON COMPTÉ MARI CONTEMPLAZ  
DISTINGUÉ PROMISE BAISSE RÉFLÉCHIE MENHISZ ESTRAYANT JOLIE VAGUE DANGER OISEAUX MEMBRES JARDIN DOIGTS COMPRIS FRANCHIE VEILLEZ MÉPRIS BIACONEX EXCUSEZ SYMPATHIE BLEU BLESSÉE SORRIEZ PROJET CHIEN ELIZA RENONCEZ TRESSEZ DÉTESTEZ SÉRÉNITE TOILETTE TEINTE TERRIBLE TAPIS OLIVER BOUCHEZ SOUAISSEZ MEMBREZ PLAINDEZ FAVSKEZ EMPÈCHER  
SINCERITÉ REMUE DISPARAÎTRE SOIE SECOND RÉSOLUTION PORTAIT POTHINE NEUF MOYEN COURIR PIÈCE NEVER MORCEAU GATE ENTE DEJUNER MILLCOTE ESSUITE SUBIR RUINE RENDY CHAUSSE ASSUREZ RENDRE BONTE PRINCIPES LARG DÉTAILS BOUCLES BOUCHE GENOU SOCIÉTÉ ACCOMPAGNEZ TUÉZ RENDAIS CONSIDÉREZ NEIGE LECTURE INSTITUTTY PAPIER LANGUANT CHANDELLE AGRÉABLE SCOTLAND  
TANTE SIX POOLE HU DENT AUGURE AH REPONSES DÉTRUIREZ RAPPORT GRAVEZ FIDÉLÉ CONNU BATTRE ACTE CIRCONSTANCE VÉGÉOZ CLASSE VIE SURPRISE LENTEMENT LEÇONS COEUR ACCÈS SOPHIE HIVER GRÉVÉE ASPIRATZ INTEROMPT FICHÉ PERSUASIF HUMIDE HÔTEL RAMBO FÉE DÉCLINAIS  
LEAH IMMÉDIATE LENDemain POURPRE OBSERVEZ EMPORTEZ RÉVÉLÉEZ PRÉCIPITEZ COMMENT ATTRAPE SERREZ HAFÉ FIXEZ GAGNE SOUFFLEZ REMERIC PLATEAU OEUVRE CORTAISIE VALLEE SOLEIL QUINZE LOF LINE IMPATIENCE SECOURS RETIRE PILOTE BLANC GROSSESSE CHASSEZ ATTRAPEZ ELEPHANT ELEPHANT CHAUFFEZ FLOTZ CONTRAIRE SENSUELLEZ ARDÉZ ANGELE PITÉ MAMIE PERDU RIEN EFFRAYEZ ARRANGEZ MÉLUCA  
CONVERSEZ VAINCZ SACRIFICE MYSTÈRE COUSINS CARTER BORD ANCIENT NORME FIVEZ EDIE CHIREHN BANC SOCIÉTÉ PRAGIQUE PERSONNES ORKNEY ENDROT DISTANCE HIER HIBERNER ABRI DÉPOSER GLASSZ ENSENGE VIOLENCE TAQUE PRETENDZ HERZOG FRÂCHEZ ANGLAIS LAUD CLOCHE AVEUGLE SENATION ÉDOUARD BURNS ENDUREZ RÉUNIREZ INSTRUZEST AMUSEZ INDIFFÉRENTEZ ERERENT CONSOLAREZ REPARTIEZ PLIS MEUBLE  
INSTINCT CHILLER ÉTROIT SCÈNE RESTA RÉCIT OSCUR EVENEMENT ASPECT ARRÊTER YOUT SCATCHER

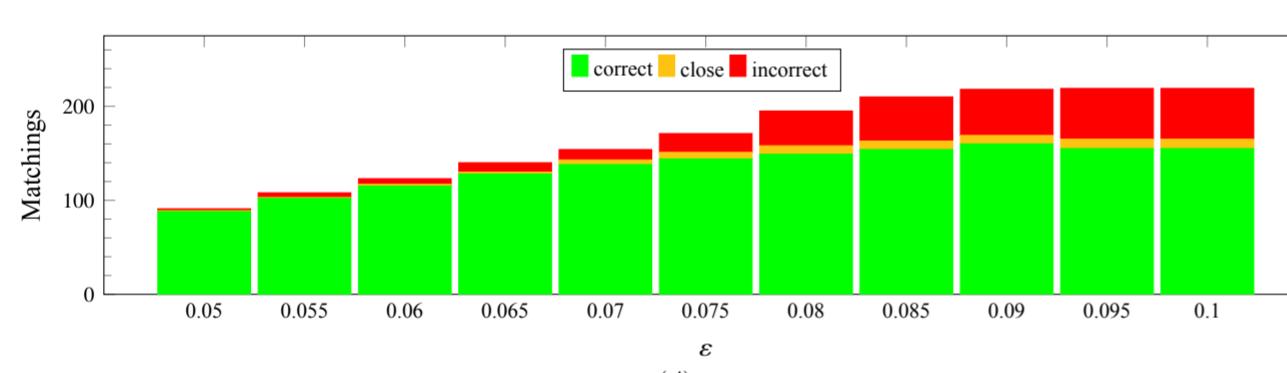
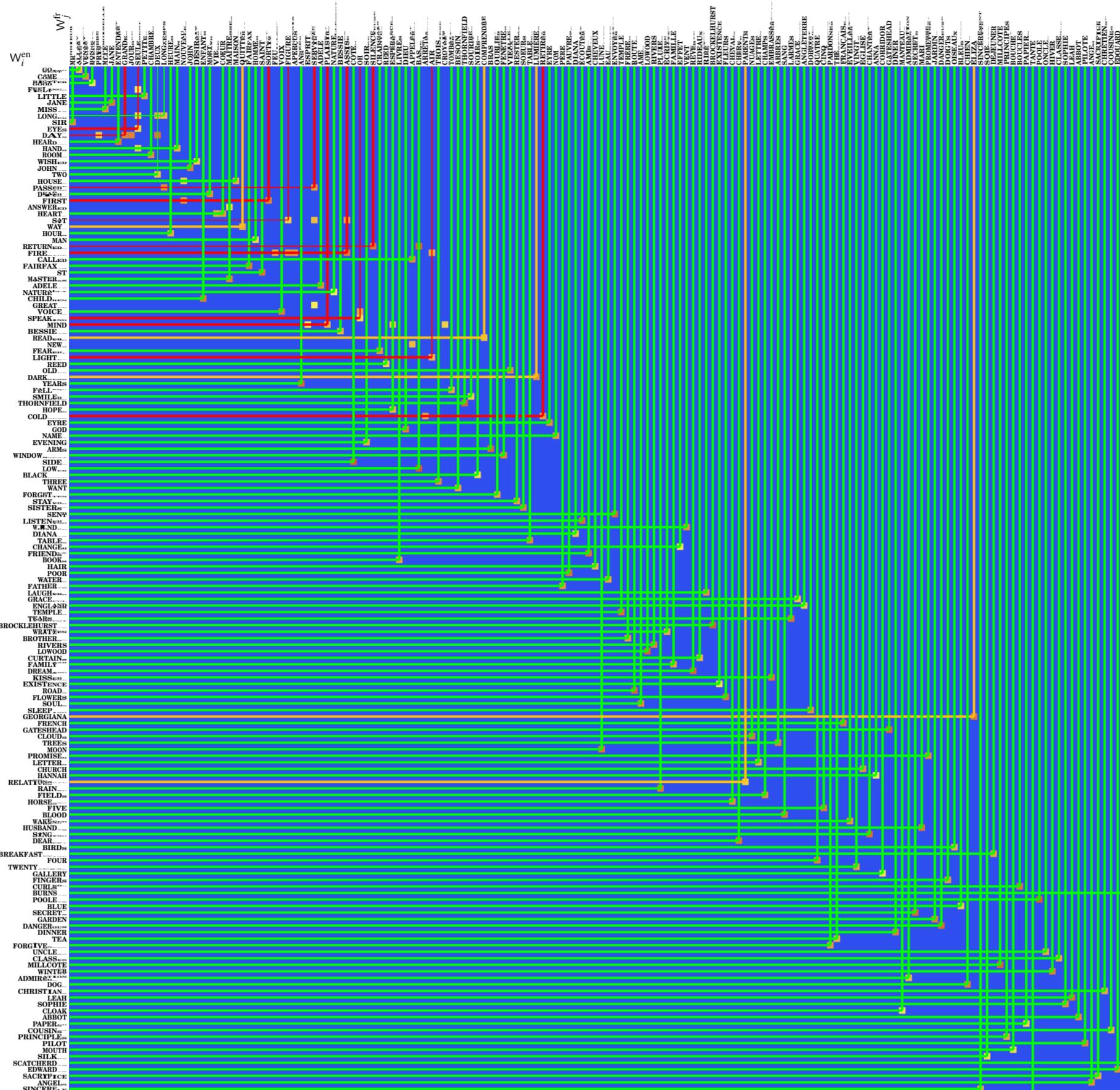


Fig. S8. Text mining in French. (Continued) (a') Statistically identified topics ( $n_{ii} \geq 20$ ) in a French version of *Jane Eyre*, with the same color encoding scheme as Fig. S3. (b') Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{fr}})$  between selected topics in English and French versions of *Jane Eyre*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) rank of the similarity score for the optimal node. Green (resp. red) cross-hair indicates a correct (resp. incorrect) match. Amber cross-hair marks a link between distinct concepts that share the same hypernyms. (c') Results from control experiments with different choices of the  $\varepsilon$ -parameter in ballpark screening criteria (1.13).

ESPÈCES FORMÉES VARIÉTÉS DIFFÉRENCES GRANDE NOMBRE CROISEMENT DEUX ORGANES CERTAINES MODIFICATIONS PARTIES EXISTENCE NATURELLE ANIMAUX SELECTION ENTRE PLANTES CAS CONSIDÉRABLE NATURELLES CARACTÈRES GENRES PRODUISSENT EXEMPLE GÉNÉRALEMENT INDIVIDUS DISTINCTES DESCENDANTS TROUVENT CONDITIONS GROUPES SEULEMENT LONGUEURS VUE IMPORTANTE NOUVELLES HABITANTS ÉTAT POINT PRÉSENTENT DIVERSES REMARQUES CLASSEIFICATION DIRE SAGE SEMBLABLES DÉVELOPPEMENT PÉRIODE EFFET RACES CONFORMATION EXPLIQUER OR CONTINENT RÉGIONS PETITES ÎLES OISEAUX FÉCONDITÉ RAPPORTS COMMUNES PRINCIPES VIVANT PREMIÈRE MR TENDANCE RÉSULTAT OBSERVATION DEGRÉ DIFFICULTÉ DOMESTIQUES TERRESTRES RESSEMBLANCE GRAINES ACTUELLEMENT FLEURS CAUSES POSSIBLES TEMPS SUPPOSES ORDINAIREMENT ANCIENNES STERILITÉ INSECTES CHANGEMENTS FACILEMENT PLACÉS ADMETTRE DOUTE ACTION INTERMÉDIAIRES HYBRIDES PAYS AVANTAGEUSES FORT HABITUDES COMPLÈTEMENT COMPRENDRE ADAPTATION TROIS MYSTÈRE ÉLEVÉES PARFAITEMENT PARFAITEMENT CONSERVÉES DÉMONTRÉ THÉORIE MANIÈRE APPARTENANT RELATIVEMENT VIE ENTIERE PORTÉE CONSTITUÉE DEVIENTS PARENTS EAU FAÇON SERVANT UTILES PRODUCTIONS CRÉATION HOMME ANCÉTRA ANALOGUES GÉOLOGIQUES PARTICULIÈRE RESTES MEMBRES SUJET FAMILLES OUTRE REPRODUCTEURS VOISINES VERTU PROBABLEMENT SORTE USAGE ÉTENDUE ÉGALEMENT LENTEMENT LUTTE MOYENEMENT AGIT SUCCESSIVES INFÉRIEURE JEUNES REPRÉSENTANTES PROUVENT NÉCESSAIREMENT CITER APPAREILS SPÉCIALES AUGMENTATION CLIMAT EUROPE OEUVRES NOURRITURE PERFECTIONNÉES RAISON FAVORABLES DÉTERMINER SUITE PRÈS AILLEURS MALES SUBIÉTÉ ARRIVÉ SEXES POISSONS CONSÉQUENT HUI AUJOURD PIGEONS ABEILLES PASSÉ DIRECTEMENT MAMMIFÈRES SÉRIE SAUVAGE ARBRES VOYONS CONSÉQUENCE HAUTE RUDIMENTAIRES ORDRES ACCUMULATION PROBABLE MONDE DISPARAÎTRE VERS PRIMITIVEMENT ÉLOIGNÉS ATTRIBUER COUCHES ACQUISIES AILES ORIGINE AUSTRALIE ENSEMBLE PROVENANT ADULTE EXTRÉMEMENT ALLIÉS ENFIN UNIS MESURE FOSSILES ÉVIDEMMENT DOCTEUR CHEVAUX CHAPITRE DÉPEND RARES EXTINCTION PREUVE SUD MILLION ACCIDENTELLEMENT DISTRIBUTION CONTRAIRE RÈGLE VOLANT RELIÉS AFFECTÉS ADMIRABLEMENT SECONDAIRES TYPE SURFACE EXAMINÉES PERMET HYPOTHÈSE OCEANIQUES FEMELLES ÉTONNANTE SOI FRAPPANT GLACIAIRE SOUCHE BRANCHES EXPOSÉ UNIFORMEMENT COULEUR AFFINITÉS CONSTATE LIGNES INTERVALLES NORD DÉCOUVERT RAPIDE PROPRE MER EXACTEMENT ENSUITE TRANSPORT PROFONDEUR PERSISTANCE GROSSEUR DIVERGENCE BASE VALEUR MONTAGNES FONCTIONS ÉTUDE GRADUELLEMENT AFFIRMÉ IMPOSSIBLE BUT BEC LIMITES LOI GLOBE GRADATIONS INFLUENCE LOIS NAISSANTE PENSES INDÉPENDANTE STRUCTURE DISPOSITION COURTE PATTES QUANTITÉ LARVES COMPLEXES ISOLÉES RECONNUS SENS VINGT RÉCÉDEMMENT CONNAISSANCE CONTRAIRE ATTEINT RÔLE IMMENSE CONCURRENCE ETÉ SÉPARÉES RÉPANDUE MARINE AJOUTÉE COURS TRAITÉ POSITION TRANSITION OFFRENT JOURS CIRCONSTANCES QUESTION OPINION DIMINUÉES ANTÉRIEURE FOURNIS SOUTIENS PROGRÈS ÉLÉMENT TERME PHYSIQUES POURVUS FRÉQUEMMENT JOUR NID SOUMISÉS SÉPARÉS CHERCHE FEUILLES MIGRATION VOUÉ ÉTABLI AMIENS MÉTIS INUTILE COQUILLAGES SPÉCIFIQUES DISTANCE CHIENS TRAVAIL SON OBJECTION ÉPAISSEUR RICHES SOURIS RAPPROCHÉES TRANSFORMATIONS RÉELLEMENT RAPPELER FOND EXTRAORDINAIRE EXCEPTION INTÉGRALEMENT VÉRITABLE TRANSMIS SOMME EXTÉRIEURE CORRESPONDANT CHAÎNONS JUGER IDENTIQUES ANGLETERRE GÉOGRAPHIQUE DÉFAUT CORPS INFINI TEMPÉRÉES CHANCES ARCHIPEL QUATRE PURE INCONNU DEMANDE QUEUE OEIL FORCE RIGoureusement REGARDÉ RÉCENTE ESPACE PHASES PERDUE PLAN MULTIPLIÉES MERIDIIONALE ÉCONOMIE AUTEURS PROFESSEUR LOIN OPPOSÉS ESSAIE EMPLOI CHANGÉ DIVISIONS CONSTRUITES CONDUIT BEAUTÉ TERTIAIRES PLAÎNES FINISSEMPÊCHÉ RÉCIPROQUES EXCLUSIVEMENT MÂCHOIRES RETOUR PERMETTRE ROCHES HEMISPHERE RAYON INSIGNIFIANT ENVIRON CONNUS SEPTENTRIONALE ESCLAVES DOUCE CORRÉLATION APITUDE YEUX TÊTE HISTOIRE SOULEVEMENT LOCALITÉ ANALOGIE ATTENDRE TENDRE RÉUNIS DÉTRUIT FROID ASSURÉ INDISSOCIABLE GÉNÉALOGIQUE RÉCOLTÉ UNIQUE CHAUVEs PEINE IMAGINÉE SIMULTANÉEMENT IMPLIQUE ENTENDU MATÉRIEL PROVIENNENT OBJET SECRÈTE COUVERTES ABONDANTES RECHERCHES ESSENTIELLEMENT CAPTIVÉES CROISSANT ACTE EMBRYON APPARITION RÉDUITE PUISSE MORT RECULÉ HUAT DOMINANTES SUCCESSION OBTENUE ENORME PLUMES TIREZ EXTERMINÉS VACUÉES CULTIVÉES CONSTANTE AIR RÉSISTE FRUITS AQUATIQUES RAREMENT POUSSÉS MANIFESTEZ EMPORTÉS AMÉLIORATION OS STIGMATE FACULTÉ RAISONS FLOTTANTE SIGNALISÉ VIGOUROUS BATTACHE PARLER EXCELLENT DISCUTER BESTIAL RAYON INDIGÈNES DIX POSSESSIFS BOTANISTES REVENU ATTENTION DESTRUCTION GREFFE RÉUSSI FIXÉ MUTUELLEMENT TRACE SEDIMENTS DÉVIATIONS CORNES SUR HEURES ALIMENTS AFRIQUE MIVART GRÂCE GARTNER APPARTIENNENT PRÉCISEMENT CESSÉ BARRETTES TROMPE RÉGULIÈREMENT PRONONCÉ PICOCHE TRANSMETTRE PARALLÉLISME NUISIBLE MASSE MARCHÉ ETAMINES RÉGARDÉS DIVERSITÉ CLAIREMENT TAILLE NOM MODE BÉTAIIL COMMENCEZ PRÉTENDRE APPELE MÉTAMORPHOSE NEUTRES ÉMIGRE VASTE IDÉE AUXQUELLES BISET MARQUES SINGULIÈRE IMPARFAITE FLORE DÉCRIT AUTORISÉ AUTOMATIQUE RECUEILLE NAGEOIRE ARGUMENT LAPPE CONVAINCU PRÉCOCÉ MOUTONS MILIEU JAMBES SUPÉRIEUR SIX HOMOLOGUES GÉRÉ FAUNE DOUZE ARCTIQUES MERS OPERES ÉPROUVE CONTENUE COMPLÈTE VÉNÉS VÉGÉTAL SURPRISE RECOUVRÉ PERY INTRODUCEZ EMPAREZ PALÉONTOLOGISTES CONCERNÉS INÉDITABLE MÈRE LARGE CENT BORD PISTILS OWEN OIE EXTREMITÉ AFFAISSEMENT FORMICA SÉGMENTÉS VISITES EFFECTUÉS EXÉCUTION COURANTS CULBUTANTS CIRRIPÈDES CAVERNES SIÈCLES PARASITES MÉTIERS COMMENCEMENT HOOKER RÉSISTERÉ MÉLANGE MAINTEINÉ ÉTRANGER SUPPLANTER HABITUELLEMENT ÉQUATRÉ ALLES TRANCHÉES REVENUS LIBRES ÉCHAPPE ARRÊTÉ ZONES TERRAIN LYELL LUMIÈRE AVANTAGE CINQ SEMAIS FINISSE ESTIME CHASSE TOTALISANT PARIÉTÉ GÉOLOGIQUES CUMBRIENNE ULTIMEURÉ RÉPONSE PHYSIOLOGIQUE PAPILLONS TIGES STATIONS SOMMET SANGUINEA RANG POILS MALAIS ENTER ANGLAIS AIMÉE REUSE ASSERTION HERBERG ATROPHIE REPOSE MIGRÉ INSISTE PONDENT NATATOIRE DÉCIDER VOLUME DONNE DISCUSSION DIAGRAMME SITUÉES GERMÉRENT RONGEURS LEVYEN IMMIGRÉS PROPORTION LIEN ZÉLANDE DISPERSION RÉSULTAMENT REMPLIE NOIRE BLEUÉ APPÉTITÉ SUPPLANTÉ HABITUELLEMENT ÉQUATRÉ ALLES TRANCHÉES REVENUS LIBRES ÉCHAPPE ARRÊTÉ ZONES TERRAIN LYELL LUMIÈRE AVANTAGE CINQ SEMAIS FINISSE ESTIME CHASSE TOTALISANT PARIÉTÉ BOURDONS DENTS AUTREMENT REVÊTEZ ATTAQUEZ PIERRE FIGURE ENFOURAIS DISSEMÉNÉS SURVIVRE SUPPORTER RESTREINT NERFS EFFICACÉ DROIT BRAS BALEINE CONCÈDE AQUÉ AMBIANCES TRAVERSÉE PROVOQUE POULE ENLEVÉS ATTACHEZ SONGÉ REMONTE PERPETUER INÉGALÉ CONFONDRE HÉRITÉ ELIMINÉE CONFONDRE HÉRITÉ PRATIQUE OCCASIONNELLEMENT INSECTE FREINS ESPÈCE DOIGTS COMPORTE PROSPÉRÉ ILLEGITIMES ANE VITÉS SAISIR PLANTS PAON METTRE ÉVANISSEMENT GLANDES DOCUMENTS SERT RÈGNE MADÈRE I EMBRYONS CANDOLLE BOUCHE CHAUME SUBSEQUENTES CHAUDEURS HÉRITÉ CONFONDRE HÉRITÉ PRATIQUE OCCASIONNELLEMENT INSECTE FREINS ESPÈCE DOIGTS COMPORTE PROSPÉRÉ LAMELLES HERMAPHRODITES GROS BRÉTAÏNE APPARENCE TOUR REPTILES NÉANMOINS KILOMÈTRES ARRANGEMENT

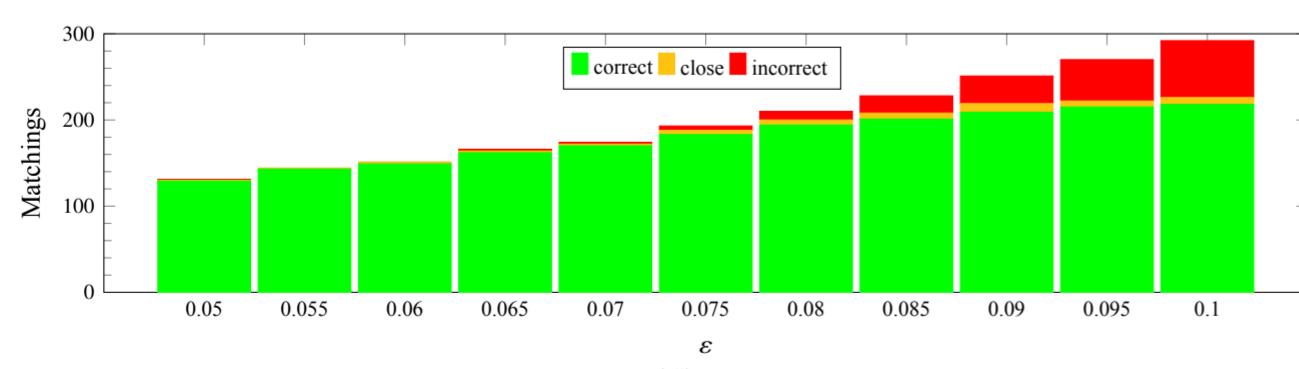
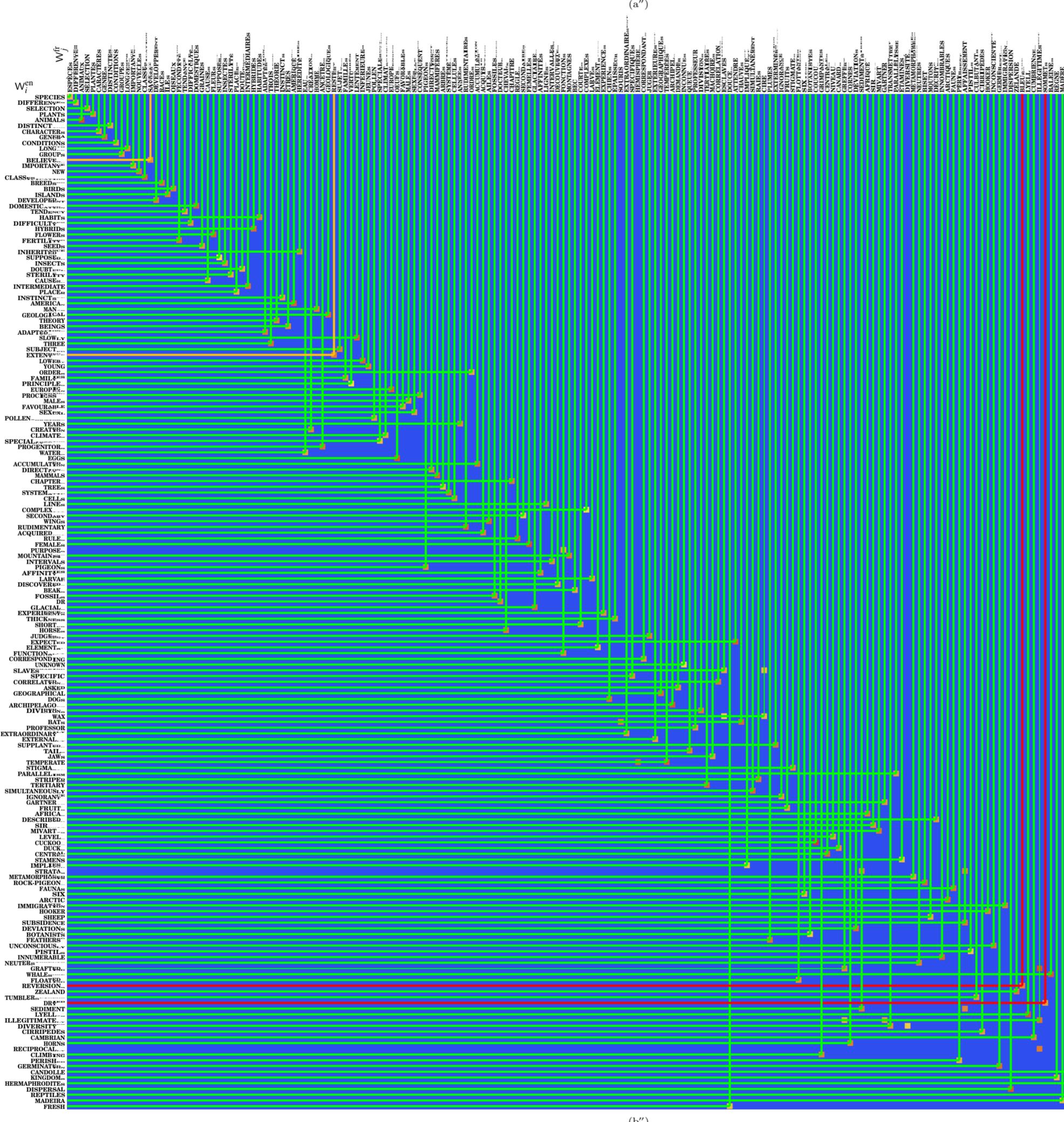


Fig. S8. Text mining in French. (Continued) (a'') Statistically identified topics ( $n_{ii} \geq 20$ ) in a French version of *Origin of Species*, with the same color encoding scheme as Fig. S3. (b'') Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{fr}})$  between selected topics in English and French versions of *Origin of Species*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) rank of the similarity score for the optimal node. Green (resp. red) cross-hair indicates a correct (resp. incorrect) match. Amber cross-hair marks a link between distinct concepts that share the same hypernyms. (c'') Results from control experiments with different choices of the  $\epsilon$ -parameter in ballpark screening criteria (1,13).

### 6.3 Modified Porter stemming algorithm for Latin

**Definition 6.21** (Latin stop words). If a word belongs to the following list<sup>79</sup>:

*a, ab, abhinc, absque, abusque, ac, acta, actu, actum, acturus, actus, ad, adaeque, adhic, adhuc, adusque, adversum, adversus, agam, agamini, agamur, agamus, agant, agantur, agar, agare, agaris, agas, agat, agatis, agatur, age, agebam, agebamin, agebamur, agebamus, agebant, agebantur, agebar, agebare, agebaris, agebas, agebat, agebatis, agebatur, agemini, agemur, agemus, agendi, agendo, agendum, agendus, agens, agent, agentur, agere, agerem, ageremini, ageremur, ageremus, agerent, agerentur, agerer, agerere, ageris, ageres, ageret, ageretis, ageretur, ageris, ages, aget, agetis, agetur, agi, agimini, agimur, agimus, agis, agit, agite, agitis, agito, agitor, agitote, agitur, ago, agor, agunt, agunto, aguntur, alia, aliae, aliam, aliarum, alias, alicui, alicuius, alicujus, alie, alii, alis, alio, aliorum, alios, aliquae, aliquando, aliquantulum, aliquarum, aliquas, aliquem, aliqui, aliquibus, aliquid, aliquis, aliquo, aliquorum, aliquos, aliquot, aliquotiens, aliter, aliud, alium, alius, alter, altera, alterae, alteram, alterarum, alteras, alteri, alteris, alterius, altero, alterorum, alteros, alterum, an, ante, antea, apsque, apud, at, atque, attorque, aut, autem, caetera, cætera, cæterae, cæteram, cæterarum, cæterarum, cæteras, cæteras, cætere, cæteri, cæteris, cætero, cætero, cæterorum, cæterorum, cæteros, cæteros, caeterum, cæterus, cæterus, causa, cetera, ceterae, ceteræ, ceteram, ceterarum, ceteras, ceter, ceteris, cetero, ceterorum, ceteros, ceterum, ceterus, circa, circum, cis, crita, com, concoque, contorque, contra, coque, coram, cui, cuidam, cuiquam, cuique, cuius, cuiusdam, cuiusquam, cuiusque, cuiusvis, cuivis, cuius, cuiusdam, cuiusquam, cuiusque, cuiusvis, cum, cur, de, decoque, dein, deinde, denique, denuo, deque, detorque, dum, e, ea, eadem, eae, eadem, eam, eandem, earum, earundem, easdem, ecquando, egeram, egeramus, egerant, egeras, egerat, egeratis, egere, egerim, egerimus, egerint, egeris, egerit, egeritis, egero, egerunt, egi, egimus, egisse, egissem, egisset, egisses, egissetis, egisti, egistis, egit, ego, ei, eidem, eis, eisdem, eius, eiusdem, ejus, ejusdem, en, enim, eo, eodem, eorum, eorundem, eos, eosdem, eram, eramus, erant, eras, erat, eratis, ere, erga, ergo, erimus, eris, erit, eritis, ero, erunt, es, esse, essem, essemus, essent, esses, esset, essetis, est, este, estis, esto, estote, et, etiam, etiamsi, etsi, eum, eundem, ex, excoque, extorque, extra, fac, face, facere, facerem, faceremus, facerent, faceret, faceretis, faciam, faciamus, faciant, facias, faciat, faciatis, faciebam, faciebamus, faciebant, faciebas, faciebat, faciebatis, faciemus, faciendi, faciendo, faciendum, faciendus, faciens, facient, facies, faciet, facietis, facimus, facio, facis, facit, facite, facitis, facito, facitote, faciunt, faciunto, facta, factae, factam, factarum, factas, facte, facti, factis, facto, factorum, factos, factu, factum, facturus, factus, feceram, feceramus, fecerant, feceras, fecerat, feceratis, fecere, fecerim, fecerimus, fecerint, feceris, fecerit, feceritis, fecero, fecerunt, feci, fecimus, fecisse, fecissem, fecissemus, fecissent, fecisses, fecisset, fecissetis, fecisti, fecistis, fecit, fi, fiam, fiamus, fiant, fias, fiat, fiatis, fiebam, fiebamus, fiebant, fiebas, fiebat, fiebatis, siemus, fiendi, fiendo, fiendum, fient, fierem, fieremus, fierent, fieres, fieret, fieretis, fieri, fies, fiet, fietis, fimus, fio, fis, fit, fite, fitis, fito, fitote, fiunt, fiunto, fore, forem, foremus, forent, fores, foret, foretis, forms, forsani, forsitan, fortasse, fortassis, frequenter, frequentissime, frequentius, fueram, fueramus, fuerant, fuerat, fueratis, fuere, fuerim, fuerimus, fuerint, fueris, fuerit, fueritis, fuero, fuerunt, fui, fuimus, suisce, suissem, suissemus, suiscent, suisces, suisset, suissetis, fuisti, fuistis, fuit, futura, futurae, futuram, futuraram, futuras, future, futuri, futuris, futuro, futurorum, futuros, futurum, futurus, habe, habeam, habeami, habeamur, habeamus, habeant, habeantur, habear, habeare, habeas, habeat, habeatis, habeatur, habebam, habebamini, habebamur, habebamus, habebant, habebantur, habebar, habebare, habebaris, habebas, habebat, habebatis, habebatur, habebere, habeberis, habebimini, habebimur, habebimus, habebis, habebit, habebitis, habebitur, habeo, habebor, habebunt, habebuntur, habemini, habemur, habemus, habenda, habendae, habendam, habendarum, habendas, habende, habendi, habendis, habendo, habendorum, habendos, habendum, habendus, habens, habent, habente, habentem, habentes, habenti, habentia, habentibus, habentis, habentium, habento, habentor, habentur, habeo, habeor, habere, haberem, haberemini, haberemur, haberemus, haberent, haberentur, haberer, haberere, habereris, haberes, haberet, haberetis, haberetur, haber, haberier, haberis, habes, habet, habete, habetis, habeto, habetore, habetote, habetur, habita, habitae, habitam, habitarum, habitas, habite, habiti, habitis, habito, habitorum, habitos, habitu, habitum, habitura, habiturae, habituram, habituram, habituras, habiture, habituri, habituris, habituro, habiturorum, habituros, habiturum, habiturus, habitus, habueram, habueramus, habuerant, habueras, habuerat, habueratis, habuere, habuerim, habuerimus, habuerint, habueris, habuerit, habueritis, habuero, habuerunt, habui, habuimus, habuisse, habuissem, habuissemus, habuissent, habuisses, habuisset, habuissetis, habuisti, habuistis, habuit, hac, hae, haec, hanc, harum, has, haud, hi, hic, hinc, his, hoc, horum, hos, huc, huic, huius, hunc, iam, ibi, id, idem, igitur, iis, illa, illae, illam, illarum, illas, ille, illi, illis, illius, illo, illorum, illos, illud, illum, immo, in, incoque, inde, infra, infrequenter,*

<sup>79</sup>Our list of Latin stop words is based on [https://wiki.digitalclassicist.org/Stopwords\\_for\\_Greek\\_and\\_Latin](https://wiki.digitalclassicist.org/Stopwords_for_Greek_and_Latin) and <http://snowball.tartarus.org/otherapps/schinke/intro.html>, with extensive additions to roughly match their counterparts in English. In particular, we have included all the conjugated forms of *faciō* “make”, *fīō* “become”, *sum* “be” and *possum* “can”. It is worth noting that the conjugation tables for *faciō* and *fīō* partially overlap.

*infrequentissime, infrequentius, inquam, inque, inquit, inquiet, inquietus, inquis, inquisti, inquit, inquitis, inquito, inquiunt, inter, interim, intorque, intra, ipsa, ipsam, ipsarum, ipsas, ipse, ipsi, ipsis, ipsius, ipso, ipsorum, ipsos, ipsum, is, isdem, ista, istae, istam, istarum, istas, iste, isti, istis, istius, isto, istorum, istos, istud, istum, ita, itaque, item, iter, iterum, iuxta, jam, juxta, magis, me, mea, meae, meam, mearam, meas, mecum, mei, meis, meo, meorum, meos, meum, meus, mi, mihi, minime, minus, modo, mox, multa, multae, multam, multarum, multas, multe, multi, multis, multo, multorum, multos, multum, multus, nam, ne, nec, necque, nemine, neminem, nemini, nemo, neque, nequeamus, nequeant, nequeas, nequeat, nequeatis, nequeo, nequeundi, nequeundo, nequeundum, nequeunt, nequeunto, nequi, nequibam, nequibamus, nequibant, nequibas, nequibat, nequibitis, nequibimus, nequibis, nequibit, nequibitis, nequibo, nequibunt, nequiens, nequieram, nequieramus, nequierant, nequieras, nequierat, nequieratis, nequierere, nequierim, nequierimus, nequierint, nequieris, nequierit, nequieritis, nequiero, nequierunt, nequii, nequimus, nequit, nequimus, nequire, nequirem, nequiremus, nequirent, nequires, nequiret, nequiretis, nequis, nequissem, nequissemus, nequissent, nequisses, nequisset, nequissetis, nequisti, nequistis, nequit, nequite, nequitis, nequito, nequitote, nequitu, nequitum, nequiturus, nequivi, nequivisti, nequivit, nihil, nihil, nihil, nihilis, nihil, nihilorum, nihilum, nil, nimie, nimis, nimium, nisi, nobis, nobiscum, nolam, nolebam, nolebamus, nolebant, nolebas, nolebat, nolebatis, nolemus, nolens, nolent, noles, nolet, noletis, noli, nolim, nolimus, nolint, nolis, nolit, nolite, nolitis, nolito, nolitote, nolle, nollem, nollemus, nollent, nolles, nollet, nolletis, nolo, nolueram, nolueramus, noluerant, nolueras, noluerat, nolueratis, noluerere, noluerim, noluerint, nolueris, noluerit, nolueritis, noluerero, noluerunt, nolui, noluius, noluisse, noluissem, noluissemus, noluissent, noluisse, noluisset, noluissetis, noluisti, nolustis, noluit, nolumus, nolunt, nolunto, non, nondum, nonne, nonnulla, nonnullae, nonnullam, nonnullarum, nonnullas, nonnulla, nonnulli, nonnullis, nonnullius, nonnullo, nonnullorum, nonnullos, nonnullum, nonnullus, nonnumquam, nos, noster, nostra, nostrae, nostram, nostrarum, nostras, nostri, nostris, nostro, nostrorum, nostros, nostrum, nulla, nullae, nullam, nullarum, nullas, nulle, nulli, nullis, nullius, nullo, nullorum, nullos, nullum, nullus, num, numquam, nunc, nunquam, nusquam, o, ob, oblique, obtorque, olim, omne, omnem, omnes, omni, omnia, omnibus, omnino, omnis, omnium, optorque, paene, parum, pauca, paucae, paucam, paucarum, paucas, pauce, pauci, paucior, pauciora, pauciore, pauciores, pauciori, paucioribus, paucioris, pauciorum, paucis, paucissima, paucissimae, paucissimam, paucissimarum, paucissimas, paucissime, paucissimi, paucissimis, paucissimo, paucissimorum, paucissimos, paucissimum, paucissimus, paucius, pauco, paucorum, paucos, paucum, paucus, penes, penitus, per, peraeque, plenisque, plerumque, plura, plure, plures, pluribus, plurima, plurimae, plurimam, plurimarum, plurimas, plurime, plurimi, plurimis, plurimo, plurimorum, plurimos, plurimum, plurimus, pluris, plurium, plus, pone, posse, possem, possemus, possent, posses, possetis, possim, possimus, possint, possis, possit, possitis, possum, possumus, possunt, post, postea, potens, potente, potentem, potentes, potenti, potentia, potentibus, potentior, potentiora, potentiore, potentiorem, potentiores, potentiori, potentioribus, potentioris, potentiorum, potentis, potentissima, potentissimae, potentissimam, potentissimarum, potentissimas, potentissime, potentissimi, potentissimis, potentissimo, potentissimorum, potentissimos, potentissimum, potentissimus, potentium, potentius, poteram, poteramus, poterant, poteras, poterat, poteratis, potere, poterimus, poteris, poterit, poteritis, potero, poterunt, potes, potest, potestis, potueram, potueramus, potuerant, potueras, potuerat, potueratis, potuere, potuerim, potuerimus, potuerint, potueris, potuerit, potueritis, potuero, potuerunt, potui, potuimus, potuisse, potuissem, potuissemus, potuisserent, potuisses, potuisset, potuissetis, potuisti, potuistis, potuit, praे, præ, praeter, praetorque, prius, pro, prope, propter, qua, quacum, quadam, quae, quaedam, quaelibet, quaequam, quaeque, quaevis, quale, qualem, quales, quali, qualia, qualibus, qualis, qualium, quam, quamquam, quamque, quandam, quando, quandoque, quanta, quantae, quantam, quantarum, quantas, quante, quanti, quantis, quanto, quantorum, quantos, quantum, quantus, quanvis, quaque, quare, quarum, quarumquam, quarumque, quarundam, quarunvis, quas, quasdam, quasi, quasquam, quasque, quasvis, quavis, queam, queamus, queant, queas, queat, queatis, queatur, quem, quemquam, quemque, quendam, quenvis, queo, queundi, queundo, queundum, queundus, queunt, queunto, qui, quia, quibam, quibamus, quibant, quibas, quibat, quibatis, quibatur, quibimus, quibis, quibit, quibitis, quibitur, quibo, quibunt, quibus, quibuscum, quibusdam, quibusquam, quibusque, quibusvis, quicum, quicunque, quid, quidam, quiddam, quidem, quidlibet, quidquam, quidque, quidvis, quiens, quieram, quieramus, quierant, quieras, quierat, quieratis, quiere, quierim, quierimus, quierint, quieris, quierit, quieritis, quiero, quierunt, qui, quiimus, quii, quilibet, quimus, quin, quiquam, quique, quire, quirem, quiremus, quirent, quires, quiret, quiretis, quiretur, quiri, quis, quisnam, quisquam, quisque, quisquis, quisse, quissem, quissemus, quissent, quisses, quisset, quissetis, quisti, quistis, quit, quite, quitis, quito, quitote, quitu, quitum, quitur, quiturus, quivi, quivis, quivisti, quivit, quo, quoad, quocum, quod, quodam, quodlibet, quomadmodum, quominus, quonodo, quonam, quondam, quoniam, quoquam, quoque, quorsum, quorum, quorumquam, quorundam, quorunvis, quos, quosdam, quosquam, quosque, quosvis, quot, quotusquisque, quousque, quovis, recoque, retorque, rursus, saepe, saepissime, saepius, saltem, sane, se, secum, secundum, sed, semper, sese, si, sibi, sic, sicut, sim, simul, simus, sine, sint, sis, sit, sitis, sive, solum, statim, sua, sua, suam, suarum, suas, sub, subter, sue, sui, suis, sum, sumus, sunt, sunto, suo, suorum, suos,*

*super, supra, susque, suum, suus, tale, talem, tales, tali, talia, talibus, talis, talium, tam, tamen, tametsi, tandem, tanta, tantae, tantam, tantarum, tantas, tante, tanti, tantis, tanto, tantorum, tantos, tantum, tantus, te, tecum, tenus, tibi, torque, tot, tota, totae, totam, totarum, totas, tote, toti, totis, totius, toto, totorum, totos, totum, totus, trans, tu, tua, tuae, tuam, tuarum, tuas, tue, tui, tuis, tum, tuo, tuorum, tuos, tuum, tuus, ubi, ubicumque, ubinam, ubique, ubiubi, uel, uero, ultra, umquam, una, unae, unam, unarum, unas, unde, undique, une, uni, unis, unius, uno, unorum, unos, unquam, unum, unus, usquam, usque, ut, uter, uterque, utique, utra, utrae, utram, utrarum, utras, utri, utribique, utris, utrius, utro, utroque, utrorum, utros, utrum, valde, vel, versum, versus, vester, vestra, vestrae, vestram, vestrarum, vestras, vestri, vestro, vestrorum, vestros, vestrum, vix, vobis, vobiscum, vos, voster,*

or is spelt as a word from the list above, followed by suffix *que*, then we consider it a Latin stop word. All the Latin stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list.  $\square$

Similar to what we did in the algorithm for French, we need to define string patterns **goLatin**, **goodLatin**, **betterLatin**, **bestLatin** to cover some irregular inflections:

*eam, eamus, eant, eas, eat, eatis, eatur, eo, eundi, eundum, eundus, eunt, eunto, i, ibam, ibamus, ibant, ibas, ibat, ibatis, ibatur, ibimus, ibis, ibit, ibitis, ibitur, ibo, ibunt, iens, ieram, ieramus, ierant, ieras, ierat, ieratis, iere, ierim, ierimus, ierint, ieris, ierit, ieritis, iero, ierunt, ii, iimus, iit, imus, ire, irem, iremus, irent, ires, iret, ieritis, ieretur, iri, is, isse, issem, issemus, issent, isses, isset, issetis, isti, istis, it, ite, itis, ito, itote, itu, itum, itur, iturus, itus, ivi, ivisti, ivit — “go”;*

*bona, bonae, bonam, bonarum, bonas, bone, boni, bonis, bono, bonorum, bonos, bonum, bonus — “good”;*

*melior, meliora, meliore, meliorem, meliores, meliori, melioribus, melioris, meliorum, melius — “better”;*

*optima, optimae, optimam, optimarum, optimas, optime, optimi, optimis, optimo, optimorum, optimos, optimum, optimus, optuma, optumae, optumam, optumarum, optumas, optume, optimi, optimis, optimo, optumorum, optumos, optumum, optumus — “best”.*

However, we are not going to treat these strings as stop words<sup>80</sup>, so as to be consistent with their counterparts in other languages.

### 6.3.1 Effective spelling and essential root

**Definition 6.22** (Latin Vowel Extensions). Hereafter in §6.3, the symbol  $V^*$  stands for any member from the list  $\{a, e, i, o, u, y\}$ , the so-called Latin vowel extensions. In line with the multiplicity notations introduced in Definition 3.3, the symbol  $V_{m_0}^*$  stands for a text string formed by consecutive appearance of one or more (not necessarily identical or distinct) members from the set of Latin vowel extensions.

Dual to the notations above, the symbol  $C^*$  stands for any character that does not belong to the list  $\{a, e, i, o, u, y\}$ , and  $C_{m_0}^*$  stands for a text string formed by consecutive appearance of zero or more (not necessarily identical or distinct) characters that do not belong to the same list.  $\square$

**Definition 6.23** (Latin protected range). Let  $\hat{\sigma}$  be a text string derived from a Latin word, its protected range  $\text{ProtRg}(\hat{\sigma})$  is an integer determined as follows:

- Try to find the string pattern  $(\emptyset|a(b|d|mb)|circum|co(m|\mu|n|rr)|de(f|l|t)|i|jnter|o(b)f|per|pr(ae|e|o)|re|su|tran)C_{m_0}^*V^*C_{m_0}^*$  in the string  $\hat{\sigma}$ ;
- If the string pattern above is found, the last position occupied by such a string defines  $\text{ProtRg}(\hat{\sigma})$ ; otherwise, set  $\text{ProtRg}(\hat{\sigma}) = 0$ .  $\square$

**Algorithm 6.24** (Latin effective spelling). For a Latin word  $\hat{\sigma}$ , its effective spelling  $\text{EffSpell}(\hat{\sigma})$  is constructed in sequential steps:

(1) *Do*<sup>81</sup>  $\bar{a} \rightarrow a$ ,  $\alpha \rightarrow ae$ ,  $\bar{e} \rightarrow e$ ,  $\bar{i} \rightarrow i$ ,  $\bar{o} \rightarrow o$ ,  $\alpha \rightarrow oe$ ,  $\bar{u} \rightarrow u$ ,  $\bar{y} \rightarrow y$ .

(2) *Do carolin~ → kapoliv, emi(s|ss|t|tt) → εμιcτ, impr~ → mpp, pass → pand, pater~ → patre, separ~ → σεπαρ, vac → βεac.*

<sup>80</sup>Nevertheless, we have included *eam*, *eo*, *is*, *iste* and *istis* in our list of Latin stop words, because they not only function as conjugated forms of the Latin verb for “to go”, but also as pronouns.

<sup>81</sup>Note that the italic forms of  $\alpha$  ( $\alpha$ ) and  $\alpha$  ( $\alpha$ ) are very similar.

## (3) Replace

$(\emptyset con)sobrin\mathbf{X}$	$(\emptyset cu)cur^{\mathbf{X}\epsilon}(r s)$	$(amit materter)^{\mathbf{X}\epsilon}(a e i o u y)\sim$				$(m p)atruel\mathbf{X}$			
$\kappa\sigma\tilde{v}$	$\kappa\mu\sigma\mathbf{X}$	$\alpha u\tilde{\nu}\tau\mathbf{X}$				$\kappa\mu\sigma\tilde{v}$			
<i>aestim~</i>	<i>ami(s tt)(\emptyset s)~</i>	<i>amiciss\mathbf{X}~</i>	<i>asp~</i>	<i>camer~</i>	<i>car(a i)o~</i>	<i>collin~</i>	<i>color~</i>	<i>delect~</i>	<i>dign~</i>
$\alpha\sigma\tau i\mu$	$\alpha\mu i\sigma$	<i>amica</i>	<i>aspa</i>	<i>ριμ</i>	<i>δεαρι</i>	<i>κολλι\tilde{v}</i>	<i>φαρβr</i>	<i>δλιγτ</i>	<i>διγn</i>
<i>divid~</i>	<i>divin~</i>	<i>er(ect ex ig)~</i>	<i>fem(ell in)~</i>	<i>haer~</i>	<i>here(d s)~</i>	<i>hestern~</i>	<i>honest~</i>	<i>horta~</i>	
$divi\delta$	$div\tilde{v}$	$\varepsilon\rho e\xi t$	$\phi\mu\lambda n$	$\eta a e r$	$\eta e i p$	$\eta e \sigma \tau$	$\eta \rho e \sigma \tau$	$\eta \rho \rho t a$	
<i>hr~</i>	<i>ingenu~</i>	<i>laet~</i>	<i>lisul~</i>	<i>locut</i>	<i>loq</i>	<i>lucas~</i>	<i>marit~</i>	<i>materi~</i>	<i>matron~</i>
$\mu\sigma\tau e r$	$in\gamma e\tilde{v}u$	$\chi\alpha\pi\pi t$	$\lambda i\zeta\zeta$	$l\lambda o q u o r$	$l\lambda o q$	$\lambda u\xi\sigma$	$\mu a \rho t$	$\mu a \tau \rho i$	$\lambda a \delta \eta$
<i>memin~</i>	<i>met~</i>	<i>mra~</i>	<i>nupt~</i>	<i>opin~</i>	<i>pact~</i>	<i>pag~</i>	<i>pet~</i>	<i>pig~</i>	<i>rheda~</i>
$\mu\tilde{v}i n$	$\phi e a \rho$	$\mu \rho \sigma a$	$\tilde{v}u \pi t$	$\omega p i n$	$p a n g$	$p a y$	$\pi e t$	$\pi i y$	$\omega o y \tilde{v}a$
<i>scri(b ps pt)~</i>	<i>senior~</i>	<i>sign~</i>	<i>soror~</i>	<i>tem~</i>	<i>van~</i>	<i>ven(a o)~</i>	<i>vene(m t)~</i>	<i>vidu~</i>	
$\omega \rho i t r$	$\sigma \gamma \tilde{v} o p$	$s i \gamma n$	$\sigma o \rho o p$	$\tau m e$	$\beta a \tilde{v}$	$v e \tilde{v} a$	$v e \tilde{v} a t$	$\omega i d o$	
<i>d(e(e i is o orum os um us))i(\emptyset i is s)</i>	<i>amic(i u)(\emptyset m s)</i>			<i>car(e um us)</i>		<i>die</i>	<i>heri</i>	<i>vener(\emptyset e)</i>	
$\gamma \rho \delta$		<i>amica</i>		$\delta e a \rho i$		<i>diebus</i>	$\eta e \sigma \tau$	<i>ve\tilde{v}atur</i>	

## (4) Replace

$\delta e a p i l$	$\hat{x} t r i (c x)$	$a(d f e r   l l a t   t t u l ) \sim$	$a b(l a t   t u l ) \sim$	$a m i c (i(a e m r t)e t o ) u t) \sim$
$k a r \lambda$	$\hat{x} t o r$	$a f f e r$	$a u f e r$	$a \mu \beta i c u$
$b i b \sim$	$c o(l l a t   n t u l ) \sim$	<i>complement</i>	<i>concor s</i>	<i>creat~</i>
$b i \beta$	<i>conf er</i>	<i>com pler</i>	<i>con cord</i>	<i>de a ~</i>
<i>domin(a)~</i>	<i>domina~</i>	<i>e(lat xtul)~</i>	<i>expos~</i>	<i>hort~</i>
$d o \mu i \tilde{v} o$	$d o m i \tilde{v} \alpha a$	<i>effer</i>	<i>expon</i>	$\gamma a \rho \delta t$
<i>memor~</i>	$mens^{\mathbf{X}\epsilon}(a e i o u y) \sim$	<i>ob(lat tul)~</i>	<i>oneros\mathbf{X}~</i>	<i>patrimon~</i>
$\mu e m$	$men\sigma\mathbf{X}$	<i>offer</i>	<i>oneris</i>	$\pi a t r i m o n$
<i>sci(am ar eb em en es et m o re ri sc t un v)\mathbf{X}</i>			<i>str</i>	<i>patru\hat{x}\#(m)\mathbf{X}~</i>
	<i>sci</i>		<i>starr</i>	<i>su(blat stul)~</i>
<i>truculent\mathbf{X}~</i>	<i>tutel~</i>	<i>us(i u)~</i>	<i>via~</i>	<i>terter</i>
$tr u c u m$	<i>tutar</i>	<i>uto</i>	<i>vi\alpha a</i>	$x$
			<i>x</i>	<i>\lambda g o o d \lambda</i>
$\mathbf{X}\epsilon(\emptyset ab ad amb circum co de in inter ob prod red suss trans)goLatin$				
<i>cor</i>		<i>d(ee ei eis eo eorum eos eum eus ia iae iam iarum ias ie ii iis io iorum ios is ium ius)</i>		
<i>cordis</i>			<i>divus</i>	
<i>lux</i>	<i>mens</i>	<i>mos</i>	<i>nix</i>	<i>ops</i>
$l u c u m$	<i>mentis</i>	<i>moris</i>	<i>nivis</i>	<i>opum</i>
				<i>rebus</i>
				<i>trah</i>
				<i>ult</i>
				<i>ire(m mus s)</i>
				<i>pater</i>
				<i>vis</i>
				<i>\sim emen(d\mathbf{X} s </i>
				<i>\sim ster</i>
				<i>~que</i>
				<i>sci</i>

(5) Do  $a^{\mathbf{X}\epsilon}(cc|dd|ff|gg|rr|ss|tt) \rightarrow a\mathbf{X}, x \rightarrow \check{x}$ .

## (6) Replace

$\hat{x}\notin(m)om$	$\hat{x}\epsilon(a e i o r u y)p(s t)$	$\mathbf{X}\epsilon(\hat{x}(a e i o r u y))r$	<i>cantor~</i>	<i>dele~</i>
$\hat{x}o\mu$	$\mathbf{X}p$	<i>Xerque</i>	<i>canere</i>	
<i>mpt</i>	<i>os\hat{x}</i>	<i>patr~</i>	<i>puls</i>	<i>tang~</i>
$m$	$o\sigma\hat{x}$	<i>tarrid</i>	<i>tact</i>	<i>tarriz</i>
		$tr\hat{x}\notin(a)~$	<i>ter\hat{x}</i>	<i>tract</i>
			<i>trah</i>	<i>ol</i>
				<i>ult</i>
				<i>ire(m mus s)</i>
				<i>pater</i>
				<i>vis</i>
				<i>\sim emen(d\mathbf{X} s </i>
				<i>\sim ster</i>
				<i>'que</i>
				<i>starr</i>

(7) Do  $\sim rs \rightarrow rt, \sim \hat{x}^{\epsilon}(\overline{a|e|i|o|r|t|u})er \rightarrow \hat{x}erquei, inter \sim \rightarrow jnter, super \sim \rightarrow suzp$ .(8) Remove apostrophe, and call the string obtained so far as  $\hat{\sigma}_{\sharp}$ .

- (9) Break down  $\hat{\sigma}_\sharp = \hat{\sigma}_1 \hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}_\sharp^{[\text{ProtRg}(\hat{\sigma}_\sharp)]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma}_\sharp)$  is equal to the protected range of  $\hat{\sigma}_\sharp$ .
- (10) Do  $\text{act} \rightarrow \text{ag}$ ,  $\text{erb} \rightarrow \text{erv}$ ,  $\text{eg} \sim \rightarrow \text{ag}$ ,  $\mathbf{X}^\epsilon(i|o|u)v \rightarrow \mathbf{X}t$ ,  $\text{ub} \rightarrow \text{uss}$  on  $\hat{\sigma}_1$  and call the result  $\hat{\sigma}'_1$ .
- (11) Do  $\sim\text{que} \rightarrow \text{'que}$  on  $\hat{\sigma}_2$  and call the result  $\hat{\sigma}'_2$ .
- (12) If the pattern  $\sim\mathbf{C}^*\mathbf{V}^*\mathbf{C}^*$  is found in  $\hat{\sigma}'_1$  and the two occurrences of  $\mathbf{C}^*$  represent the same letter, then remove the last letter from  $\hat{\sigma}'_1$ , and do  $\mathbf{V}^* \sim \rightarrow \emptyset$  on  $\hat{\sigma}'_2$ .<sup>82</sup> Call the results after these operations  $\hat{\sigma}''_1$  and  $\hat{\sigma}''_2$ , respectively.
- (13) On  $\hat{\sigma}''_2$ , perform the following substitutions in a sequel:
- (10.1) Do  $\hat{\chi}_\sharp(c|t)ul \rightarrow \hat{\chi}$ ,  $(cul|lent|mon) \rightarrow \emptyset$ ,  $(errim|issim) \rightarrow \emptyset$ ,  $\text{men}(\emptyset|t) \rightarrow \emptyset$ .
  - (10.2) Do  $\sim(m|ni|nt)(\emptyset|\mathbf{V}^*)(\emptyset|r)r\mathbf{V}^*|s \rightarrow \emptyset$ .
  - (10.3) Do  $\sim\hat{\chi}(mi|mu|mur|ntur|sti|ri) \rightarrow \hat{\chi}$ ,  $\hat{\chi}_\sharp(r)t \rightarrow \hat{\chi}$ .

The result after these three steps of operations is called  $\hat{\sigma}''_2$ .

- (14) Concatenate  $\hat{\sigma}''_1$  and  $\hat{\sigma}''_2$ .
- (15) Do  $\text{die}(bu|m|ru|s) \rightarrow \text{dieque}$ ,  $\text{ign} \rightarrow \text{en}$ ,  $\text{re}(bu|m|rum|s) \rightarrow \text{reque}$ ,  $ct \rightarrow \check{c}$ ,  $\text{pro}\sigma \sim \rightarrow \text{pros}$ ,  $\mathbf{V}^*\text{men}(\emptyset|t) \rightarrow \mathbf{V}^*$ ,  $\sim \rightarrow \emptyset$ .

**Algorithm 6.25** (Latin essential root). Let  $\hat{\sigma}$  be the effective spelling of a Latin word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:

- (1) Break down  $\hat{\sigma} = \hat{\sigma}_1 \hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}^{[\text{ProtRg}(\hat{\sigma})]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma})$  is equal to the protected range of  $\hat{\sigma}$ .
- (2) Do  $\mathbf{C}^*(a|ae)\mathbf{C}^* \rightarrow \mathbf{C}^*e\mathbf{C}^*$  on  $\hat{\sigma}_1$ ,<sup>83</sup> and call the result  $\hat{\sigma}'_1$ .

- (3) On  $\hat{\sigma}_2$ , perform the following substitutions in a sequel:

- (3.1) Do  $\sim\text{que} \rightarrow \emptyset$ .
- (3.2) Do  $(b|id|nd|t|v)\mathbf{V}^*(\emptyset|m|n|r|(s)_m|t) \rightarrow \emptyset$ .
- (3.3) Do  $\sim(a|e|i|in|o|r|ss|u)_m \rightarrow \emptyset$ .

The result after these three steps of operations is called  $\hat{\sigma}'_2$ .

- (4) Concatenate  $\hat{\sigma}'_1$  and  $\hat{\sigma}'_2$ .

- (5) Replace<sup>84</sup>

$sens \sim$	$\mathbf{X}^\epsilon(\emptyset a circum con di e\check{c} in jnter ob per prae re suzp)st(a ae ans ant e eb em end ens ent er es ess est et o)$	$\mathbf{X}_{\text{stand}}$
$sent$		
	$\mathbf{X}^\epsilon(\emptyset ab ad con in per prae pro red suss)d(a and ans ant e eb em end ens ent er erg es ess est et etis etur o or)$	$\mathbf{X}_{\text{datt}}$
	$\mathbf{X}^\epsilon(\emptyset ad af au circum con de dif in jnter(\emptyset queo) of prae pro re suf trans)fer(q r t)$	$\mathbf{X}_{\text{fer}}$
	$\mathbf{X}^\epsilon(\emptyset ad circum com de di e im jnter(\emptyset queo) o per prae pro re suss suzp trans)mis(\emptyset s)$	$\mathbf{X}^\epsilon(\emptyset cou prae pro)pass$
	$\mathbf{X}_{\text{mitt}}$	$\mathbf{X}_{\text{et}}$
$ci(e eam eant ear eas eat eb em end ens ent eo eor er es et t$	$es(\emptyset s t)$	$i(\check{c} c)\mathbf{V}^*\mathbf{C}_{m_0}$
$cit$	$ed$	$ic$
		$m(all el ev)$
		$noll$
		$vol$
$\sim(lev lit)$	$\simaperien$	$\simno\sigma c(\emptyset en)$
$lin$	$ap$	$not$
		$or$
$\sim p(ast ev)$	$\simpel$	$\simvinc(\emptyset en)$
$pasc$	$pell$	$vic$
	$\mathbf{X}^\epsilon(\emptyset ab abs al at circum con de dis in jnter(\emptyset queo) ob prae pro re sub trans)(let tul)$	
		$\mathbf{X}'_{\text{fer}}$

<sup>82</sup>This peculiar step is tailored for reduplications in the perfect forms of certain Latin verbs. Such reduplications are not attested in most modern Indo-European languages, with Greek being a notable exception.

<sup>83</sup>Here, the two occurrences of  $\mathbf{C}^*$  may or may not represent the same letter.

<sup>84</sup>All the patches below are devoted to some common verbs that exhibit highly irregular conjugations.

where in the last step, one constructs  $\mathbf{X}'$  from  $\mathbf{X}$  by doing  $\sim ab(\emptyset|s) \rightarrow au$ ,  $\sim (\check{c}|l|t) \rightarrow f$ ,  $dis \rightarrow dif$ ,  $ob \rightarrow of$ ,  $suss \rightarrow suf$ .

(6) Set

$$\mathbf{X}_1 = (eam|eant|east|eat|eo|eund|eunt|i|ib|ibant|ibem|ibes|ibet|ibim|ibis|ibit|ibo|ibunt|iens|ieq|ier|i|iim|iit|im|imus|ir|ire|irent|ireq|iret|iri|is|iss|isse|issem|issent|isser|ist|isti|istis|ite|iti|itis|itist|ito|itu|itum|itur|itus),$$

and

$$\mathbf{X}_2 = (er|es|less|est|for|fuer|fui|fuiim|fuiiss|fuisst|fuit|fut|si|sim|sint|sit|sum|sumus|sunt|sunto).$$

Replace<sup>85</sup>

$$\begin{array}{c} \mathbf{X}^\epsilon(\emptyset|ab|ad|amb|circum|co|de|in|jnter(\emptyset|queo)|ob|prod|red|suss|trans)\mathbf{X}_1 \\ \mathbf{X}_{eo} \\ \mathbf{X}^\epsilon(ab|ad|de|in|jnter|ob|prae|pro|suss|suzp)\mathbf{X}_2 \\ \mathbf{X}_{sum} \end{array}$$

### 6.3.2 Admissible mutation and approximate clustering

Like the case of French, vowel blotting is unnecessary for Latin. The only significant pattern of vowel alternation in Latin verb conjugations has already been taken care of by Algorithm 6.25(2).

In what follows, we will construct a bivariate Boolean-valued function  $HrdTest(\hat{\alpha}, \hat{\beta})$  on a “simple heredity test function” in Algorithm 6.26, and a set of “admissible suffix mismatch” rules in Algorithm 6.27.

**Algorithm 6.26** (Simple heredity test). *Let  $\hat{\alpha}'$  be the result from doing  $\sim \mathbf{V}_{m_0}^*(\emptyset|erquei|que|s|st) \rightarrow \emptyset$  on  $\hat{\alpha}$ , and define  $\hat{\beta}'$  similarly. The Boolean-valued function  $SimpHrdTest(\hat{\alpha}, \hat{\beta})$  returns TRUE if the lowercase form of  $\hat{\alpha}$  contains at least one instance of  $\mathbf{V}^*$  (Definition 6.22) AND at least one of the following three conditions holds:<sup>86</sup>*

- (i)  $\hat{\alpha}' = \hat{\beta}'$ ;
- (ii)  $\hat{\beta} = \hat{\alpha}t$ ;
- (iii)  $\ell(\hat{\beta}) > \ell(\hat{\alpha}) \geq \frac{\ell(\hat{\beta})}{2}$  AND  $\hat{\alpha} = \hat{\beta}^{[\ell(\hat{\alpha})]}$  AND  $\hat{\beta}^{[\ell(\hat{\alpha}')+1]} = \mathbf{V}^*$ . (See Definition 3.1 for the notations  $\hat{\beta}^{[n]}$  and  $\hat{\beta}^{\{n\}}$ .)

In what follows, we define  $SuffixNW(\hat{\alpha}, \hat{\beta})$ ,  $RootNW(\hat{\alpha}, \hat{\beta})$ ,  $NW^*(\hat{\alpha}, \hat{\beta})$  and  $SuffixSW(\hat{\alpha}, \hat{\beta})$ ,  $SuffixSW(\hat{\alpha}, \hat{\beta})$ ,  $SW^*(\hat{\alpha}, \hat{\beta})$  as what is done in the Danish case (Algorithm 5.7), which was also applied to two other Germanic languages treated in §5.

**Algorithm 6.27** (Admissible suffix mismatch). *For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function*

$$AdmSM(RootNW(\hat{\alpha}, \hat{\beta}), SuffixNW(\hat{\alpha}, \hat{\beta}), NW^*(\hat{\alpha}, \hat{\beta}))$$

*returns TRUE if  $NW^*(\hat{\alpha}, \hat{\beta}) = \emptyset$  AND  $RootNW(\hat{\alpha}, \hat{\beta})$  contains at least one instance of  $\mathbf{V}^*$  AND at least one of the following four conditions holds:*

- (i)  $SuffixNW(\hat{\alpha}, \hat{\beta}) = [(\emptyset|(a|e|i|o|u)_m(\emptyset|n)), (\emptyset|(a|e|i|o|u)_m(\emptyset|n))]$ ;
- (ii)  $SuffixNW(\hat{\alpha}, \hat{\beta}) = [c, \check{c}][\check{c}, g][c, q][d, s][den, s]$ ;
- (iii)  $SuffixNW(\hat{\alpha}, \hat{\beta}) = [d, t]$  AND  $\mathcal{Q}(RootNW(\hat{\alpha}, \hat{\beta})) = r$ ;
- (iv)  $SuffixNW(\hat{\alpha}, \hat{\beta}) = [\check{c}, c \mathbf{V}^* \mathbf{X}]$ .

Similarly, one can evaluate another Boolean-valued function

$$AdmSM(RootSW(\hat{\alpha}, \hat{\beta}), SuffixSW(\hat{\alpha}, \hat{\beta}), SW^*(\hat{\alpha}, \hat{\beta}))$$

by trading all the occurrences of NW in the statements above with SW.

**Algorithm 6.28** (Heredity test function). *The structure of the Latin heredity test function  $HrdTest(\hat{\alpha}, \hat{\beta})$  is identical to the German version (Algorithm 8.1.2), except that the functions  $SimpHrdTest$ ,  $RootNW$ ,  $SuffixNW$ ,  $NW^*$ ,  $RootSW$ ,  $SuffixSW$ ,  $SW^*$  must follow the Latin rules stated above.*

<sup>85</sup>All the patches below are devoted to suppletive verbs derived from *eō* “go” and *sum* “be”.

<sup>86</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical OR.

**Algorithm 6.29** (Approximate clustering of Latin words). *The algorithm is essentially the same as Algorithm 6.20, except that Latin rules (instead of French rules) apply to all the tags (effective spelling, essential root etc.).*

*Example 6.29.1.* There are five declensions in Latin. The first three declensions cover both nouns and adjectives, while the last two are reserved for nouns only. We extract sample nouns and adjectives from the following Wiktionary and Wikipedia pages:

[https://en.wiktionary.org/wiki/Appendix:Latin\\_first\\_declension](https://en.wiktionary.org/wiki/Appendix:Latin_first_declension)

*comētae, comētārum, comētās, comētē, comētem, comētēn, comētēs, comētīs* — “comet”;  
*dynastae, dynastārum, dynastās, dynastē, dynastēn, dynastēs, dynastīs* — “ruler”;  
*nauta, nautā, nautae, nautam, nautārum, nautās, nautīs* — “sailor”;  
*nymphae, nymphārum, nymphās, nymphē, nymphēn, nymphēs, nymphīs* — “bride”;  
*Rōma, Rōmā, Rōmae, Rōmam* — “Rome”;  
*rosa, rosā, rosae, rosam, rosārum, rosās, rosīs* — “rose”;  
*stēlla, stellā, stellae, stēllam, stēllārum, stellās, stellīs* — “star”;  
*xiphiā, xiphiae, xiphiān, xiphiārum, xiphiās, xiphiīs* — “swordfish”.

[https://en.wiktionary.org/wiki/Appendix:Latin\\_second\\_declension](https://en.wiktionary.org/wiki/Appendix:Latin_second_declension)

*ager, agrī, agrīs, agrō, agrōrum, agrōs, agrum* — “field”;  
*ampele, ampelī, ampelīs, ampelō, ampelon, ampelōrum, ampelos, ampelōs, ampelum* — “vine”;  
*atome, atomī, atomīs, atomō, atomōrum, atomōs, atomum, atomus* — “atom”;  
*bella, bellī, bellīs, bellō, bellōrum, bellum* — “war”;  
*fīlī, fīliī, fīliīs, fīliō, fīliōrum, fīliōs, fīliūm, fīlius* — “son”;  
*magister, magistrī, magistrīs, magistrō, magistrōrum, magistrōs, magistrum* — “teacher”;  
*mūre, mūrī, mūrīs, mūrō, mūrōrum, mūrōs, mūrum, mūrus* — “wall”;  
*mȳthe, mȳthī, mȳthīs, mȳthō, mȳthon, mȳthōrum, mȳthōs, mȳthōs, mȳthum* — “myth”;  
*phaenomena, phaenomenī, phaenomenīs, phaenomenō, phaenomenon, phaenomenōrum* — “phenomenon”;  
*puer, puerī, puerīs, puerō, puerōrum, puerōs, puerum* — “child”;  
*templā, templī, templīs, templō, templōrum, templūm* — “temple”.

[https://en.wiktionary.org/wiki/Appendix:Latin\\_third\\_declension](https://en.wiktionary.org/wiki/Appendix:Latin_third_declension)

*āer, āera, āere, āerem, āerēs, āerī, āeribus, āeris, āeros, āerum* — “air”;  
*animal, animālī, animālia, animālibus, animālis, animālium* — “animal”;  
*base, basem, basēs, basī, basibus, basim, basis, basīs, basium* — “pedestal”;  
*haerese, haeresem, haeresēs, haeresī, haeresibus, haeresim, haeresis, haeresīs, haeresium* — “sect”;  
*homine, hominem, hominēs, hominī, hominibus, hominis, hominum, homō* — “human”;  
*nocte, noctem, noctēs, noctī, noctibus, noctis, noctium, nox* — “night”;  
*nōmen, nōmina, nōmine, nōminī, nōminibus, nōminis, nōminum* — “name”;

*tigre, tigrem, tigrēs, tigrī, tigribus, tigride, tigridem, tigridēs, tigridī, tigridibus, tigridis, tigridum, tigrim, tigris, tigrīs, tīgris, tigrum* — “tiger”;

*ture, turrem, turrēs, turri, turribus, turrim, turris, turri, turrium* — “tower”.

[https://en.wiktionary.org/wiki/Appendix:Latin\\_fourth\\_declension](https://en.wiktionary.org/wiki/Appendix:Latin_fourth_declension)

*cornibus, cornū, cornua, cornūs, cornuum* — “horn”;

*Dīdō, Dīdōne, Dīdōnem, Dīdōnēs, Dīdōnī, Dīdōnibus, Dīdōnis, Dīdōnum* — “Dido”;

*ēchibus, ēcho, ēchū, ēchuī, ēchum, ēchūs, ēchuum* — “echo”;

*manibus, manū, manuī, manum, manus, manūs, manuum* — “hand”.

[https://en.wiktionary.org/wiki/Appendix:Latin\\_fifth\\_declension](https://en.wiktionary.org/wiki/Appendix:Latin_fifth_declension)

*diē, diēbus, diētī, diem, diērum, diēs* — “day”;

*fidē, fidēbus, fideī, fidēm, fidērum, fidēs* — “faith”;

*rē, rēbus, reī, rem, rērum, rēs* — “thing”;

*speciē, speciēbus, speciētī, speciem, speciērum, speciēs* — “view”.

[https://en.wikipedia.org/wiki/Latin\\_declension#Adjectives](https://en.wikipedia.org/wiki/Latin_declension#Adjectives)

Here are some Latin adjectives whose positive forms belong to the first and second declensions. (We also incorporate, in our list below, comparative and superlative forms of certain sample adjectives, which may or may not share the same declension pattern with the original adjectives in positive forms.)

*alta, altā, altae, altam, altārum, altās, alte, altī, altior, altiōra, altiōrem, altiōrēs, altiōrī, altiōribus, altiōris, altiōrum, altīs, altissima, altissimā, altissimae, altissimam, altissimārum, altissimās, altissime, altissimī, altissimīs, altissimō, altissimōrum, altissimōs, altissimum, altissimus, altius, altō, altōrum, altōs, altum, altus* — “tall”;

*atoma, atomā, atomae, atomam, atomārum, atomās, atome, atomī, atomīs, atomō, atomōrum, atomōs, atomum, atomus* — “indivisible”;

*miser, misera, miserā, miserae, miseram, miserārum, miserās, miserī, miserō, miserōra, miserōre, miserōrem, miserōrēs, miserōrī, miserōribus, miserōris, miserōrum, miserīs, miserius, miserō, miserōrum, miserōs, miserrima, miserrimā, miserrimae, miserrimam, miserrimārum, miserrimās, miserrime, miserrimī, miserrimīs, miserrimō, miserrimōrum, miserrimōs, miserrimum, miserrimus, miserum* — “poor”;

*sacer, sacra, sacrā, sacrae, sacram, sacrārum, sacrās, sacrī, sacrīs, sacrō, sacrōrum, sacrōs, sacram* — “sacred”;

*ūlla, ūllā, ūllae, ūllam, ūllārum, ūllās, ūlle, ūllī, ūllīs, ūllītī, ūllō, ūllōrum, ūllōs, ūllum, ūllus* — “any”.

Here are some Latin adjectives whose positive forms belong to the third declension.

*agile, agilem, agilēs, agilī, agilia, agilibus, agilis, agilium* — “agile”;

*alacer, alacre, alacrem, alacrēs, alacrī, alacia, alacribus, alacris, alacrium* — “lively”;

*atrōcem, atrōcēs, atrōcī, atrōcia, atrōcibus, atrōcior, atrōciōra, atrōciōre, atrōciōrem, atrōciōrēs, atrōciōrī, atrōciōribus, atrōciōris, atrōciōrum, atrōcis, atrōcissima, atrōcissimā, atrōcissimae, atrōcissimam, atrōcissimārum, atrōcissimās, atrōcissime, atrōcissimī, atrōcissimīs, atrōcissimō, atrōcissimōrum, atrōcissimōs, atrōcissimum, atrōcissimus, atrōcium, atrōcius, atrōx* — “fierce”;

*celer, celere, celerem, celerēs, celerī, celeria, celeribus, celerior, celerīra, celerīre, celerīrem, celerīrēs, celerīrī, celerīribus, celerīris, celerīrum, celeris, celerium, celerius, celerrima, celerrimā, celerrimae, celerrimam, celerrimārum, celerrimās, celerrime, celerrītī, celerrimīs, celerrimō, celerrimōrum, celerrimōs, celerrimūm, celerrimus* — “fast”;

*meliōr, meliōra, meliōre, meliōrem, meliōrēs, meliōrī, meliōribus, meliōris, meliōrum, melius* — “better”;

*trīste, trīstem, trīstēs, trīstī, trīstia, trīstibus, trīstior, trīstiōra, trīstiōre, trīstiōrem, trīstiōrēs, trīstiōrī, trīstiōribus, trīstiōris, trīstiōrum, trītis, trītissima, trītissimā, trītissimae, trītissimam, trītissimārum, trītissimās, trītissime, trītissimī, trītissimīs, trītissimō, trītissimōrum, trītissimōs, trītissimum, trītissimus, trītium, trītius* — “unhappy”;

*vetera, vetere, veterem, veterēs, veterī, veteribus, veteris, veterrima, veterrimā, veterrimae, veterrimam, veterrimārum, veterrimās, veterrime, veterrimī, veterrimīs, veterrimō, veterrimōrum, veterrimōs, veterrimum, veterrimus, veterum, vetus, vetustior, vetustiōra, vetustiōre, vetustiōrem, vetustiōrēs, vetustiōrī, vetustiōribus, vetustiōris, vetustiōrum, vetustissima, vetustissimā, vetustissimae, vetustissimam, vetustissimārum, vetustissimās, vetustissime, vetustissimī, vetustissimīs, vetustissimō, vetustissimōrum, vetustissimōs, vetustissimum, vetustissimus, vetustius* — “old”.

Our clustering algorithm yields the following result:

{*āer, āera, āere, āerem, āerēs, āerī, āeribus, āeris, āeros, āerum*},

{*ager, agrī, agrīs, agrō, agrōrum, agrōs, agrum*},

{*agile, agilem, agilēs, agilī, agilia, agilibus, agilis, agilium*},

{*alacer, alacre, alacrem, alacrēs, alacrī, alacia, alacribus, alacris, alacrium*},

{*alta, altā, altae, altam, altārum, altās, alte, altī, altior, altiōra, altiōre, altiōrem, altiōrēs, altiōrī, altiōribus, altiōris, altiōrum, altīs, altissima, altissimā, altissimae, altissimam, altissimārum, altissimās, altissime, altissimī, altissimīs, altissimō, altissimōrum, altissimōs, altissimum, altissimus, altius, altō, altōrum, altōs, altum, altus*},

{*ampele, ampelī, ampelīs, ampelō, ampelon, ampelōrum, ampelos, ampelōs, ampelum*},

{*animal, animālī, animālia, animālibus, animālis, animālum*},

{*atoma, atomā, atomae, atomam, atomārum, atomās, atome, atomī, atomīs, atomō, atomōrum, atomōs, atomum, atomus*},

{*atrōcem, atrōcēs, atrōcī, atrōcia, atrōcibus, atrōcior, atrōciōra, atrōciōre, atrōciōrem, atrōciōrēs, atrōciōrī, atrōciōribus, atrōciōris, atrōciōrum, atrōcis, atrōcissima, atrōcissimā, atrōcissimae, atrōcissimam, atrōcissimārum, atrōcissimās, atrōcissime, atrōcissimī, atrōcissimīs, atrōcissimō, atrōcissimōrum, atrōcissimōs, atrōcissimūm, atrōcissimus, atrōcium, atrōcius, atrōx*},

{*base, basem, basēs, basī, basibus, basim, basis, basīs, basium*},

{*bella, bellī, bellīs, bellō, bellōrum, bellum*},

{*celer, celere, celerem, celerēs, celerī, celeria, celeribus, celerior, celeriōra, celeriōre, celeriōrem, celeriōrēs, celeriōrī, celeriōribus, celeriōris, celeriōrum, celeris, celerium, celerius, celerrima, celerrimā, celerrima, celerrimam, celerrimārum, celerrimās, celerrime, celerrimī, celerrimīs, celerrimō, celerrimōrum, celerrimōs, celerrimum, celerrimus*},

{*comētae, comētarum, comētās, comētē, comētem, comētēn, comētēs, comētīs*},

{*cornibus, cornū, cornua, cornūs, cornuum*},

{*Dīdō, Dīdōne, Dīdōnem, Dīdōnēs, Dīdōnī, Dīdōnibus, Dīdōnis, Dīdōnum, diē, diēbus, diētī, diem, diērum, diēs*},

{*dynastae, dynastārum, dynastās, dynastē, dynastēn, dynastēs, dynastīs*},

{*ēchibus, ēcho, ēchū, ēchuī, ēchum, ēchūs, ēchuum*},

{*fidē, fidēbus, fidētī, fidēm, fidērum, fidēs*},

{*fīlī, fīliī, fīliīs, fīliō, fīliōrum, fīliōs, fīlium, fīlius*},

{haerese, haerem, haeresēs, haerest̄, haeresibus, haeresim, haeresis, haeresīs, haeresium},  
 {hominē, hominem, hominēs, hominīt̄, hominibus, hominis, hominum, homō},  
 {magister, magistrīt̄, magistrīs, magistrō, magistrōrum, magistrōs, magistrum},  
 {manibus, manū, manuū, manum, manus, manūs, manuum},  
 {melior, melius},  
 {meliōra, meliōre, meliōrem, meliōrēs, meliōrīt̄, meliōribus, meliōris, meliōrum},  
 {misēr, misera, miserā, miserae, miseram, miserārum, miserās, miserīt̄, miserīor, miserīora, miserīore, miserīōrem, miserīōrēs, miserīōrīt̄, miserīōribus, miserīōris, miserīōrum, miserīs, miserīus, miserō, miserōrum, miserōs, miserrima, miserrimā, miserrimae, miserrimam, miserrimārum, miserrimās, miserrime, miserrimīt̄, miserrimīs, miserrimō, miserrimōrum, miserrimōs, miserrimum, miserrimus, miserum},  
 {mūre, mūrīt̄, mūrīs, mūrō, mūrōrum, mūrōs, mūrum, mūrus},  
 {mȳthe, mȳthīt̄, mȳthīs, mȳthō, mȳthon, mȳthōrum, mȳthōs, mȳthōs, mȳthum},  
 {nautā, nautā, nautae, nautam, nautārum, nautās, nautīs},  
 {nocte, noctem, noctēs, noctīt̄, noctibus, noctis, noctium, nox},  
 {nōmen, nōmina, nōmine, nōminīt̄, nōminibus, nōminis, nōminum},  
 {nymphae, nymphārum, nymphās, nymphē, nymphēn, nymphēs, nymphīs},  
 {phaenomena, phaenomenīt̄, phaenomenīs, phaenomenō, phaenomenon, phaenomenōrum},  
 {puer, puerīt̄, puerīs, puerō, puerōrum, puerōs, puerum},  
 {rē, rēbus, reīt̄, rem, rērum, rēs},  
 {Rōma, Rōmā, Rōmae, Rōmam},  
 {rosā, rosā, rosae, rosam, rosārum, rosās, rosīs},  
 {sacer, sacra, sacrā, sacrae, sacram, sacrārum, sacrās, sacrīt̄, sacrīs, sacrō, sacrōrum, sacrōs, sacrum},  
 {speciē, speciēbus, speciēt̄, speciem, speciērum, speciēs},  
 {stēlla, stēllā, stēllae, stēllam, stēllārum, stēllās, stēllīs},  
 {templa, templīt̄, templīs, templō, templōrum, templum},  
 {tigre, tigrem, tigrēs, tigrīt̄, tigrībus, tigrīde, tigrīdem, tigrīdēs, tigrīdīt̄, tigrīdībus, tigrīdis, tigrīdūm, tigrīm, tigrīs, tigrīs, tigrīm},  
 {trīste, trīstem, trīstēs, trīstīt̄, trīstia, trīstibus, trīstīor, trīstīora, trīstīore, trīstīōrem, trīstīōrēs, trīstīōrīt̄, trīstīōribus, trīstīōrīs, trīstīōrum, trīstīs, trīstīssimā, trīstīssimā, trīstīssimāe, trīstīssimām, trīstīssimārum, trīstīssimās, trīstīssimē, trīstīssimīt̄, trīstīssimīs, trīstīssimō, trīstīssimōrum, trīstīssimōs, trīstīssimō, trīstīssimūs, trīstīssimūs, trīstīssimūs},  
 {ture, turre, turrem, turrēs, turrīt̄, turribus, turrim, turris, turrīs, turrium},  
 {ūlla, ūllā, ūllae, ūllam, ūllārum, ūllās, ūlle, ūllīt̄, ūllīs, ūllīt̄, ūllō, ūllōrum, ūllōs, ūllum, ūllus},  
 {vetera, veterē, veterem, veterēs, veterīt̄, veteribus, veteris, veterīma, veterīmā, veterīmāe, veterīmām, veterīmārum, veterīmās, veterīmē, veterīmīt̄, veterīmīs, veterīmō, veterīmōrum, veterīmōs, veterīmūm, veterīmūs, veterū, veterūm, vetus, vetustīor, vetustīora, vetustīore, vetustīōrem, vetustīōrēs, vetustīōrīt̄, vetustīōribus, vetustīōrīs, vetustīōrum, vetustīssimā, vetustīssimā, vetustīssimāe, vetustīssimām, vetustīssimārum, vetustīssimās, vetustīssimē, vetustīssimīt̄, vetustīssimīs, vetustīssimō, vetustīssimōrum, vetustīssimōs, vetustīssimūs, vetustīssimūs, vetustīssimūs}

{*xiphiā*, *xiphiae*, *xiphiān*, *xiphiārum*, *xiphiās*, *xiphiīs*}.

*Example 6.29.2.* Latin verbs are divided into four conjugation groups. We test our algorithm on the following selections of verbs, based primarily on the Wikipedia links below:

[https://en.wikipedia.org/wiki/Latin\\_conjugation#First\\_conjugation](https://en.wikipedia.org/wiki/Latin_conjugation#First_conjugation)

*amā, amābam, amābāminī, amābāmur, amābāmus, amābant, amābantur, amābar, amābare, amābāris, amābās, amābat, amābātis, amābātūr, amābere, amāberis, amābimīnī, amābimur, amābimus, amābis, amābit, amābitis, amābitur, amābō, amābor, amābunt, amābuntur, amāminī, amāmur, amāmus, amandī, amandō, amandum, amandus, amāns, amant, amantō, amantor, amantur, amāre, amārem, amārēminī, amārēmur, amārēmus, amārent, amārentur, amārer, amārēre, amārēris, amārēs, amāret, amārētis, amārētur, amārī, amāris, amās, amat, amāta, amātā, amātae, amātam, amātārum, amātās, amāte, amātī, amātis, amātīs, amātō, amātor, amātōrum, amātōs, amātōte, amātū, amātum, amātūr, amātūs, amāveram, amāverāmus, amāverant, amāverās, amāverat, amāverātis, amāvēre, amāverim, amāverīmus, amāverint, amāveris, amāverīs, amāverit, amāveritis, amāverītis, amāverō, amāvērunt, amāvī, amāvīmus, amāvīsse, amāvīssēmus, amāvīssent, amāvīssēs, amāvīsset, amāvīssētis, amāvīstī, amāvīstis, amāvit, amem, amēminī, amēmur, amēmus, ament, amentur, amer, amēre, amēris, amēs, amet, amētis, amētūr, amō, amor — “love”;*

*dā, dabam, dabāminī, dabāmur, dabāmus, dabant, dabantur, dabar, dabāre, dabāris, dabās, dabat, dabātis, dabātūr, dabere, daberis, dabimīnī, dabimur, dabimus, dabis, dabit, dabitis, dabitur, dabō, dabor, dabunt, dabuntur, damīnī, damur, damus, dandī, dandō, dandum, dandus, dāns, dant, dantō, dantor, dantur, dare, darem, darēminī, darēmur, darēmus, darent, darentur, darer, darēre, darēris, darēs, daret, darētis, darētur, dari, das, dat, data, datā, datae, datam, datārum, datās, date, datī, datis, datīs, datō, dator, datōrum, datōs, datōtē, datū, datum, datur, datūr, datus, dederam, dederāmus, dederant, dederat, dederātis, dedēre, dederim, dederimus, dederīmus, dederint, dederis, dederīs, dederit, dederitis, dederītis, dederō, dedērunt, dedī, dedimus, dedisse, dedissēmus, dedissent, dedissēs, dedisset, dedissētis, dedistī, dedistis, dedit, dem, dēminī, dēmur, dēmus, dent, dentur, der, dēre, dēris, dēs, det, dētis, dētūr, dō, dor — “give”;*

*fricā, fricābam, fricābāminī, fricābāmur, fricābāmus, fricābant, fricābantur, fricābar, fricābāre, fricābāris, fricābās, fricābat, fricābātis, fricābātūr, fricābere, fricāberis, fricābimīnī, fricābimur, fricābimus, fricābis, fricābit, fricābitis, fricābitur, fricābō, fricābor, fricābunt, fricābuntur, fricāminī, fricāmur, fricāmus, fricāndī, fricāndō, fricāndum, fricāndus, fricāns, fricānt, fricāntō, fricāntor, fricāntur, fricāre, fricārem, fricārēminī, fricārēmur, fricārēmus, fricārent, fricārentur, fricārer, fricārēre, fricārēs, fricāret, fricārētis, fricārētur, fricārī, fricāris, fricās, fricāt, fricātē, fricātibūs, fricātīs, fricātō, fricātōr, fricātōte, fricātū, fricātuī, fricātūm, fricātūr, fricātūr, fricātūs, fricātūm, fricātūum, fricēm, fricēminī, fricēmur, fricēmus, fricēnt, fricēntur, fricēr, fricēre, fricēris, fricēs, fricēt, fricētē, fricō, fricōr, fricueram, fricuerāmus, fricuerant, fricuerās, fricuerat, fricuerātis, fricuerē, fricuerim, fricuerīmus, fricuerint, fricueris, fricuerīs, fricuerit, fricuerītis, fricuerō, fricuerūt, fricūt, fricūmus, fricūsse, fricūssēmus, fricūssēs, fricūssent, fricūssētis, fricūstī, fricūstis, fricuit — “rub”;*

*iūta, iūtā, iūtae, iūtam, iūtārum, iūtās, iūte, iūtī, iūtō, iūtōrum, iūtū, iūtūm, iūtūr, iūtūs, iuvā, iuvābam, iuvābāminī, iuvābāmur, iuvābāmus, iuvābant, iuvābantur, iuvābar, iuvābāre, iuvābāris, iuvābās, iuvābat, iuvābātis, iuvābātūr, iuvābere, iuvāberis, iuvābimīnī, iuvābimur, iuvābimus, iuvābis, iuvābit, iuvābitis, iuvābitur, iuvābō, iuvābor, iuvābunt, iuvābuntur, iuvāminī, iuvāmur, iuvāmus, iuvāndī, iuvāndō, iuvāndum, iuvāndus, iuvāns, iuvānt, iuvāntō, iuvāntor, iuvāntur, iuvāre, iuvārem, iuvārēminī, iuvārēmur, iuvārēmus, iuvārent, iuvārentur, iuvārer, iuvārēre, iuvārēris, iuvārēs, iuvāret, iuvārētis, iuvārētur, iuvārī, iuvāris, iuvās, iuvāt, iuvātē, iuvātis, iuvātō, iuvātōr, iuvātōte, iuvātūr, iuvāt, iuvēmīnī, iuvēmur, iuvēmus, iuvēt, iuvēnt, iuvēntur, iuvēr, iūveram, iūverāmus, iūverās, iūverat, iūverātis, iūverātūr, iuvēr, iūverim, iūverīmus, iūverint, iuvēris, iūverīs, iūverit, iūverītis, iūverō, iūverūt, iuvēs, iuvēt, iuvētūr, iūvī, iūvīmīnī, iuvīsse, iuvīssēmus, iuvīssēs, iuvīssent, iuvīssētis, iūvīstī, iūvīstis, iūvīt, iuvō, iuvor — “help”;*

*lavā, lavābam, lavābāminī, lavābāmur, lavābāmus, lavābant, lavābantur, lavābar, lavābare, lavābāris, lavābās, lavābat, lavābātis, lavābātūr, lavābere, lavāberis, lavābimīnī, lavābimur, lavābimus, lavābis, lavābit, lavābitis, lavābitur, lavābō, lavābor, lavābunt, lavābuntur, lavāminī, lavāmur, lavāmus, lavāndī, lavāndō, lavāndum, lavāndus, lavāns, lavānt, lavāntō, lavāntor, lavāntur, lavāre, lavārem, lavārēminī, lavārēmur, lavārēmus, lavārent, lavārentur, lavārer, lavārēre, lavārēris, lavārēs, lavāret, lavārētis, lavārētur, lavārī, lavāris, lavās, lavāt, lavātā, lavātāe, lavātam, lavātārum, lavātās, lavāte, lavātī, lavātīs, lavātō, lavātor, lavātōrum, lavātōs, lavātōte, lavātū, lavātūm, lavātūr, lavātūs, lavātūs, lavēmīnī, lavēmur, lavēmus, lavēnt, lavēntur, laver, lavēre, lavēris, lavēs, lavet, lavētis, lavētūr, lavō, lavor — “wash”;*

*mīrābāminī, mīrābāmur, mīrābantur, mīrābar, mīrābāre, mīrābāris, mīrābātur, mīrābere, mīrāberis, mīrābimīnī, mīrābimur, mīrābitur, mīrābor, mīrābuntur, mīrāminī, mīrāmur, mīrandī, mīrandō, mīrandum, mīrandus, mīrāns, mīrantor, mīrantur, mīrāre, mīrārēminī, mīrārēmur, mīrārentur, mīrārer, mīrārēre, mīrārēris, mīrārētur, mīrārī, mīrāris, mīrāta, mīrātā, mīrātae, mīrātam, mīrātārum, mīrātās, mīrāte, mīrātī, mīrātīs, mīrātō, mīrātor, mīrātōrum, mīrātōs, mīrātū, mīrātum, mīrātūrus, mīrātus, mīrēminī, mīrēmur, mīrentur, mīrer, mīrēre, mīrēris, mīrētur, mīror — “admire”;*

*portā, portābam, portābāminī, portābāmur, portābāmus, portābāntur, portābar, portābāre, portābāris, portābās, portābat, portābātis, portābātur, portābere, portāberis, portābiminī, portābimur, portābimus, portābis, portābit, portābitis, portābitur, portābō, portābor, portābuntur, portāminī, portāmur, portāmus, portandī, portandō, portandum, portandus, portāns, portant, portantō, portantor, portantur, portārem, portārēminī, portārēmur, portārēmus, portārent, portārentur, portārer, portārēre, portārēris, portārēs, portāret, portārētis, portārētur, portārī, portāris, portās, portat, portātā, portātae, portātam, portātārum, portātās, portātē, portātī, portātīs, portātō, portātor, portātōrum, portātōs, portātōte, portātū, portātum, portātūr, portātūrus, portātūs, portāveram, portāverāmus, portāverant, portāverās, portāverat, portāverātis, portāvēre, portāverim, portāverimus, portāverīmus, portāverint, portāveris, portāverīs, portāverit, portāveritis, portāverītis, portāverō, portāvērunt, portāvī, portāvīmus, portāvisse, portāvissem, portāvissēmus, portāvissent, portāvissēs, portāvissset, portāvissētis, portāvistī, portāvistis, portāvit, portem, portēminī, portēmur, portēmus, portent, portentur, porter, portēre, portēris, portēs, portet, portētis, portētur, portō, portor — “convey”;*

*secā, secābam, secābāminī, secābāmur, secābāmus, secābāntur, secābar, secābāre, secābāris, secābās, secābat, secābātis, secābātur, secābere, secāberis, secābiminī, secābimur, secābimus, secābīs, secābit, secābitis, secābitur, secābō, secābor, secābuntur, secāminī, secāmur, secāmus, secandī, secandō, secandum, secandus, secāns, secant, secantō, secantor, secantur, secāre, secārem, secārēminī, secārēmur, secārēmus, secārent, secārentur, secārer, secārēre, secārēris, secārēs, secāret, secārētis, secārētur, secārī, secāris, secās, secat, secāte, secātis, secātō, secātor, secātōte, secātūr, secēt, secēmīnī, secēmur, secēmus, secēnt, secēntur, secēr, secēre, secēris, secēs, secet, secētis, secētūr, secō, secor, secta, sectā, sectae, sectam, sectārum, sectās, secte, sectī, sectīs, sectō, sectōrum, sectōs, sectū, sectum, sectūrus, sectus, secueram, secuerāmus, secuerant, secuerās, secuerat, secuerātis, secuēre, secuerim, secuerimus, secuerīmus, secuerint, secueris, secuerīs, secuerit, secueritis, secuerītis, secuerō, secuērunt, secuī, secuīmus, secuisse, secuisse, secuissēmus, secuissent, secuissēs, secuisset, secuissētis, secuistī, secuistis, secuit — “cut”;*

*stā, stābam, stābāmus, stābāntur, stābās, stābat, stābātis, stābātur, stābimus, stābis, stābitis, stābitur, stābō, stābunt, stāmus, standī, standō, standum, standus, stāns, stant, stantō, stāre, stārem, stārēmus, stārent, stārēs, stāret, stārētis, stārētur, stārī, stās, stat, stata, statā, statae, statam, statārum, statās, state, stāte, statī, statīs, stātis, statō, stātō, statōrum, statōs, stātōte, statū, statum, stātūr, stātūrus, status, stem, stēmus, stent, stēs, stet, steteram, steterāmus, steterant, steterās, steterat, steterātis, stetēre, steterim, steterīmus, steterīmus, steterint, steteris, steterīs, steterit, steteritis, steterītis, steterō, stetērunt, stetī, stetimus, stētis, stetisse, stetissēmus, stetissent, stetissēs, stetisset, stetissētis, stetistī, stetistis, stetit, stētūr, stētū — “stand”;*

*vetā, vetābam, vetābāminī, vetābāmur, vetābāmus, vetābāntur, vetābar, vetābāre, vetābāris, vetābās, vetābat, vetābātis, vetābātur, vetābere, vetāberis, vetābiminī, vetābimur, vetābīs, vetābit, vetābitis, vetābitur, vetābō, vetābor, vetābuntur, vetāminī, vetāmur, vetāmus, vetandī, vetandō, vetandum, vetāndus, vetāns, vetant, vetantō, vetantor, vetantur, vetāre, vetārem, vetārēminī, vetārēmur, vetārēmus, vetārent, vetārentur, vetārer, vetārēre, vetārēris, vetārēs, vetāret, vetārētis, vetārētur, vetārī, vetāris, vetās, vetat, vetāte, vetātis, vetātō, vetātor, vetātōte, vetātūr, vetem, vetēminī, vetēmur, vetēmus, vetent, vetentur, veter, vetēre, vetēris, vetēs, vetet, vetētis, vetētūr, vetita, vetitā, vetitae, vetitam, vetitārum, vetitās, vetite, vetītī, vetitō, vetitōrum, vetitōs, vetitū, vetitum, vetitūrus, vetitus, vetō, vedor, vetueram, vetuerāmus, vetuerant, vetuerās, vetuerat, vetuerātis, vetuēre, vetuerim, vetuerimus, vetuerīmus, vetuerint, vetueris, vetuerīs, vetuerit, vetueritis, vetuerītis, vetuerō, vetuērunt, vetuī, vetuīmus, vetuisse, vetuisse, vetuissēmus, vetuissent, vetuissēs, vetuisset, vetuissētis, vetuistī, vetuistis, vetuit — “forbid”.*

[https://en.wikipedia.org/wiki/Latin\\_conjugation#Second\\_conjugation](https://en.wikipedia.org/wiki/Latin_conjugation#Second_conjugation)

*auctā, auctā, auctae, auctam, auctārum, auctās, auctī, auctīs, auctō, auctōrum, auctōs, auctū, auctum, auctūrus, auctus, augē, augeām, augeāminī, augeāmur, augeāmus, augeant, augeantur, augear, augeāris, augeās, augeat, augeātis, augeātūr, augēbam, augēbāminī, augēbāmur, augēbāmus, augēbāntur, augēbar, augēbāre, augēbāris, augēbās, augēbat, augēbātis, augēbātūr, augēbere, augēberis, augēbiminī, augēbimur, augēbimus, augēbis, augēbit, augēbitis, augēbitur, augēbō, augēbuntur, augēminī,*

*augēmūr, augēmūs, augēndī, augēndō, augēndum, augēndus, augēns, augēnt, augēntō, augēntor, augēntur, augēō, augēor, augēre, augērem, augērēminī, augērēmūr, augērēmūs, augērent, augērentur, augērer, augērēre, augērēris, augērēs, augēret, augērētis, augērētūr, augērī, augēris, augēs, augēt, augētē, augētis, augētō, augētor, augētōtē, augētūr, auxeram, auxerāmūs, auxerant, auxerās, auxerat, auxerātis, auxēre, auxerim, auxerimus, auxerīmūs, auxerint, auxeris, auxerīs, auxerit, auxeritis, auxerītis, auxerō, auxērunt, auxī, auximus, auxisse, auxissem, auxissēmūs, auxissent, auxissēs, auxisset, auxissētis, auxistī, auxistis, auxit — “increase”;*

*cīē, cieam, cieāminī, cieāmur, cieāmus, cieant, cieantur, clear, cieāre, cieāris, cieās, cieat, cieātis, cieātur, cieābam, ciebāminī, ciebāmur, ciebāmus, ciebant, ciebantur, ciebār, ciebāre, ciebāris, ciebās, ciebat, ciebātis, ciebātūr, ciebere, cieberis, ciebiminī, ciebimur, ciebimus, ciebis, ciebit, ciebitis, ciebitur, ciebō, ciebor, ciebunt, ciebuntur, ciebūnī, ciebūr, ciebūmus, ciendī, ciendō, ciendum, ciendus, ciēns, cient, ciento, cientor, cientur, cieō, cieor, cieere, ciearem, cieremīnī, cieremur, cieremus, ciérent, ciérentur, ciérer, ciéreē, ciéreēris, ciéreēs, ciéret, ciérettis, ciéretur, ciért, ciéris, ciēs, ciet, ciéte, ciéts, ciéto, ciétor, ciétoē, ciéetur, cita, citā, citae, citam, citārum, citās, cite, citī, citīs, citō, citōrum, citōs, citū, citum, citūrus, citus, cíveram, cíverāmus, cíverant, cíverās, cíverat, cíveratis, cívēre, cíverim, cíverimus, cíverūmus, cíverint, cíveris, cíverīs, cíverit, cíveritis, cíverūtis, cíverō, cíverunt, cívī, cívimus, cívisse, cívissem, cívissēmus, cívissent, cívissēs, cívisset, cívissētis, cívistī, cívistis, cívit — “arouse”;*

*docē, doceam, doceāminī, doceāmur, doceāmus, doceant, doceantur, docear, doceāre, doceāris, doceās, doceat, doceātis, doceātur, docēbam, docēbāminī, docēbāmur, docēbāmus, docēbant, docēbantur, docēbar, docēbāre, docēbāris, docēbās, docēbat, docēbātis, docēbātur, docēbere, docēberis, docēbiminī, docēbimur, docēbimus, docēbis, docēbit, docēbitis, docēbitur, docēbō, docēbor, docēbunt, docēminī, docēmūr, docēmus, docendī, docendō, docendum, docendus, docēns, docent, docentō, docentor, docentur, doceō, doceor, docēre, docērem, docērēminī, docērēmur, docērēmus, docērent, docērentur, docērer, docērēre, docērēris, docērēs, docēret, docērētis, docērētur, docērī, docēris, docēs, docet, docēte, docētis, docētō, docētor, docētōte, docētur, docta, doctā, doctae, doctam, doctārum, doctās, docte, doctī, doctīs, doctō, doctōrum, doctōs, doctū, doctum, doctūrus, doctus, docueram, docuerāmus, docuerant, docuerās, docuerat, docuerātis, docuēre, docuerim, docuerimus, docuerīmus, docuerint, docueris, docuerīs, docuerit, docueritis, docuerītis, docuerō, docuērunt, docuī, docuimus, docuisse, docuissēm, docuissēmus, docuissent, docuissēs, docuisset, docuissētis, docuistī, docuistis, docuit — “teach”;*

*ferbueram, ferbuerāmus, ferbuerant, ferbuerās, ferbuerat, ferbuerātis, ferbuēre, ferbuerim, ferbuerimus, ferbuerīmus, ferbuerint, ferbueris, ferbuerīs, ferbuerit, ferbueritis, ferbuerītis, ferbuerō, ferbuērunt, ferbuī, ferbuīmus, ferbuisse, ferbuīsem, ferbuīsēmus, ferbuīscent, ferbuīssēs, ferbuīsset, ferbuīsētis, ferbuīstī, ferbuīstis, ferbuit, fervê, ferveam, ferveāmus, ferveant, ferveās, ferveat, ferveātis, fervēbam, fervēbāmus, fervēbant, fervēbās, fervēbat, fervēbātis, fervēbimus, fervēbis, fervēbit, fervēbitis, fervēbō, fervēbunt, fervēmus, servendī, servendō, servendum, servēns, servent, serventō, serveō, servēre, servērem, servērēmus, servērent, servērēs, servēret, servērētis, servēs, servet, servēte, servētis, servētō, servētōte, servitū, servitum, servitūrus — “boil”;*

*fōta, fōtā, fōtae, fōtam, fōtarum, fōtās, fōte, fōtī, fōtīs, fōtō, fōtōrum, fōtōs, fōtū, fōtum, fōtūrus, fōtus, fōvē, fōveam, fōveāmī, fōveāmur, fōveāmus, fōveant, fōveantur, fōvear, fōveāre, fōveāris, fōveās, fōveat, fōveātis, fōveātūr, fōvēbam, fōvēbāmī, fōvēbāmur, fōvēbāmus, fōvēbant, fōvēbantur, fōvēbar, fōvēbāre, fōvēbāris, fōvēbās, fōvēbat, fōvēbātis, fōvēbātūr, fōvēbere, fōvēberis, fōvēbimī, fōvēbimur, fōvēbis, fōvēbit, fōvēbitis, fōvēbitur, fōvēbō, fōvēbor, fōvēbunt, fōvēbuntur, fōvēminī, fōvēmur, fōvēmus, fōvendī, fōvendō, fōvendum, fōvendus, fōvēns, fōvent, fōventō, fōventor, fōventur, fōveō, fōveor, fōveram, fōverāmus, fōverant, fōverās, fōverat, fōverātis, fōvēre, fōvēre, fōvērem, fōvēremī, fōvēremur, fōvēremus, fōvērent, fōvērentur, fōvērer, fōvērēre, fōvērēris, fōvērēs, fōvēret, fōvērētis, fōvērētūr, fōvērī, fōverim, fōverimus, fōverīmus, fōverint, fōveris, fōverīs, fōverit, fōverit,*

*fōveritis, fōverītis, fōverō, fōvērunt, fovēs, fovent, fovēte, fovētis, fovētō, fovētor, fovētōte, fovētur, fōvī, fōvīmus, fōvisse, fōvissem, fōvissēmus, fōvissent, fōvissēs, fōvisset, fōvissētis, fōvistī, fōvistis, fōvit — “caress”;*

*iubē, iubeāminī, iubeāmur, iubeāmus, iubeant, iubeantur, iubēar, iubeāris, iubeās, iubeat, iubēatis, iubeātur, iubēbam, iubēbāminī, iubēbāmur, iubēbāmus, iubēbant, iubēbantur, iubēbar, iubēbāre, iubēbāris, iubēbās, iubēbat, iubēbātis, iubēbātūr, iubēbere, iubēberis, iubēbimīnī, iubēbimur, iubēbimus, iubēbis, iubēbit, iubēbitis, iubēbitur, iubēbō, iubēbor, iubēbunt, iubēbuntur, iubēminī, iubēmur, iubēmus, iubēndī, iubēndō, iubēndum, iubēndus, iubēns, iubēnt, iubēntō, iubēntor, iubēntur, iubēō, iubēor, iubēre, iubērem, iubērēmīnī, iubērēmur, iubērēmus, iubērent, iubērentur, iubērer, iubērēre, iubērēris, iubērēs, iubēret, iubērētis, iubērētūr, iubērī, iubēris, iubēs, iubēte, iubētis, iubētō, iubētor, iubētōte, iubētūr, iussa, iussā, iussae, iussam, iussārum, iussās, iusse, iusseram, iusserāmus, iusserant, iusserās, iusserat, iusserātis, iusserē, iusserim, iusserimus, iusserīmus, iusserint, iusseris, iusserīs, iusserit, iusserītis, iusserō, iusserunt, iussī, iussīmus, iussīs, iussisse, iussissem, iussissēmus, iussissent, iussissēs, iussisset, iussissētis, iussistī, iussistis, iussit, iussō, iussōrum, iussōs, iussū, iussum, iussūrus, iussus — “order”;*

*momorderam, momorderāmus, momorderant, momorderās, momorderat, momorderātis, momordēre, momordērim, momorderimus, momorderīmus, momorderint, momorderis, momorderīs, momorderit, momorderitis, momorderītis, momorderō, momordērunt, momordī, momordīmus, momordīsse, momordīssēmus, momordīssent, momordīssēs, momordīsset, momordīssētis, momordīstī, momordīstis, momordīt, mordē, mordeām, mordeāminī, mordeāmur, mordeāmus, mordeant, mordeantur, mordear, mordeāre, mordeāris, mordeās, mordeat, mordeātis, mordeātūr, mordēbam, mordēbāminī, mordēbāmur, mordēbāmus, mordēbant, mordēbantur, mordēbar, mordēbāre, mordēbāris, mordēbās, mordēbat, mordēbātis, mordēbātūr, mordēbere, mordēberis, mordēbimīnī, mordēbimur, mordēbimus, mordēbis, mordēbit, mordēbitis, mordēbitur, mordēbō, mordēbor, mordēbunt, mordēbuntur, mordēminī, mordēmur, mordēmus, mordēndī, mordēndō, mordēndum, mordēndus, mordēns, mordēnt, mordēntō, mordēntor, mordēntur, mordēō, mordēor, mordēre, mordērem, mordērēmīnī, mordērēmūr, mordērēmus, mordērent, mordērentur, mordērer, mordērēre, mordērēris, mordērēs, mordēret, mordērētis, mordērētūr, mordērī, mordēris, mordēs, mordēte, mordētis, mordētō, mordētor, mordētōte, mordētūr, morsa, morsā, morsae, morsam, morsārum, morsās, morse, morsī, morsīs, morsō, morsōrum, morsōs, morsū, morsum, morsūrus, morsus — “bite”;*

*monē, moneam, moneāminī, moneāmur, moneāmus, moneant, moneantur, monear, moneāre, moneāris, moneās, moneat, moneātis, moneātur, monēbam, monēbāminī, monēbāmur, monēbāmus, monēbant, monēbantur, monēbar, monēbāre, monēbāris, monēbās, monēbat, monēbātis, monēbātūr, monēbere, monēberis, monēbimīnī, monēbimur, monēbimus, monēbis, monēbit, monēbitis, monēbitur, monēbō, monēbor, monēbunt, monēbuntur, monēminī, monēmur, monēmus, monendī, monendō, monendum, monendus, monēns, monent, monentō, monentor, monentur, moneō, moneor, monēre, monērem, monērēmīnī, monērēmur, monērēmus, monērent, monērentur, monērer, monērēre, monērēris, monēret, monērētis, monērētūr, monērī, monēris, monēs, monet, monēte, monētis, monētō, monētor, monētōte, monētūr, monitū, monitūrūs, monitus, monueram, monuerām, monuerant, monuerās, monuerat, monuerātis, monuēre, monuerim, monuerimus, monuerītis, monuerint, monueris, monuerīs, monuerit, monueritis, monuerītis, monuerō, monuērunt, monuī, monuīmus, monuīsse, monuīssem, monuīssēmus, monuīssent, monuīssēs, monuīsset, monuīssētis, monuīstī, monuīstis, monuīt — “warn”;*

*polliceāminī, policeāmur, policeantur, pollicear, policeāre, policeāris, policeātūr, pollicēbāminī, pollicēbāmur, pollicēbantur, pollicēbar, pollicēbāre, pollicēbāris, pollicēbātūr, pollicēbere, pollicēberis, pollicēbimīnī, pollicēbimur, pollicēbitur, pollicēbor, pollicēbuntur, pollicēminī, pollicēmur, pollicēndī, pollicēndō, pollicēndum, pollicēndus, pollicēns, pollicēntor, pollicēntur, polliceor, pollicēre, pollicērēmīnī, pollicērēmur, pollicērentur, pollicērer, pollicērēre, pollicērēris, pollicērētūr, pollicērī, pollicēris, pollicētor, pollicētūr, pollicita, pollicitā, pollicitae, pollicitam, pollicitārum, pollicitās, pollicite, pollicitī, pollicitīs, pollicitō, pollicitōrum, pollicitōs, pollicitū, pollicitum, pollicitūrus, pollicitus — “promise”;*

*spondē, spondeam, spondeāminī, spondeāmur, spondeāmus, spondeant, spondeantur, spondear, spondeāre, spondeāris, spondeās, spondeat, spondeātis, spondeātūr, spondēbam, spondēbāminī, spondēbāmur, spondēbāmus, spondēbant, spondēbantur, spondēbar, spondēbāre, spondēbāris, spondēbās, spondēbat, spondēbātis, spondēbātūr, spondēbere, spondēberis, spondēbimīnī, spondēbimur, spondēbimus, spondēbis, spondēbit, spondēbitis, spondēbitur, spondēbō, spondēbor, spondēbunt, spondēbuntur, spondēminī, spondēmur, spondēmus, spondēndī, spondēndō, spondēndum, spondēndus, spondēns, spondēnt, spondēntō, spondēntor, spondēntur, spondeō, spondēre, spondērem, spondērēmīnī, spondērēmur, spondērent, spondērer,*

*spondērēre, spondērēris, spondērēs, spondēret, spondētēs, spondērētur, spondērī, spondēris, spondēs, spondet, spondēte, spondētis, spondētō, spondētor, spondētōte, spondētur, spōnse, spōnsī, spōnsīs, spōnsō, spōnsōrum, spōnsōs, spōnsū, spōnsum, spōnsūrus, spōnsus, spoponderam, spoponderāmus, spoponderant, spoponderās, spoponderat, spoponderātis, spopondēre, spoponderim, spoponderimus, sponderītus, spoponderint, spoponderis, spoponderīs, spoponderit, spoponderitis, spoponderītis, spoenderō, spōndērunt, spōndī, spōndīmus, spōndisse, spōndissem, spōndissēmus, spōndissent, spōndissēs, spōndisset, spōndissētis, spōndistī, spōndistis, spōndit — “vow”;*

*stridē, strideam, strideāminī, strideāmur, strideāmus, strideant, strideantur, stridear, strideāre, strideāris, strideās, strideat, strideātis, strideātur, stridēbam, stridēbāminī, stridēbāmur, stridēbāmus, stridēbānt, stridēbāntur, stridēbar, stridēbāre, stridēbāris, stridēbās, stridēbat, stridēbātis, stridēbātur, stridēbere, stridēberis, stridēbīminī, stridēbīmur, stridēbīmus, stridēbīs, stridēbit, stridēbitis, stridēbitur, stridēbō, stridēbor, stridēbunt, stridēbuntur, stridēminī, stridēmur, stridēmus, stridēndī, stridēndō, stridēndum, stridēndus, stridēns, strident, stridentō, stridentor, stridentur, stridēō, strideor, stridēre, stridērem, stridērēminī, stridērēmur, stridērēmus, stridērent, stridērentur, stridērer, stridērēre, stridērēris, stridērēs, stridēret, stridērētis, stridērētur, stridērī, stridēris, stridēs, stridet, stridēte, stridētis, stridētō, stridētor, stridētōte, stridētūr, stridū, stridum, stridūrus — “hiss”;*

*tenē, teneam, teneāminī, teneāmur, teneāmus, teneant, teneantur, tenear, teneāre, teneāris, teneās, teneat, teneātis, teneātur, tenēbam, tenēbāminī, tenēbāmur, tenēbāmus, tenēbant, tenēbantur, tenēbar, tenēbāre, tenēbāris, tenēbās, tenēbat, tenēbātis, tenēbātur, tenēbere, tenēberis, tenēbiminī, tenēbimur, tenēbimus, tenēbis, tenēbit, tenēbitis, tenēbitur, tenēbō, tenēbor, tenēbunt, tenēbuntur, tenēminī, tenēmur, tenēmus, tenendī, tenendō, tenendum, tenendus, tenēns, tenent, tenentō, tenentor, tenentur, teneō, teneor, tenēre, tenērem, tenērēminī, tenērēmur, tenērēmus, tenērent, tenērentur, tenērer, tenērēre, tenērēris, tenērēs, tenēret, tenērētis, tenērētur, tenērī, tenēris, tenēs, tenet, tenēte, tenētis, tenētō, tenētor, tenētōte, tenētūr, tenueram, tenuerāmus, tenuerant, tenuerās, tenuerat, tenuerātis, tenuēre, tenuerim, tenuerimus, tenuerīmus, tenuerint, tenueris, tenuerīs, tenuerit, tenueritis, tenuerītis, tenuerō, tenuerunt, tenuī, tenuimus, tenuisse, tenuissem, tenuissēmus, tenuissent, tenuissēs, tenuisset, tenuissētis, tenuistī, tenuistis, tenuit — “hold”;*

*terrē, terream, terreāmīnī, terreāmūr, terreāmūs, terreat, terreant, terreantur, terrear, terreāre, terreāris, terreās, terreat, terreātis, terreātur, terrēbam, terrēbāmīnī, terrēbāmūr, terrēbāmūs, terrēbant, terrēbantur, terrēbar, terrēbāre, terrēbāris, terrēbās, terrēbat, terrēbātis, terrēbātūr, terrēbere, terrēberis, terrēbimīnī, terrēbimūr, terrēbimūs, terrēbis, terrēbit, terrēbitis, terrēbitur, terrēbō, terrēbor, terrēbunt, terrēbuntur, terrēmīnī, terrēmur, terrēmus, terrendī, terrendō, terrendum, terrendus, terrēns, torrent, terrentō, torrentor, torrentur, terreō, terreor, terrēre, terrērem, terrērēmīnī, terrērēmūr, terrērēmūs, terrērent, terrērentur, terrērer, terrērēre, terrērēris, terrērēs, terrēret, terrēretis, terrēretur, terrērī, terrēris, terrēs, terret, terrēte, terrētis, terrētō, terrētor, terrētōte, terrētūr, territa, territā, territae, territam, territārum, territās, territe, territī, territīs, territō, territōrum, territōs, territū, territum, territūrus, territus, terrueram, terruerāmūs, terruerant, terruerās, terruerat, terruerātis, terruēre, terruerim, terruerimūs, terruerīmūs, terruerint, terrueris, terruerīs, terruerit, terrueritis, terruerītis, terruerō, terruerunt, terruū, terruimus, terruisse, terruissem, terruissēmūs, terruissent, terruissēs, terruisset, terruissētis, terruistī, terruistis, terruit — “frighten”;*

*vidē, videam, videāminī, videāmur, videāmus, videant, videantur, videar, videāris, videās, videat, videātis, videātur, vidēbam, vidēbāminī, vidēbāmur, vidēbāmus, vidēbant, vidēbantur, vidēbar, vidēbāre, vidēbāris, vidēbās, vidēbat, vidēbātis, vidēbātūr, vidēbere, vidēberis, vidēbiminī, vidēbimur, vidēbimus, vidēbis, vidēbit, vidēbitis, vidēbitur, vidēbō, vidēbor, vidēbunt, vidēbuntur, vidēminī, vidēmur, vidēmus, videndī, videndō, videndum, videndus, vidēns, vidēnt, vidēntō, vidēntor, vidēntur, videō, videor, vīderam, vīderāmus, vīderant, vīderās, vīderat, vīderātis, vidēre, vīdērem, vidēreminī, vidēremur, vidēremus, vidērent, vidērentur, vidērer, vidērēre, vidērēris, vidērēs, vidēret, vidērētis, vidērētur, vidērī, vīderim, vīderimus, vīderīmus, vīderint, vidēris, vīderis, vīderīs, vīderit, vīderitis, vīderītis, vīderō, vīdērunt, vidēs, videt, vidēte, vidētis, vidētō, vidētor, vidētōte, vidētūr, vīdī, vīdimus, vīdisse, vīdissem, vīdissēmus, vīdissent, vīdissēs, vīdisset, vīdissētis, vīdistī, vīdistis, vīdit, vīsa, vīsā, vīsae, vīsam, vīsārum, vīsās, vīse, vīsī, vīsīs, vīsō, vīsōrum, vīsōs, vīsū, vīsum, vīstūrus, vīsus — “see”.*

[https://en.wikipedia.org/wiki/Latin\\_conjugation#Third\\_conjugation](https://en.wikipedia.org/wiki/Latin_conjugation#Third_conjugation)

*ācta, āctā, āctae, āctam, āctārum, āctās, ācte, āctī, āctīs, āctō, āctōrum, āctōs, āctū, āctum, āctūrus, āctus, agam, agāminī, agāmur, agāmus, agant, agantur, agar, agāre, agāris, agās, agat, agātis, agātur, age, agēbam, agēbāminī, agēbāmur, agēbāmus, agēbant, agēbantur, agēbar, agēbāre, agēbāris, agēbās, agēbat, agēbātis, agēbātur, agēminī, agēmur, agēmus, agendī, agendō, agendum, agendus, agēns, agent, agentur, agere, agēre, agerem, agerēminī, agerēmur, agerēmus, agerent, agerentur, agerer, agerēre, agerēris, agerēs, ageret, agerētis,*

*agerētur, ageris, agēris, agēs, agēt, agētur, agī, agiminī, agimur, agimus, agis, agit, agite, agitis, agitō, agitor, agitōte, agitur, agō, agor, agunt, aguntō, aguntor, aguntur, ēgerāmus, ēgerās, ēgerat, ēgerātis, ēgēre, ēgerim, ēgerimus, ēgerīmus, ēgerint, ēgeris, ēgerīs, ēgerit, ēgeritis, ēgerītis, ēgerō, ēgerūnt, ēgī, ēgimus, ēgisse, ēgissēm, ēgissēmus, ēgissent, ēgissēs, ēgisset, ēgissētis, ēgistī, ēgistis, ēgit* — “drive”;

*adhaerēscam, adhaerēscāmīnī, adhaerēscāmūr, adhaerēscāmūs, adhaerēscānt, adhaerēscāntur, adhaerēscār, adhaerēscāre, adhaerēscāris, adhaerēscās, adhaerēscāt, adhaerēscātās, adhaerēscātār, adhaerēscātār, adhaerēscēbam, adhaerēscēbāmīnī, adhaerēscēbāmūr, adhaerēscēbāmūs, adhaerēscēbānt, adhaerēscēbāntur, adhaerēscēbar, adhaerēscēbāre, adhaerēscēbāris, adhaerēscēbās, adhaerēscēbat, adhaerēscēbātīs, adhaerēscēbātūr, adhaerēscēmīnī, adhaerēscēmūr, adhaerēscēmūs, adhaerēscēndī, adhaerēscēndō, adhaerēscēndūr, adhaerēscēndus, adhaerēscēns, adhaerēscēnt, adhaerēscēntur, adhaerēscērē, adhaerēscērē, adhaerēscērēmīnī, adhaerēscērēmūr, adhaerēscērēmūs, adhaerēscērēnt, adhaerēscērēntur, adhaerēscērērē, adhaerēscērēris, adhaerēscērēs, adhaerēscērēt, adhaerēscērētīs, adhaerēscērētūr, adhaerēscērēs, adhaerēscēris, adhaerēscēs, adhaerēscēt, adhaerēscētīs, adhaerēscētūr, adhaerēscī, adhaerēscimīnī, adhaerēscimūr, adhaerēscimūs, adhaerēscis, adhaerēscit, adhaerēscite, adhaerēscitīs, adhaerēscitō, adhaerēscitor, adhaerēscitōte, adhaerēscitūr, adhaerēscō, adhaerēscor, adhaerēscunt, adhaerēscuntō, adhaerēscuntor, adhaerēscuntur — “adhere”;*

*adolēscam, adolēscāmīnī, adolēscāmūr, adolēscāmūs, adolēscānt, adolēscāntur, adolēscār, adolēscāre, adolēscāris, adolēscās, adolēscāt, adolēscātis, adolēscātūr, adolēscē, adolēscēbam, adolēscēbāmīnī, adolēscēbāmūr, adolēscēbāmūs, adolēscēbānt, adolēscēbāntur, adolēscēbar, adolēscēbārē, adolēscēbāris, adolēscēbās, adolēscēbat, adolēscēbātīs, adolēscēbātūr, adolēscēmīnī, adolēscēmūr, adolēscēmūs, adolēscēndāl, adolēscēndō, adolēscēndum, adolēscēndus, adolēscēns, adolēscēnt, adolēscēntur, adolēscērē, adolēscērem, adolēscērēmīnī, adolēscērēmūr, adolēscērēmūs, adolēscērent, adolēscērentur, adolēscērer, adolēscērerē, adolēscērerīs, adolēscērēs, adolēscēret, adolēscērētīs, adolēscērētūr, adolēscēris, adolēscērīs, adolēscēt, adolēscētīs, adolēscētūr, adolēscī, adolēscimīnī, adolēscimūr, adolēscimūs, adolēscis, adolēscit, adolēscite, adolēscitis, adolēscitō, adolēscitor, adolēscitōtē, adolēscitur, adolēscō, adolēscor, adolēscunt, adolēscuntō, adolēscuntor, adolēscuntur, adolēveram, adolēverāmūs, adolēverant, adolēverās, adolēverat, adolēverātīs, adolēvērē, adolēverim, adolēverimus, adolēverīmūs, adolēverint, adolēveris, adolēverīs, adolēverit, adolēveritis, adolēverītīs, adolēverō, adolēverunt, adolēvī, adolēvīmūs, adolēvisse, adolēvissem, adolēvissēmūs, adolēvissent, adolēvissēs, adolēvisset, adolēvissētīs, adolēvistī, adolēvistis, adolēvit, adulta, adultā, adultae, adultam, adultārum, adultās, adulte, adultī, adultīs, adultō, adultōrum, adultōs, adultū, adultum, adultūrus, adultus — “mature”;*

caedam, caedāminī, caedāmur, caedāmus, caedant, caedantur, caedar, caedāre, caedāris, caedās, caedat, cae-  
dātis, caedātur, caede, caedēbam, caedēbāminī, caedēbāmur, caedēbāmus, caedēbant, caedēbantur, caedēbar,  
caedēbāre, caedēbāris, caedēbās, caedēbat, caedēbātis, caedēbātur, caedēminī, caedēmur, caedēmus, caedendī,  
caedendō, caedendum, caedendus, caedēns, caedent, caedentur, caedere, caedēre, caederem, caederēminī, cae-  
derēmur, caederēmus, caederent, caederentur, caederer, caederēre, caederēris, caederēs, caederet, caederētis,  
caederētūr, caederis, caedēris, caedēs, caedet, caedētis, caedētūr, caedī, caedimīnī, caedimur, caedimus, cae-  
dis, caedit, caedite, caeditis, caeditō, caeditor, caeditōte, caeditur, caedō, caedor, caedunt, caeduntō, caeduntor,  
caeduntur, caesa, caesā, caesae, caesam, caesārum, caesās, caese, caesī, caesīs, caesō, caesōrum, caesōs, cae-  
sū, caesum, caesūrus, caesus, cecīderam, cecīderāmus, cecīderant, cecīderās, cecīderat, cecīderātis, cecīdere,  
cecīderim, cecīderimus, cecīderīmus, cecīderint, cecīderis, cecīderīs, cecīderit, cecīderitis, cecīderītis, cecīderō,  
cecīderunt, cecīdī, cecīdimus, cecīdisse, cecīdissem, cecīdissēmus, cecīdissent, cecīdissēs, cecīdisset, cecīdissētis,  
cecīdistī, cecīdistis, cecīdit — “kill”;

*carpam, carpāminī, carpāmur, carpāmus, carpant, carpantur, carpar, carpāre, carpāris, carpās, carpat, carpātis, carpātur, carpe, carpēbam, carpēbāminī, carpēbāmur, carpēbāmus, carpēbant, carpēbantur, carpēbar, carpēbāre, carpēbāris, carpēbās, carpēbat, carpēbātis, carpēbātur, carpēminī, carpēmur, carpēmus, carpēndī, carpēndō, carpēndum, carpēndus, carpēns, carpēnt, carpēntur, carpēre, carpēre, carpērem, carpēremīnī, carpēremur, carpēremus, carpērent, carpērentur, carpērer, carpērēre, carpērēris, carpērēs, carpēret, carpērētis, carpērētūr, carpēris, carpēris, carpēs, carpētis, carpētūr, carpī, carpimīnī, carpimur, carpimus, carpis, carpit, carpite, carpitis, carpītō, carpitor, carpītōte, carpitur, carpō, carpor, carpseram, carpserāmus, carpserant, carpserās, carpserat, carpserātis, carpsēre, carpserim, carpserimus, carpserīmus, carpserint, carpse-  
ris, carpsēris, carpsērit, carpseritis, carpserītis, carpserō, carpsērunt, carpsī, carpsimus, carpsisse, carpsissem, carpsissēmus, carpsissent, carpsissēs, carpsisset, carpsissētis, carpsistī, carpsistis, carpsit, carpta, carptā, carptae, carptam, carptārum, carptās, carpte, carptī, carptis, carptō, carptōrum, carptōs, carptū, carptum, carptūrus, carptus, carpunt, carpuntō, carpuntor, carpuntur — “pluck”;*

*colam, colāminī, colāmur, colāmus, colant, colantur, collar, colāre, colāris, colās, colat, colātis, colātur, cole, colēbam, colēbāminī, colēbāmur, colēbāmus, colēbant, colēbantur, colēbar, colēbāre, colēbāris, colēbās, colēbat, colēbātis, colēbātur, colēminī, colēmur, colēmus, colendī, colendō, colendum, colendus, colēns, co lent, co lentur, colere, colēre, colerem, colerēminī, colerēmur, colerēmus, colerent, colerentur, colerer, colerēre, colerēris, colerēs, coleret, colerētis, colerētur, coleris, colēris, colēs, colet, colētis, colētur, colī, colimīnī, colimur, colimus, colis, colit, colite, colitis, colitō, colitor, colitōte, colitur, colō, color, colueram, coluerāmus, coluerant, coluerās, coluerat, coluerātis, coluēre, coluerim, coluerimus, coluerīmus, coluerint, colueris, coluerīs, coluerit, colueritis, coluerītis, coluerō, coluērunt, coluī, coluimus, coluisse, coluisse, coluissēmus, coluissent, coluissēs, coluisset, coluissētis, coluistī, coluistis, coluit, colunt, coluntō, coluntor, coluntur, culta, cultā, cultae, cultam, cultārum, cultās, culte, cultī, cultīs, cultō, cultōrum, cultōs, cultū, cultum, cultūrus, cultus — “cultivate”;*

*cucurreram, cucurrerāmus, cucurrerant, cucurrerās, cucurrerat, cucurrerātis, cucurrēre, cucurrerim, cucurrerimus, cucurrerīmus, cucurrerint, cucurreris, cucurrerīs, cucurrerit, cucurreritis, cucurrerītis, cucurrerō, cucurrerūrunt, cucurrī, cucurrimus, cucurrisse, cucurrissem, cucurrissemus, cucurrissent, cucurrissēs, cucurrisset, cucurrissētis, cucurristū, cucurristis, cucurrit, currām, currāminī, currāmūr, currāmūs, currānt, currāntur, currar, currāre, currāris, currās, currāt, currātis, currātur, curre, currēbam, currēbāmīnī, currēbāmūr, currēbāmūs, currēbānt, currēbāntur, currēbar, currēbāre, currēbāris, currēbās, currēbat, currēbātis, currēbātur, currēminī, currēmūr, currēmūs, currēndī, currēndō, currēndum, currēndus, currēns, currēnt, currēntur, currēre, currērem, currēminī, currēmūr, currēmūs, currērent, currērentur, currērer, currērēre, currērēis, currērēs, currēret, currētis, currētēr, currēris, currēs, currēt, currētis, currētēr, currī, currīmīnī, currīmūr, currīmūs, currīmūs, currīs, currīt, currīte, currītis, currītō, currītor, currītōtē, currītūr, currō, currōr, currūnt, currūntō, currūntor, currūntur, cursā, cursāt, cursās, cursām, cursārūm, cursās, curse, cursī, cursīs, cursō, cursōrūm, cursōs, cursū, cursum, cursūrūs, cursus — “run”;*

*emam, emāminī, emāmur, emāmus, emant, emantur, emar, emāre, emāris, emās, emat, emātis, emātur,eme, emēbam, emēbāminī, emēbāmur, emēbāmus, emēbant, emēbantur, emēbar, emēbāre, emēbāris, emēbās, emēbat, emēbātis, emēbātur, emēminī, emēmur, emēmus, emendī, emendō, emendum, emendus, emēns, ement, ementur, emeram, emerāmus, emerant, emerās, emerat, emerātis, emere, emēre, emēre, emerem, emerēminī, emerēmur, emerēmus, emerent, emerentur, emerer, emerēre, emerēris, emerēs, emeret, emerētis, emerētur, emerim, emerimus, emerīmus, emerint, emeris, emēris, emerīs, emerit, emeritis, emerītis, emerō, emērunt, emēs, emet, emētis, emētur, emī, emīt, emimīnī, emimur, emimus, emimūs, emis, emīs, emīsse, emīssēmūs, emīssēnt, emīsēs, emīsset, emīssētis, emīstī, emīstis, emit, emīt, emite, emitis, emitō, emitōr, emitōte, emitur, emō, emor, empia, emptā, emptae, emptam, emptārum, emptās, empē, emptī, emptīs, emptō, emptōrum, emptōs, emptū, emptum,emptum, emptūrus, emptus, ēmpetus, emunt, emuntō, emuntor, emuntur — “buy”;*

*flectam, flectāminī, flectāmūr, flectāmus, flectant, flectantur, flectar, flectāre, flectāris, flectās, flectat, flectātis, flectātur, flecte, flectēbam, flectēbāminī, flectēbāmūr, flectēbāmus, flectēbant, flectēbantur, flectēbar, flectēbāre, flectēbāris, flectēbās, flectēbat, flectēbātis, flectēbātūr, flectēminī, flectēmūr, flectēmus, flectendī, flectendō, flectendum, flectendus, flectēns, flectent, flectentur, flectere, flectēre, flecterem, flecterēminī, flecterēmūr, flecterēmus, flecterent, flecterentur, flecterer, flecterēre, flecterēris, flecterēs, flecteret, flecterētis, flecterētūr, flecteris, flectēris, flectēs, flectet, flectētūr, flectī, flectimīnī, flectimūr, flectimus, flectis, flectit, flectite, flectitis, flectitō, flectitor, flectitōtē, flectitūr, flectō, flector, flectunt, flectuntō, flectuntor, flectuntur, flexa, flexā, flexae, flexam, flexārum, flexās, flexe, flexeram, flexerāmus, flexerant, flexerās, flexerat, flexerātis, flexēre, flexerim, flexerimus, flexerīmus, flexerint, flexeris, flexerīs, flexerit, flexeritis, flexerītis, flexerō, flexerunt, flexī, fleximus, flexīs, flexisse, flexissem, flexissēmus, flexissent, flexissēs, flexisset, flexissētis, flexistī, flexistis, flexit, flexō, flexōrum, flexōs, flexū, flexum, flexūrus, flexus — “bend”;*

*flōrēscam, flōrēscāminī, flōrēscāmur, flōrēscāmus, flōrēscant, flōrēscantur, flōrēscar, flōrēscāre, flōrēscāris, flōrēscās, flōrēscat, flōrēscātis, flōrēscātur, flōrēsce, flōrēscēbam, flōrēscēbāminī, flōrēscēbāmur, flōrēscēbāmus, flōrēscēbant, flōrēscēbantur, flōrēscēbar, flōrēscēbāre, flōrēscēbāris, flōrēscēbās, flōrēscēbat, flōrēscēbātis, flōrēscēbātūr, flōrēscēminī, flōrēscēmur, flōrēscēmus, flōrēscēndi, flōrēscēndō, flōrēscēndum, flōrēscēndus, flōrēscēns, flōrēscēnt, flōrēscēntur, flōrēscēre, flōrēscērē, flōrēscērēm, flōrēscērēmūr, flōrēscērēmūs, flōrēscērēnt, flōrēscērēntur, flōrēscērērē, flōrēscērēris, flōrēscērēs, flōrēscērēt, flōrēscērētis, flōrēscērētūr, flōrēscērēris, flōrēscēris, flōrēscēs, flōrēscēt, flōrēscētis, flōrēscētūr, flōrēscētī, flōrēscimī, flōrēscimūr, flōrēscimūs, flōrēscis, flōrēscit, flōrēscite, flōrēscitis, flōrēscitō, flōrēscitor, flōrēscitōte, flōrēscitur, flōrēscō, flōrēscor, flōrēscunt, flōrēscuntō, flōrēscuntor, flōrēscuntur — “blossom”;*

*fūderam, fūderāmus, fūderant, fūderās, fūderat, fūderatis, fūdere, fūderim, fūderimus, fūderīmus, fūderint, fūderis, fūderīs, fūderit, fūderitis, fūderītis, fūderō, fūdērunt, fūdī, fūdimus, fūdisse, fūdissem, fūdissēmus, fūdissent,*

*fūdissēs, fūdisset, fūdissētis, fūdistī, fūdistis, fūdit, fundam, fundāmus, fundant, fundantur, fundās, fundat, fundātis, fundātur, funde, fundēbam, fundēbāmus, fundēbant, fundēbantur, fundēbās, fundēbat, fundēbātis, fundēbātur, fundēmus, fundendī, fundendō, fundendum, fundendus, fundēns, fundent, fundentur, fundere, funderem, funderēmus, funderent, funderentur, funderēs, funderet, funderētis, funderētur, fundēs, fundet, fundētis, fundētūr, fundī, fundimus, fundis, fundit, fundite, funditis, funditō, funditōte, funditur, fundō, fundunt, funduntō, funduntur, fūsa, fūsā, fūsae, fūsam, fūsārum, fūsās, fūse, fūst, fūsīs, fūsō, fūsōrum, fūsōs, fūsū, fūsum, fūsūrus, fūsus — “pour”;*

*genita, genitā, genitae, genitam, genitārum, genitās, genite, genitī, genitīs, genitō, genitōrum, genitōs, genitū, genitum, genitūrus, genitus, genueram, genuerāmus, genuerant, genuerās, genuerat, genuerātis, genuēre, genuērim, genuerimus, genuerīmus, genuerint, genueris, genuerīs, genuerit, genueritis, genuerītis, genuerō, genuerunt, genuī, genuimus, genuisse, genuissem, genuissēmus, genuissent, genuissēs, genuisset, genuissētis, genuistī, genuistis, genuit, gignam, gignāmīnī, gignāmūr, gignāmūs, gignant, gignantur, gignar, gignāre, gignāris, gignās, gignat, gignātis, gignātūr, gigne, gignēbam, gignēbāmīnī, gignēbāmūr, gignēbāmūs, gignēbānt, gignēbāntur, gignēbar, gignēbāre, gignēbāris, gignēbās, gignēbat, gignēbātis, gignēbātūr, gignēmīnī, gignēmūr, gignēmūs, gignendī, gignendō, gignendum, gignendus, gignēns, gignent, gignentur, gignere, gignēre, gignerem, gignerēmīnī, gignerēmūr, gignerēmūs, gignerent, gignerentur, gignerer, gignerēre, gignerēris, gignerēs, gigneret, gignerētis, gignerētūr, gigneris, gignēris, gignēs, gignet, gignētis, gignētūr, gignī, gignimīnī, gignimūr, gignimūs, gignis, gignit, gignite, gignitō, gignitor, gignitōte, gignitur, gignō, gignor, gignunt, gignuntō, gignuntur, gignuntur — “beget”;*

*geram, gerāmīnī, gerāmūr, gerāmūs, gerant, gerantur, gerar, gerāre, gerāris, gerās, gerat, gerātis, gerātūr, gere, gerēbam, gerēbāmīnī, gerēbāmūr, gerēbāmūs, gerēbānt, gerēbāntur, gerēbār, gerēbāris, gerēbās, gerēbāt, gerēbātis, gerēbātūr, gerēminī, gerēmūr, gerēmūs, gerēndī, gerēndō, gerēndum, gerēndus, gerēns, gerent, gerentur, gerere, gerēre, gererem, gererēmīnī, gererēmūr, gererēmūs, gererent, gererentur, gererer, gererēre, gererēris, gererēs, gereret, gererētis, gererētūr, gereris, gerēris, gerēs, geret, gerētis, gerētūr, gerētī, gerimīnī, gerimur, gerimus, geris, gerit, gerite, geritis, geritō, geritor, geritōte, geritūr, gerō, geror, gerunt, geruntō, geruntor, geruntur, gesseram, gesserāmūs, gesserant, gesserās, gesserat, gesserātis, gessēre, gesserim, gesserimūs, gesserīmūs, gesserint, gesseris, gesserīs, gesserit, gesseritis, gesserītis, gesserō, gesserēunt, gessī, gessimus, gessisse, gessissem, gessissēmus, gessissent, gessissēs, gessisset, gessissētis, gessistī, gessistis, gessit, gesta, gestā, gestae, gestam, gestārum, gestās, geste, gestī, gestīs, gestō, gestōrum, gestōs, gestū, gestum, gestūrus, gestus — “wear”;*

*īcam, īcāmīnī, īcāmūr, īcāmūs, īcant, īcantur, īcar, īcāre, īcāris, īcās, īcat, īcātis, īcātūr, īce, īcēbam, īcēbāmīnī, īcēbāmūr, īcēbāmūs, īcēbānt, īcēbāntur, īcēbar, īcēbāre, īcēbāris, īcēbās, īcēbat, īcēbātis, īcēbātūr, īcēmīnī, īcēmūr, īcēmūs, īcēndī, īcēndō, īcēndum, īcēndus, īcēns, īcent, īcentur, īceram, īcerāmūs, īcerant, īcerās, īcerat, īcerātis, īcere, īcēre, īcerem, īcerēmīnī, īcerēmūr, īcerēmūs, īcerent, īcerentur, īcerer, īcerēre, īcerēris, īcerēs, īceret, īcerētis, īcerētūr, īcerim, īcerimus, īcerīmūs, īcerint, īceris, īcerīs, īceris, īcerit, īceritis, īcerō, īcērunt, īcēs, īcēt, īcētis, īcētūr, īcētī, īcimīnī, īcimur, īcimūs, īcīs, īcīsse, īcīssēmūs, īcīssētis, īcīssēs, īcīsset, īcīstī, īcīstis, īcīt, īcīte, īcītis, īcītō, īcītor, īcītōte, īcītūr, īcō, īcor, īcta, īctā, īctae, īctam, īctārum, īctās, īcte, īctī, īctīs, īctō, īctōrum, īctōs, īctū, īctūrum, īctūs, īcūnt, īcūntō, īcūntor, īcūntur — “strike”;*

*lēcta, lēctā, lēctae, lēctam, lēctārum, lēctās, lēcte, lēctī, lēctīs, lēctō, lēctōrum, lēctōs, lēctū, lēctum, lēctūrus, lēctus, legam, legāmīnī, legāmūr, legāmūs, legant, legantur, legar, legār, legāris, legās, legat, legātis, legātūr, lege, legēbam, legēbāmīnī, legēbāmūr, legēbāmūs, legēbānt, legēbāntur, legēbar, legēbāre, legēbāris, legēbās, legēbat, legēbātis, legēbātūr, legēminī, legēmūr, legēmūs, legendī, legendō, legendum, legendus, legēns, legent, legentur, lēgeram, lēgerāmūs, lēgerant, lēgerās, lēgerat, lēgerātis, legere, legēr, lēgerē, legerem, legerēmīnī, legerēmūr, legerēmūs, legerent, legerentur, legerer, legerēre, legerēris, legerēs, legeret, legerētis, legerētūr, lēgerim, lēgerimūs, lēgerimūs, lēgerint, legeris, legēris, lēgeris, lēgerīs, lēgerit, lēgeritis, lēgerītis, lēgerō, lēgerēunt, legēs, leget, legētis, legētūr, legī, lēgī, legimīnī, legimur, legimus, lēgimus, legis, lēgisse, lēgissēmūs, lēgissētis, lēgissēs, lēgissētis, lēgīstī, lēgīstis, legit, lēgit, legite, legitis, legitō, legitor, legitōte, legitūr, legō, legor, legunt, leguntō, leguntor, leguntur — “read”;*

*lēveram, lēverāmūs, lēverant, lēverās, lēverat, lēverātis, lēvēre, lēverim, lēverimūs, lēverint, lēveris, lēverīs, lēverit, lēveritis, lēverītis, lēverō, lēvērunt, lēvī, lēvīmus, lēvisse, lēvissēmūs, lēvissētis, lēvissēs, lēvissētis, lēvīstī, lēvīstis, lēvit, linam, lināmīnī, lināmūr, lināmūs, linant, linantur, linear, lināre, lināris, linās, linat, linātis, linātūr, line, linēbam, linēbāmīnī, linēbāmūr, linēbāmūs, linēbānt, linēbāntur, linēbar, linēbāre, linēbāris, linēbās, linēbat, linēbātis, linēbātūr, linēminī, linēmur, linēmus, linendī, linendō, linendum, linendus, linēns, linent, linentur, linere, linēre, linerem, linerēmīnī, linerēmūr, linerēmūs, linerent, linerentur, linerer, linerēre, linerēris, linerēs, lineret, linerētis, linerētūr, lineris, linēris, linēs, linet, linētis, linētūr, linī, linimīnī, linimur, linimus, linis, linit, linitis, linitō, linitor, linitōte, linitur, linō, linor, linunt, linuntō, linuntor, linuntur, litī, litum, litūrus, litus — “smear”;*

*locūta, locūtā, locūtae, locūtam, locūtarum, locūtās, locūte, locūtī, locūtīs, locūtō, locūtōrum, locūtōs, locūtū, locūtūm, locūtūrus, locūtūs, loquāminī, loquāmur, loquantur, loquar, loquāre, loquāris, loquātur, loquēbāminī, loquēbāmur, loquēbantur, loquēbar, loquēbāre, loquēbāris, loquēbātur, loquēminī, loquēmur, loquendī, loquendō, loquendum, loquendus, loquēns, loquentur, loquere, loquēre, loquerēminī, loquerēmur, loquerentur, loquerer, loquerēre, loquerēris, loquerētur, loqueris, loquēris, loquētur, loquī, loquimīnī, loquimur, loquitōr, loquitur, loquor, loquuntōr, loquuntur — “speak”;*

*messa, messā, messae, messam, messārum, messās, messe, messī, messīs, messō, messōrum, messōs, messū, messueram, messuerāmus, messuerant, messuerās, messuerat, messuerātis, messuēre, messuerim, messuerimus, messuerīmus, messuerint, messueris, messuerīs, messuerit, messueritis, messuerītis, messuerō, messuērunt, messū, messuīmus, messuisse, messuissem, messuīssemus, messuīssent, messuīsses, messuīsset, messuīssetis, messuīstī, messuīstis, messuit, messum, messūrus, messus, metam, metāmīnī, metāmūr, metāmūs, metant, metantur, metar, metārē, metāris, metās, metat, metātis, metātūr, mete, metēbam, metēbāmīnī, metēbāmūr, metēbāmūs, metēbant, metēbantur, metēbar, metēbāre, metēbāris, metēbās, metēbat, metēbātis, metēbātūr, metēminī, metēmur, metēmus, metendī, metendō, metendum, metendus, metēns, metent, metentur, metere, metērē, meterem, meterēmīnī, meterēmur, meterēmus, meterent, meterentur, meterer, meterērē, meterēris, meterēs, meteret, meterētis, mete-  
rētur, meteris, metēris, metēs, metet, metētis, metētūr, metī, metimīnī, metimur, metimus, metis, metit, metite, metitis, metitō, metitor, metitōte, metitur, metō, metor, metunt, metuntō, metuntor, metuntur — “reap”;*

*miseram, miseramus, miserant, miserās, miserat, miseratis, misere, miserim, miserimus, miserīmus, miserint, miseris, miserīs, miserit, miseritis, miserītis, miserō, miserunt, misī, misimus, misisse, misissēmus, misissent, misissēs, misisset, misissētis, misistī, misistis, misit, missa, missā, missae, missam, missārum, missās, misse, missī, missīs, missō, missōrum, missōs, missū, missum, missūrus, missus, mittam, mittāmī, mittāmur, mittāmus, mittant, mittantur, mittar, mittāre, mittāris, mittās, mittat, mittātis, mittātur, mitte, mittēbam, mittēbāmī, mittēbāmur, mittēbāmus, mittēbānt, mittēbāntur, mittēbar, mittēbare, mittēbāris, mittēbās, mittēbat, mittēbātis, mittēbātur, mittēmī, mittēmur, mittēmus, mittendī, mittendō, mittendum, mittendus, mittēns, mittent, mitten-  
tur, mittere, mittēre, mitterem, mitterēmī, mitterēmur, mitterēmus, mitterent, mitterentur, mitterer, mitterēre, mitterēris, mitterēs, mitteret, mitterētis, mitterētur, mitteris, mittēs, mittet, mittētis, mittētur, mittī, mittīmī, mittimur, mittimus, mittis, mittit, mittite, mittitis, mittitō, mittitor, mittitōte, mittitur, mittō, mittor, mittunt, mittuntō, mittuntor, mittuntur — “send”;*

*nōscam, nōscāminī, nōscāmur, nōscāmus, nōscant, nōscantur, nōscar, nōscāre, nōscāris, nōscās, nōscat, nōscātis, nōscātūr, nōsce, nōscēbam, nōscēbāminī, nōscēbāmūr, nōscēbāmus, nōscēbānt, nōscēbāntur, nōscēbar, nōscēbāre, nōscēbāris, nōscēbās, nōscēbat, nōscēbātis, nōscēbātūr, nōscēminī, nōscēmūr, nōscēmus, nōscēdī, nōscēdō, nōscēdūm, nōscēdūs, nōscēns, nōscēnt, nōscēntur, nōscēre, nōscērē, nōscērem, nōscēremīnī, nōscērēmūr, nōscērēmūs, nōscērent, nōscērentur, nōscērer, nōscērērē, nōscērēris, nōscērēs, nōscēret, nōscērētis, nōscērētūr, nōscēris, nōscēs, nōscet, nōscētis, nōscētūr, nōscī, nōscimīnī, nōscimūr, nōscimus, nōscis, nōscit, nōscite, nōscitis, nōscitō, nōscitor, nōscitōte, nōscitūr, nōscō, nōscor, nōscunt, nōscuntō, nōscuntōr, nōscuntur, nōta, nōtā, nōtae, nōtam, nōtārum, nōtās, nōte, nōtī, nōtīs, nōtō, nōtōrum, nōtōs, nōtū, nōtūm, nōtūrus, nōtus, nōveram, nōverāmus, nōverant, nōverās, nōverat, nōverātis, nōvēre, nōverim, nōverimus, nōverītis, nōverint, nōveris, nōverīs, nōverit, nōveritīs, nōverītis, nōverō, nōvērunt, nōtī, nōvīmus, nōvisse, nōvissem, nōvissēmus, nōvissent, nōvissēs, nōvisset, nōvissētis, nōvistī, nōvistis, nōvit — “know”;*

pāscam, pāscāminī, pāscāmūr, pāscāmus, pāscant, pāscantur, pāscar, pāscāre, pāscāris, pāscās, pāscat, pāscātis, pāscātūr, pāscē, pāscēbam, pāscēbāminī, pāscēbāmūr, pāscēbāmus, pāscēbānt, pāscēbāntur, pāscēbar, pāscēbārē, pāscēbāris, pāscēbās, pāscēbat, pāscēbātis, pāscēbātūr, pāscēminī, pāscēmūr, pāscēmus, pāscēndī, pāscēndō, pāscēndum, pāscēndus, pāscēns, pāscēnt, pāscēntur, pāscēre, pāscēre, pāscērem, pāscēremīnī, pāscēremur, pāscēremūs, pāscērent, pāscērentur, pāscērer, pāscērērē, pāscērēris, pāscērēs, pāscēret, pāscērētīs, pāscērētūr, pāscēris, pāscēris, pāscēs, pāscet, pāscētīs, pāscētūr, pāscētī, pāscimīnī, pāscimūr, pāscimūs, pāscis, pāscit, pāscite, pāscitis, pāscitō, pāscitor, pāscitōtē, pāscitūr, pāscō, pāscor, pāscunt, pāscuntō, pāscuntor, pāscuntur, pasta, pastā, pastae, pastam, pastārum, pastās, paste, pastī, pastīs, pastō, pastōrum, pastōs, pastū, pastum, pastūrus, pastus, pāveram, pāverāmus, pāverant, pāverās, pāverat, pāverātīs, pāvēre, pāverim, pāverimūs, pāverīmus, pāverint, pāveris, pāverīs, pāverit, pāveritis, pāverītīs, pāverō, pāverunt, pāvī, pāvīmus, pāvisse, pāvissem, pāvissēmus, pāvissent, pāvissēs, pāvisset, pāvissētīs, pāvistī, pāvistis, pāvit — “feed”;

*pellam, pellāminī, pellāmur, pellāmus, pellant, pellantur, pellar, pellāre, pellāris, pellās, pellat, pellātis, pellātur, pelle, pellēbam, pellēbāminī, pellēbāmur, pellēbāmus, pellēbant, pellēbantur, pellēbar, pellēbāre, pellēbāris, pellēbās, pellēbat, pellēbātis, pellēbātūr, pellēminī, pellēmur, pellēmus, pellēndī, pellēndō, pellēndum, pellēndus, pellēns, pellēnt, pellēntur, pellēre, pellēre, pellerem, pellerēminī, pellerēmur, pellerēmus, pellerent, pellerentur,*

*pellerer, pellerēre, pellerēris, pellerēs, pelleret, pellerētis, pellerētur, pelleris, pellēris, pelles, pellet, pellētis, pellētur, pelli, pelliminī, pellimur, pellimus, pellis, pellit, pellite, pellitis, pellitō, pellitor, pellitōtē, pellitur, pellō, pellor, pellunt, pelluntō, pelluntor, pelluntur, pepuleram, pepulerāmus, pepulerant, pepulerās, pepulerat, pepulerātis, pepulēre, pepulerim, pepulerimus, pepulerīmus, pepulerint, pepuleris, pepulerīs, pepulerit, pepuleritis, pepulerītis, pepulerō, pepulērunt, pepulī, pepulimus, pepulisse, pepulissēm, pepulissent, pepulisēs, pepulisset, pepulissētis, pepulistī, pepulistis, pepulit, pulsā, pulsae, pulsam, pulsārum, pulsās, pulse, pulsī, pulsīs, pulsō, pulsōrum, pulsōs, pulsū, pulsum, pulsūrus, pulsus — “beat”;*

*petam, petāminī, petāmur, petāmus, petant, petantur, petar, petāre, petāris, petās, petat, petātis, petātur, pete, petēbam, petēbāminī, petēbāmur, petēbāmus, petēbant, petēbantur, petēbar, petēbāre, petēbāris, petēbās, petēbat, petēbātis, petēbātur, petēminī, petēmur, petēmus, petendī, petendō, petendum, petendus, petēns, petent, petentur, petere, petēre, peterem, peterēminī, peterēmur, peterēmus, peterent, peterentur, peterer, peterēre, peterēris, peterēs, peteret, peterētis, peterētur, peteris, petēris, petēs, petet, petētis, petētur, petī, petimintī, petimur, petimus, petis, petit, petīta, petītā, petītae, petītam, petītarum, petītās, petite, petīte, petītī, petitis, petītīs, petitō, petītō, petitor, petītōrum, petītōs, petitōte, petītū, petītūm, petitur, petītūrus, petītūs, petīveram, petīverāmus, petīverant, petīverās, petīverat, petīverātis, petīvēre, petīverim, petīverimus, petīverīmus, petīverint, petīveris, petīverīs, petīverit, petīveritis, petīverītis, petīverō, petīvērunt, petīvī, petīvīmus, petīvisse, petīvissem, petīvissēmus, petīvissent, petīvissēs, petīvissset, petīvissētis, petīvistī, petīvistis, petīvit, petō, petor, petunt, petuntō, petuntor, petuntur — “seek”;*

sata, satā, satae, satam, satārum, satās, sate, satī, satīs, satō, satōrum, satōs, satū, satum, satūrus, satus, seram, serāminī, serāmur, serāmus, serant, serantur, serar, serāre, serāris, serās, serat, serātis, serātur, sere, serēbam, serēbāminī, serēbāmur, serēbāmus, serēbant, serēbantur, serēbar, serēbāre, serēbāris, serēbās, serēbat, serēbātis, serēbātur, serēminī, serēmur, serēmus, serendī, serendō, serendum, serendus, serēns, serent, serentur, serere, serēre, sererem, sererēminī, sererēmur, sererēmus, sererent, sererentur, sererer, sererēre, sererēris, sererēs, sereret, sererētis, sererētut, sereris, serēris, serēs, seret, serētis, serētut, serī, serimīnī, serimur, serimus, seris, serit, serite, seritis, seritō, seritor, seritōte, seritut, serō, seror, serunt, seruntō, seruntor, seruntur, sēveram, sēverāmus, sēverant, sēverās, sēverat, sēverātis, sēvēre, sēverim, sēverimus, sēverītīs, sēverint, sēveris, sēvērīs, sēverit, sēveritis, sēverītīs, sēverō, sēvērunt, sēvī, sēvīmus, sēvisse, sēvissem, sēvissēmus, sēvissent, sēvissēs, sēvisset, sēvissētis, sēvīstī, sēvītis, sēvīt — “sow”;

*tācta, tāctā, tāctae, tāctam, tāctārum, tāctās, tācte, tāctī, tāctīs, tāctō, tāctōrum, tāctōs, tāctū, tāctum, tāctūrus, tāctus, tangam, tangāminī, tangāmur, tangāmus, tangant, tangantur, tangar, tangāre, tangāris, tangās, tangat, tangātis, tangātur, tange, tangēbam, tangēbāminī, tangēbāmur, tangēbāmus, tangēbant, tangēbantur, tangēbar, tangēbāre, tangēbāris, tangēbās, tangēbat, tangēbātis, tangēbātur, tangēminī, tangēmur, tangēmus, tangendī, tangendō, tangendum, tangendus, tangēns, tangent, tangentur, tangere, tangēre, tangerem, tangerēminī, tangērēmur, tangerēmus, tangerent, tangerentur, tangerer, tangerēre, tangerēris, tangerēs, tangeret, tangerētis, tangērētur, tangeris, tangēris, tangēs, tanget, tangētis, tangētur, tangī, tangimīnī, tangimur, tangimus, tangis, tangit, tangite, tangitis, tangitō, tangitor, tangitōte, tangitur, tangō, tangor, tangunt, tanguntō, tanguntor, tanguntur — “touch”;*

*teram, terāminī, terāmūr, terāmūs, terant, terantur, terar, terāre, terāris, terās, terat, terātis, terātūr, tere, terēbam, terēbāminī, terēbāmūr, terēbāmūs, terēbant, terēbantur, terēbar, terēbāre, terēbāris, terēbās, terēbat, terēbātis, terēbātūr, terēminī, terēmūr, terēmūs, terendī, terendō, terendum, terendus, terēns, terent, terentur, terere, terēre, tererem, tererēminī, tererēmūr, tererēmūs, tererent, tererentur, tererer, tererērē, tererēris, tererēs, tereret, tererētis, tererētūr, tereris, terēris, terēs, teret, terētis, terētūr, terī, terimīnī, terimur, terimus, teris, terit, terite, teritis, teritō, teritor, teritōtē, teritur, terō, teror, terunt, teruntō, terunctor, teruntur, trīta, trītā, trītāe,*

*tritam, tritārum, tritās, trīte, trītī, trītīs, trītō, trītōrum, trītōs, trītū, trītūm, trītūrus, trītūs, trīveram, trīverāmus, trīverant, trīverās, trīverat, trīverātis, trīvēre, trīverim, trīverimus, trīverīmus, trīverint, trīveris, trīverīs, trīverit, trīveritis, trīverītis, trīverō, trīvērunt, trīvī, trīvīmus, trīvisse, trīvissem, trīvissēmus, trīvissent, trīvissēs, trīvisset, trīvissētis, trīvistī, trīvistis, trīvit — “triturate”;*

*texam, texāminī, texāmūr, texāmus, texant, texantur, texar, texāre, texāris, texās, texat, texātis, texātur, texēbam, texēbāminī, texēbāmūr, texēbāmus, texēbant, texēbantur, texēbar, texēbāre, texēbāris, texēbās, texēbat, texēbātis, texēbātur, texēminī, texēmur, texēmus, texendī, texendō, texendum, texendus, texēns, texent, texentur, texere, texēre, texerem, texerēminī, texerēmūr, texerēmus, texerent, texerentur, texerer, texerēre, texerēris, texerēs, texeret, texerētis, texerētur, texeris, texēris, texēs, texet, texētis, texētur, texī, teximinī, teximur, teximus, texis, texit, texite, texitis, texitō, texitor, texitōte, texitur, texō, texor, texta, textā, textae, textam, textārum, textās, texte, textī, textīs, textō, textōrum, textōs, textū, textum, textūrus, textus, texueram, texuerāmus, texuerant, texuerās, texuerat, texuerātis, texuēre, texuerim, texuerimus, texuerīmus, texuerint, texueris, texuerīs, texuerit, texueritis, texuerītis, texuerō, texuērunt, texuī, texuimus, texuisse, texuissem, texuissēmus, texuissent, texuissēs, texuisset, texuissētis, texuistī, texuistis, texuit, texunt, texuntō, texuntor, texuntur — “weave”;*

*tracta, tractā, tractae, tractam, tractārum, tractās, tracte, tractī, tractīs, tractō, tractōrum, tractōs, tractū, tractum, tractūrus, tractus, traham, trahāminī, trahāmur, trahāmus, trahant, trahantur, trahar, trahāre, trahāris, trahās, trahat, trahātis, trahātur, trahe, trahēbam, trahēbāminī, trahēbāmur, trahēbāmus, trahēbant, trahēban-  
tur, trahēbar, trahēbāre, trahēbāris, trahēbās, trahēbat, trahēbātis, trahēbātur, trahēminī, trahēmur, trahēmus, trahendī, trahendō, trahendum, trahendus, trahēns, trahent, trahentur, trahere, trahēre, traherem, traherēminī, traherēmur, traherēmus, traherent, traherentur, traherer, traherēre, traherēris, traherēs, traheret, traherētis, traherēt-  
ur, traheris, trahēris, trahēs, trahet, trahētis, trahētur, trahī, trahimīnī, trahimur, trahimus, trahis, trahit, trahite, trahitis, trahitō, trahitor, trahitōte, trahitur, trahō, trahor, trahunt, trahuntō, trahuntor, trahuntur —  
“drag”;*

*versa, versā, versae, versam, versārum, versās, verse, versī, versīs, versō, versōrum, versōs, versū, versum, ver-  
sūrus, versus, vertam, vertāminī, vertāmur, vertāmus, vertant, vertantur, vertar, vertāre, vertāris, vertās, vertat,  
vertātis, vertātur, verte, vertēbam, vertēbāminī, vertēbāmur, vertēbāmus, vertēbant, vertēbantur, vertēbar, ver-  
tēbare, vertēbāris, vertēbās, vertēbat, vertēbātis, vertēbātūr, vertēminī, vertēmūr, vertēmus, vertendī, vertendō,  
vertendum, vertendus, vertēns, vertent, vertentur, verteram, verterāmus, verterant, verterās, verterat, verterā-  
tis, vertere, vertēre, verterem, verterēminī, verterēmūr, verterēmus, verterent, verterentur, verterer, verterēre,  
verterēris, verterēs, verteret, verterētis, verterētūr, verterim, verterimus, verterīmus, verterint, verteris, verterīs,  
vertēris, verterit, verteritis, verterītis, verterō, vertērunt, vertēs, vertet, vertētis, vertētūr, vertē, vertimīnī, verti-  
mur, vertimus, vertis, vertisse, vertissem, vertissēmus, vertissent, vertissēs, vertisset, vertissētis, vertistī, vertistis,  
vertit, vertite, vertitis, vertitō, vertitor, vertitōte, vertitur, vertō, vertor, vertunt, vertuntō, vertuntor, vertuntur —  
“turn”;*

*vīceram, vīcerāmus, vīcerant, vīcerās, vīcerat, vīcerātis, vīcēre, vīcerim, vīcerimus, vīcerint, vīceris, vīcerīs, vīcerit, vīceritis, vīcerītis, vīcerō, vīcērunt, vīcī, vīcīmus, vīcīsse, vīcīssem, vīcīssēmus, vīcīssent, vīcīssēs, vīcīsset, vīcīssētis, vīcīstī, vīcīstis, vīcīt, vīcta, vīctā, vīctae, vīctam, vīctārum, vīctās, vīcte, vīctī, vīctīs, vīctō, vīctōrum, vīctōs, vīctū, victum, vīctūm, vīctūrus, vīctus, vīctus, vincam, vincāmīnī, vincāmūr, vincāmūs, vincānt, vincāntur, vincār, vincārē, vincārīs, vincās, vincāt, vincātis, vincātūr, vīnce, vīncebam, vīncebāmīnī, vīncebāmūr, vīncebāmūs, vīncebānt, vīncebāntur, vīncebār, vīncebārē, vīncebārīs, vīncebās, vīncebat, vīncebātis, vīncebātūr, vīncemīnī, vīncemūr, vīncemūs, vīcēndī, vīcēndō, vīcēndum, vīcēndus, vīcēns, vīcēnt, vīcēntur, vīcēntr, vīcēntr, vīcērē, vīcērēmīnī, vīcērēmūr, vīcērēmūs, vīcērentur, vīcērer, vīcērērē, vīcērērīs, vīcērēs, vīcēret, vīcērētis, vīcērētūr, vīcēris, vīcērīs, vīcēs, vīcēt, vīcētis, vīcētūr, vīcī, vīcīmīnī, vīcīmūr, vīcīmūs, vīcīs, vīcīt, vīcītē, vīcītūr, vīcītō, vīcītōr, vīcītōr, vīcītōr — “conquer”;*

*vīsa, vīsā, vīsae, vīsam, vīsāminī, vīsāmūr, vīsāmus, vīsānt, vīsāntur, vīsār, vīsāre, vīsāris, vīsārum, vīsās, vīsat, vīsātis, vīsātur, vīse, vīsēbam, vīsēbāminī, vīsēbāmūr, vīsēbāmus, vīsēbānt, vīsēbāntur, vīsēbar, vīsēbāre, vīsēbāris, vīsēbās, vīsēbat, vīsēbātis, vīsēminī, vīsēmur, vīsēmus, vīsendī, vīsendō, vīsendum, vīsendus, vīsēns, vīsent, vīsentur, vīseram, vīserāmus, vīserās, vīserat, vīserātis, vīsere, vīsēre, vīserem, vīserēminī, vīserēmur, vīserēmus, vīserent, vīserentur, vīserer, vīserēre, vīserēris, vīserēs, vīseret, vīserētis, vīseretur, vīserim, vīserimus, vīserīmus, vīserint, vīseris, vīserīs, vīserit, vīseritis, vīserītis, vīserō, vīserunt, vīses, vīset, vīsētis, vīsētur, vīstī, vīsiminī, vīsimur, vīsimus, vīsis, vīsts, vīsisse, vīsissem, vīsissēmus, vīsissent, vīsissēs, vīsisset, vīsissētis, vīsistī, vīsistis, vīsit, vīsite, vīsitis, vīsitō, vīsitor, vīsitōte, vīsitur, vīsō, vīsor, vīsōrum, vīsōs, vīsū, vīsum, vīsunt, vīsuntō, vīsuntor, vīsuntur, vīsūrus, vīsus — “visit”;*

vomam, vomāminī, vomāmur, vomāmus, vomant, vomantur, vomar, vomāre, vomāris, vomās, vomat, vomātis, vomātur, vomē, vomēbam, vomēbāminī, vomēbāmur, vomēbāmus, vomēbānt, vomēbāntur, vomēbar, vomēbāre, vomēbāris, vomēbās, vomēbat, vomēbātis, vomēbātūr, vomēminī, vomēmur, vomēmus, vomendī, vomendō, vomendum, vomendus, vomēns, voment, vomentur, vomere, vomerem, vomerēminī, vomerēmur, vomerēmus, vomerent, vomerentur, vomerer, vomerēre, vomerēris, vomerēs, vomeret, vomerētis, vomerētūr, vomeris, vomēris, vomēs, vomet, vomētis, vomētūr, vomī, vomimī, vomimur, vomimus, vomis, vomit, vomita, vomitā, vomitae, vomitam, vomitārum, vomitās, vomite, vomitī, vomitis, vomitīs, vomitō, vomitor, vomitōrum, vomitōs, vomitōte, vomitū, vomitum, vomitur, vomitūrus, vomitus, vomō, vomor, vomueram, vomuerāmus, vomuerant, vomuerās, vomuerat, vomuerātis, vomuēre, vomuerim, vomuerimus, vomuerīmus, vomuerint, vomueris, vomuerīs, vomuerit, vomueritis, vomuerītis, vomuerō, vomuērunt, vomū, vomuimus, vomuisse, vomuissem, vomuissēmus, vomuis-  
sent, vomuissēs, vomuisset, vomuissētis, vomuistī, vomuistis, vomuit, vomunt, vomuntō, vomuntor, vomuntur — “vomit”.

[https://en.wikipedia.org/wiki/Latin\\_conjugation#Fourth\\_conjugation](https://en.wikipedia.org/wiki/Latin_conjugation#Fourth_conjugation)

*aperī, aperiam, aperiāminī, aperiāmur, aperiāmus, aperiant, aperiantur, aperiar, aperiāre, aperiāris, aperiās, aperiat, aperiātis, aperiātūr, aperiēbam, aperiēbāminī, aperiēbāmur, aperiēbāmus, aperiēbant, aperiēbantur, aperiēbar, aperiēbāre, aperiēbāris, aperiēbās, aperiēbat, aperiēbātis, aperiēbātur, aperiēminī, aperiēmur, aperiēmus, aperiendī, aperiendō, aperiendum, aperiendus, aperiēns, aperient, aperientur, aperiēre, aperiēris, aperiēs, aperiet, aperiētis, aperiētūr, aperiēminī, aperiēmur, aperiēmus, aperiō, aperior, aperiēre, aperiērem, aperiēremīnī, aperiēremur, aperiēremus, aperiērent, aperiērentur, aperiērer, aperiērēre, aperiērēris, aperiērēs, aperiēret, aperiēretis, aperiēretur, aperiērī, aperiēris, aperiō, aperit, aperiēte, aperiētis, aperiētō, aperiōtor, aperiētōte, aperiētūr, aperiunt, aperiuntō, aperiuntor, aperiuntur, aperta, apertā, apertae, apertam, apertārum, apertās, aperte, aperiī, aperiīs, aperiōtō, aperiōrum, aperiōs, aperiū, aperitum, aperiūrus, apertus, aperueram, aperuerāmus, aperuerant, aperuerās, aperuerat, aperuerātis, aperuerēre, aperuerim, aperuerimus, aperuerīmus, aperuerint, aperueris, aperuerīs, aperuerit, aperueritis, aperuerītis, aperuerō, aperiuerērunt, aperiū, aperiūmus, aperiūsse, aperiūssem, aperiūssēmus, aperiūssent, aperiūssēs, aperiūsset, aperiūssētis, aperiūstī, aperiūstis, aperiuit — “open”;*

*audī, audīam, audiāminī, audiāmūr, audiāmūs, audīant, audīantur, audīar, audiāre, audiāris, audiās, audīat, audiātis, audiātur, audiēbam, audiēbāminī, audiēbāmūr, audiēbāmūs, audiēbant, audiēbantur, audiēbar, audiēbāre, audiēbāris, audiēbās, audiēbat, audiēbātis, audiēbātūr, audiēminī, audiēmūr, audiēmūs, audiendī, audiendō, audiendum, audiendus, audiēns, audīent, audīentur, audiēre, audiēris, audiēs, audīet, audiētis, audiētūr, audīminī, audīmūr, audīmūs, audiō, audīor, audīre, audīrem, audīremīnī, audīremūr, audīremūs, audīrent, audīrentur, audīrer, audīrēre, audīrēris, audīrēs, audīret, audīrētis, audīrētūr, audīrī, audīris, audīs, audit, audīta, audītā, audītāe, audītam, audītārum, audītās, audīte, audītī, audītis, audītīs, audītō, audītor, audītōrum, audītōs, audītōte, audītū, audītum, audītur, audītūrus, audītūs, audiunt, audiuntō, audiuntor, audiuntur, audīveram, audīverāmūs, audīverant, audīverās, audīverat, audīverātis, audīvēre, audīverim, audīverimus, audīverīmus, audīverint, audīveris, audīverīs, audīverit, audīveritis, audīverītis, audīverō, audīverunt, audīvī, audīvīmus, audīvisse, audīvissem, audīvissēmūs, audīvissent, audīvissēs, audīvisset, audīvissētis, audīvistī, audīvistis, audīvit — “hear”;*

*oriāminī, oriāmur, orientur, oriar, oriāre, oriāris, oriātur, oriēbāminī, oriēbāmur, oriēbantur, oriēbar, oriēbāre, oriēbāris, oriēbātūr, oriēminī, oriēmur, oriendī, oriendō, oriendum, oriendus, oriēns, orientur, oriēre, oriēris, oriētūr, orīminī, orīmur, orior, orīre, orīrēminī, orīrēmur, orīrentur, orīrer, orīrērē, orīrēris, orīrētūr, orīrī, orīris, orītor, orītūr, orītūrus, oriuntor, oriuntur, orta, ortā, ortae, ortam, ortārum, ortās, orte, ortī, ortīs, ortō, ortōrum, ortōs, ortū, ortum, ortus — “rise”;*

*sancī, sanciam, sanciāminī, sanciāmur, sanciāmus, sanciant, sanciantur, sanciar, sanciāre, sanciāris, sanciās, sanciat, sanciātis, sanciātur, sanciēbam, sanciēbāminī, sanciēbāmur, sanciēbāmus, sanciēbant, sanciēbantur, sanciēbar, sanciēbāre, sanciēbāris, sanciēbās, sanciēbat, sanciēbātis, sanciēbātur, sanciēminī, sanciēmur, sanciēmus, sanciēndī, sanciēndō, sanciēndum, sanciēndus, sanciēns, sancient, sancientur, sanciēre, sanciēris, sanciēs, sanciet, sanciētis, sanciētur, sanciēminī, sanciēmur, sanciēmus, sanciō, sancior, sanciēre, sanciērem, sanciērēmi-  
nī, sanciērēmur, sanciērēmus, sanciērent, sanciērentur, sanciērer, sanciērēre, sanciērēris, sanciērēs, sanciēret, sanciērētis, sanciērētur, sanciērī, sanciēris, sanciēs, sancit, sanciēte, sanciētis, sanciētō, sanciētor, sanciētōte, sanciētūr, sanciētūnt, sanciētūntor, sanciētūntur, sāncta, sānctā, sānctae, sānctam, sānctārum, sānctās, sāncte, sānctī, sānctīs, sānctō, sānctōrum, sānctōs, sānctū, sānctum, sānctūrus, sānctus, sānixeram, sānixerāmus, sānixerant, sānixerās, sānixerat, sānixerātis, sānixerē, sānixerim, sānixerimus, sānixerīmus, sānixerint, sānixeris, sānixerīs, sānixerit, sānxe-  
ritis, sānixerītis, sānixerō, sānixerunt, sānxit, sānxiimus, sānxisse, sānxissem, sānxiissēmus, sānxiissent, sānxiissēs, sānxiisset, sānxiissētis, sānxitī, sānxitis, sānxit — “confirm”;*

*vēneram, vēnerāmus, vēnerant, vēnerās, vēnerat, vēnerātis, vēnēre, vēnerim, vēnerimus, vēnerīmus, vēnerint, vēneris, vēnerīs, vēnerit, vēneritis, vēnerītis, vēnerō, vēnērunt, venī, vēnī, veniam, veniāmus, veniant, veniās, veniat, veniātis, veniātur, veniēbam, veniēbāmus, veniēbant, veniēbās, veniēbat, veniēbātis, veniēbātur, veniēmus, veniēndī, veniēndō, veniēndum, veniēndus, veniēns, venient, veniēs, veniet, veniētis, veniētūr, veniētūs, vēniēmus, vēniō, vēnēre, vēnērem, vēnērēmus, vēnērent, vēnērēs, vēnēret, vēnērētis, vēnērētūr, vēnērētūs, vēnēsse, vēnēssem, vēnissēmus, vēnissent, vēnissēs, vēnisset, vēnissētis, vēnistī, vēnistis, venit, vēnit, venīte, vēnītis, vēnītō, vēnītōte, vēnītur, vēniunt, vēniuntō, vēnta, vēntā, vēntae, vēntam, vēntārum, vēntās, vēnte, vēntī, vēntīs, vēntō, vēntōrum, vēntōs, vēntū, vēntum, vēntūrus, vēntus — “come”.*

As we throw each group of test verbs into our algorithm, the verb forms from the same conjugation groups are clustered correctly.

We note, however, that the inflected forms of *vetus* “old” and *vetō* “forbid” will be partly conflated if they appear simultaneously:

{vetā, vetābam, vetābāminī, vetābāmur, vetābāmus, vetābant, vetābantur, vetābar, vetābāre, vetābāris, vetābās,  
vetābat, vetābātis, vetābātur, vetābere, vetāberis, vetābimīnī, vetābimur, vetābimus, vetābis, vetābit, vetābitis,  
vetābitur, vetābō, vetābor, vetābunt, vetābuntur, vetāminī, vetāmūr, vetāmus, vetāndī, vetāndō, vetāndum, ve-  
tāndus, vetāns, vetānt, vetāntō, vetāntor, vetāntur, vetārē, vetārēmīnī, vetārēmūr, vetārēmūs, vetārent,  
vetārentur, vetārer, vetārērē, vetārērīs, vetārēs, vetāret, vetārētīs, vetārētūr, vetārī, vetārīs, vetās, vetat, vetāte,  
vetātīs, vetātō, vetātor, vetātōtē, vetātūr, vetem, vetēmīnī, vetēmūr, vetēmūs, vetent, vetentur, veter, vetera, ve-  
tere, vetēre, veterem, veterēs, veterī, veteribus, veteris, vetēris, vetērima, vetērimā, vetērimae, vetērimam,

veterimārum, veterimās, veterime, veterimī, veterimīs, veterimō, veterimōrum, veterimōs, veterimum, veterimus, veterum, vetēs, vetet, vetētis, vetētur, vetita, vetitā, vetitae, vetitam, vetitārum, vetitās, vetite, vetitī, vetitīs, vetitō, vetitōrum, vetitōs, vetitū, vetitum, vetitūrus, vetitus, vetō, vedor, vetueram, vetuerāmus, vetuerant, vetuerās, vetuerat, vetuerātis, vetuēre, vetuerim, vetuerimus, vetuerīmus, vetuerint, vetueris, vetuerīs, vetuerit, vetueritis, vetuerītis, vetuerō, vetuērunt, vetuī, vetuīmus, vetuisse, vetuissēmus, vetuissent, vetuissēs, vetuisset, vetuissētis, vetuistī, vetuistis, vetuit, vetus, vetustissimī, vetustissimōs},

{vetustior, vetustiōra, vetustiōre, vetustiōrem, vetustiōrēs, vetustiōrī, vetustiōribus, vetustiōris, vetustiōrum, vetustissima, vetustissimā, vetustissimae, vetustissimam, vetustissimārum, vetustissimās, vetustissime, vetustissimō, vetustissimōrum, vetustissimōs, vetustissimum, vetustissimus, vetustius}.

Such unfortunate conflations seem unavoidable. Observe that *veteris* is the genitive singular form of *vetus*, while *vetēris* is the second-person singular present passive subjunctive of *vetō*. Both *veteris* and *vetēris* share the same essential root *vet*, and they are practically indistinguishable without context, if the document in question does not employ the diacritical mark for long vowels. Some other Latin words also behave similarly: *portā* is both the ablative singular of *porta* “gate” and the singular present active imperative of *portō* “convey”.

*Example 6.29.3.* We pick a list of typical irregular verbs from the Wikipedia links indicated below.

[https://en.wikipedia.org/wiki/Latin\\_conjugation#Irregular\\_verbs](https://en.wikipedia.org/wiki/Latin_conjugation#Irregular_verbs)

[https://en.wikipedia.org/wiki/Latin\\_conjugation#Third\\_conjugation\\_.E2.80.93i.C5.8D\\_verbs](https://en.wikipedia.org/wiki/Latin_conjugation#Third_conjugation_.E2.80.93i.C5.8D_verbs)

*cape, capere, caperem, caperēminī, caperēmur, caperēmus, caperent, caperentur, caperer, caperēre, caperēris, caperēs, caperet, caperētis, caperētur, caperis, capī, capiam, capiāminī, capiāmur, capiāmus, capiant, capiantur, capiar, capiāre, capiāris, capiās, capiat, capiātis, capiātur, capiēbam, capiēbāminī, capiēbāmur, capiēbāmus, capiēbānt, capiēbāntur, capiēbar, capiēbāre, capiēbāris, capiēbās, capiēbat, capiēbātis, capiēbātūr, capiēminī, capiēmur, capiēmus, capiēndī, capiēndō, capiēndum, capiēndus, capiēns, capient, capientur, capiēre, capiēris, capiēs, capiet, capiētis, capiētur, capiēminī, capimur, capimus, capiō, capiōr, capis, capit, capite, capitis, capiōtō, capitor, capitōte, capitur, capiunt, capiuntō, capiuntor, capiuntur, capta, captā, captae, captam, captārum, captās, capte, captī, captīs, captō, captōrum, captōs, captū, captum, captūrus, captus, cēperam, cēperāmus, cēperant, cēperās, cēperat, cēperātis, cēpēre, cēperim, cēperimus, cēperīmus, cēperint, cēperis, cēperīs, cēperit, cēperitis, cēperītis, cēperō, cēpērunt, cēpī, cēpīmus, cēpisce, cēpissem, cēpissēmus, cēpissent, cēpissēs, cēpisset, cēpissētis, cēpīstī, cēpīstis, cēpīt — “capture”;*

*cupe, cupere, cuperem, cuperēminī, cuperēmur, cuperēmus, cuperent, cuperentur, cuperer, cuperēre, cuperēris, cuperēs, cuperet, cuperētis, cuperētur, cuperis, cupī, cupiam, cupiāminī, cupiāmur, cupiāmus, cupiant, cupiantur, cupiar, cupiāre, cupiāris, cupiās, cupiat, cupiātis, cupiātur, cupiēbam, cupiēbāminī, cupiēbāmur, cupiēbāmus, cupiēbānt, cupiēbāntur, cupiēbar, cupiēbāre, cupiēbāris, cupiēbās, cupiēbat, cupiēbātis, cupiēbātūr, cupiēmi-nī, cupiēmur, cupiēmus, cupiēndī, cupiēndō, cupiēndum, cupiēndus, cupiēns, cupient, cupientur, cupiēre, cupiēris, cupiēs, cupiet, cupiētis, cupiētur, cupiēminī, cupimur, cupimus, cupiō, cupiōr, cupis, capit, capite, capitis, cupiōtō, capitor, capitōte, capitur, capiunt, capiuntō, capiuntor, capiuntur, capiūverāram, capiūverāmus, capiūverānt, capiūverās, capiūverāt, capiūverātis, capiūverēre, capiūverim, capiūverimus, capiūverīmus, capiūverint, capiūveris, capiūverīs, capiūverit, capiūverītis, capiūverō, capiūverērunt, capiūverī, capiūvīmus, capiūvisse, capiūvissem, capiūvissēmus, capiūvissent, capiūvissēs, capiūvissētis, capiūvīstī, capiūvīstis, capiūvīt — “desire”;*

*eam, eamus, eant, eas, eat, eatis, eatur, eo, eundi, eundum, eundus, eunt, eunto, i, ibam, ibamus, ibant, ibas, ibat, ibatis, ibatur, ibimus, ibis, ibit, ibitis, ibitur, ibo, ibunt, iens, ieram, ieramus, ierant, ieras, ierat, ieratis, iere, ierim, ierimus, ierint, ieris, ierit, ieritis, iero, ierunt, ii, iīmus, iīt, imus, ire, irem, iremus, irent, ires, iret, iretis, iretur, iri, is, isse, issē, issēmus, issent, isses, issēt, issētis, istī, istis, it, ite, itis, ito, itote, itu, itum, itur, iturus, itus, ivi, ivisti, ivit — “go”;*

*edam, edāminī, edāmur, edāmus, edant, edantur, edar, edāre, edāris, edās, edat, edātis, edātur, ede, edēbam, edē-bāminī, edēbāmur, edēbāmus, edēbānt, edēbāntur, edēbar, edēbāre, edēbās, edēbat, edēbātis, edēbātūr, edēminī, edēmur, edēmus, edendī, edendō, edendum, edendus, edēns, edent, edentur, ēderam, ēderāmus, ēderant, ēderās, ēderat, ēderātis, edere, edēre, ēdēre, ederem, ederēminī, ederēmur, ederēmus, ederent, ederentur, ederer, ederēre, ederēris, ederēs, ederet, ederētis, ederētur, ēderim, ēderimus, ēderīmus, ēderint, ederis, edēris, ēderīs, ēderit, ēderitis, ēderītis, ēderō, ēdērunt, edēs, edet, edētis, edētur, edī, ēdī, edimur, edimus,*

*edīmus, ēdimus, edint, edis, edīs, ēdisse, ēdissem, ēdissēmus, ēdissent, ēdissēs, ēdissētis, ēdistī, ēdistis, edit, ēdit, edite, editis, editīs, editō, editor, editōte, editur, edō, edor, edunt, eduntō, eduntor, eduntur, ēs, ēsa, ēsā, ēsae, ēsam, ēsārum, ēsās, ēse, ēsī, ēsīs, ēsō, ēsōrum, ēsōs, esse, ēsse, ēssem, ēssēmus, ēssent, ēssēs, ēsset, ēssētis, ēst, ēste, ēstis, ēstō, ēstōte, ēstur, ēsū, ēsum, ēsūrus, ēsus — “eat”;*

*fer, feram, ferāminī, ferāmur, ferāmus, ferant, ferantur, ferar, ferāre, ferāris, ferās, ferat, ferātis, ferātur, ferēbam, ferēbāminī, ferēbāmur, ferēbāmus, ferēbant, ferēbantur, ferēbar, ferēbāre, ferēbāris, ferēbās, ferēbat, ferēbātis, ferēbātūr, ferēminī, ferēmur, ferēmus, ferendī, ferendō, ferendum, ferendus, ferēns, ferent, ferentur, ferēre, ferēris, ferēs, feret, ferētis, ferētūr, ferimīnī, ferimur, ferimus, ferō, feror, ferre, ferrem, ferrēminī, ferrēmur, ferrēmus, ferrent, ferrentur, ferrer, ferrēre, ferrēris, ferrēs, ferret, ferrētis, ferrētūr, ferrī, ferris, fers, fert, ferte, fertis, fertō, fertor, fertōte, fertur, ferunt, ferunto, feruntor, feruntur, lāta, lātā, lātae, lātam, lātārum, lātās, lāte, lāti, lātis, lātō, lātōrum, lātōs, lātū, lātum, lātūrus, lātus, tuleram, tulerāmus, tulerant, tulerās, tulerat, tulerātis, tulēre, tulerim, tulerimus, tulerāmus, tulerint, tuleris, tulerīs, tulerit, tuleritis, tulerītis, tulerō, tulērunt, tulī, tulimus, tulisse, tulissem, tulissēmus, tulissent, tulissēs, tulisset, tulissētis, tulistī, tulistis, tulit — “carry”;*

*mālam, mālēbam, mālēbāmus, mālēbant, mālēbās, mālēbat, mālēbātis, mālēmus, mālēnt, mālēs, mālet, mālētis, mālim, mālīmus, mālint, mālīs, mālit, mālītis, mālle, māllem, māllēmus, māllēnt, māllēs, māllet, māllētis, mālō, mālueram, māluerāmus, māluerant, māluerās, māluerat, māluerātis, māluēre, māluerim, māluerimus, māluerīmus, māluerint, mālueris, māluerīs, māluerit, mālueritis, māluerītis, māluerō, māluērunt, māluī, māluīmus, māluisse, māluissem, māluissēmus, māluissent, māluissēs, māluisset, māluissētis, māluistī, māluistis, māluit, mālumus, mālunt, māvīs, māvult, māvultis — “prefer”;*

*morere, morerēminī, morerēmur, morerentur, morerer, morerēre, morerēris, morerētur, moreris, morī, moriāminī, moriāmur, moriantur, moriar, moriāre, moriāris, moriātur, moriēbāminī, moriēbāmur, moriēbantur, moriēbar, moriēbāre, moriēbāris, moriēbātur, moriēminī, moriēmur, moriendī, moriendō, moriendum, moriendus, moriēns, morientur, moriēre, moriēris, moriētur, morimīnī, morimur, morior, moritor, moritur, moritūrus, moriuntor, moriuntur, mortua, mortuā, mortuae, mortuam, mortuārum, mortuās, mortue, mortuī, mortuīs, mortuō, mortuōrum, mortuōs, mortuū, mortuum, mortuus — “die”;*

*nōlam, nōlēbam, nōlēbāmus, nōlēbant, nōlēbās, nōlēbat, nōlēbātis, nōlēmus, nōlēns, nōlēnt, nōlēs, nōlēt, nōlētis, nōlēt, nōlim, nōlēmus, nōlēnt, nōlēts, nōlēlit, nōlēlite, nōlēltis, nōlēlitō, nōlēlitōte, nōlēlle, nōlēlem, nōlēlēmus, nōlēllent, nōlēlēs, nōlēlēt, nōlēlētis, nōlēlō, nōlēueram, nōlēuerāmus, nōlēuerant, nōlēuerās, nōlēuerat, nōlēuerātis, nōlēuerē, nōlēuerim, nōlēuerimus, nōlēuerīmus, nōlēuerint, nōlēueris, nōlēuerīs, nōlēuerit, nōlēueritis, nōlēuerītis, nōlēuerō, nōlēuerunt, nōlēluīt, nōlēluīmus, nōlēuisse, nōlēuissem, nōlēuiressēmus, nōlēuisseent, nōlēuiressēs, nōlēuisset, nōlēuiressētis, nōlēuiressēt, nōlēuistīt, nōlēuistis, nōlēuit, nōlēumus, nōlēunt, nōlēuntō — “refuse”;*

*passa, passā, passae, passam, passārum, passās, passe, passī, passīs, passō, passōrum, passōs, passū, passum, passūrus, passus, patere, paterēminī, paterēmur, paterentur, paterer, paterēre, paterēris, paterētur, pateris, patī, patiāminī, patiāmur, patientur, patiar, patiāre, patiāris, patiātur, patiēbāminī, patiēbāmur, patiēbāntur, patiēbar, patiēbāre, patiēbāris, patiēbātur, patiēminī, patiēmur, patientdī, patientdō, patientdum, patientdus, patientēns, patientur, patiēre, patiēris, patiētur, patimīnī, patimur, patior, patitor, patitūr, patiuntor, patiuntur — “suffer”;*

*rape, rapere, raperem, raperēminī, raperēmur, raperēmus, raperent, raperentur, raperer, raperēre, raperēris, raperēs, raperet, raperētis, raperētur, raperis, rapī, rapiam, rapiāminī, rapiāmur, rapiāmus, rapiant, rapiantur, rapiar, rapiāre, rapiāris, rapiās, rapiat, rapiātis, rapiātur, rapiēbam, rapiēbāminī, rapiēbāmur, rapiēbāmus, rapiēbant, rapiēbantur, rapiēbar, rapiēbāre, rapiēbāris, rapiēbās, rapiēbat, rapiēbātis, rapiēbātur, rapiēminī, rapiēmur, rapiēmus, rapiendī, rapiendō, rapiendum, rapiendus, rapiēns, rapient, rapientur, rapiēre, rapiēris, rapiēs, rapiet, rapiētis, rapiētur, rapimīnī, rapimur, rapimus, rapiō, rapior, rapis, rapit, rapite, rapitis, rapiōtō, rapitor, rapitōte, rapitur, rapiunt, rapiuntō, rapiuntor, rapiuntur, rapta, rapiā, raptae, raptam, raptārum, raptās, rapte, raptī, raptīs, raptō, raptōrum, raptōs, raptū, raptum, raptūrus, raptus, rapueram, rapuerāmus, rapuerant, rapuerās, rapuerat, rapuerātis, rapuēre, rapuerim, rapuerimus, rapuerīmus, rapuerint, rapueris, rapuerīs, rapuerit, rapueritis, rapuerītis, rapuerō, rapuērunt, rapuī, rapuimus, rapuisse, rapuisse, rapuissēmus, rapuissent, rapuissēs, rapuisset, rapuissētis, rapuistī, rapuistis, rapuit — “snatch”;*

*velim, velīmus, velint, velīs, velit, velītis, velle, vellem, vellēmus, vellent, vellēs, vellet, vellētis, vīs, volam, volēbam, volēbāmus, volēbant, volēbās, volēbat, volēbātis, volēmus, volēns, volent, volēs, volet, volētis, volō, volt, voltis, volueram, voluerāmus, voluerant, voluerās, voluerat, voluerātis, voluēre, voluerim, voluerimus, voluerīmus, voluerint, volueris, voluerīs, voluerit, volueritis, voluerītis, voluerō, voluērunt, volū, voluīmus, voluisse, voluissem, voluissēmus, voluissent, voluissēs, voluisset, voluissētis, voluistī, voluistis, voluit, volumus, volunt, vult, vultis* — “want”.

The clustering result is satisfactory:

{*cape, capere, caperem, caperēminī, caperēmur, caperēmus, caperent, caperentur, caperer, caperēre, caperēris, caperēs, caperet, caperētis, caperētur, caperis, capī, capiam, capiāminī, capiāmur, capiāmus, capiant, capiantur, capiar, capiāre, capiāris, capiās, capiat, capiātis, capiātur, capiēbam, capiēbāminī, capiēbāmur, capiēbāmus, capiēbant, capiēbantur, capiēbar, capiēbāre, capiēbāris, capiēbās, capiēbat, capiēbātis, capiēbātūr, capiēminī, capiēmur, capiēmus, capiēndī, capiēndō, capiēndum, capiēndus, capiēns, capient, capientur, capiēre, capiēris, capiēs, capiet, capiētis, capiētur, capiēndī, capimur, capimus, capiō, capior, capis, capit, capite, capitis, capiōtō, capitor, capitōte, capitūr, capiunt, capiuntō, capiuntor, capiuntur, capta, captā, captae, captam, captārum, captās, capte, captī, captīs, captō, captōrum, captōs, captū, captum, captūrus, captus, cēperam, cēperāmus, cēperant, cēperās, cēperat, cēperātis, cēpēre, cēperim, cēperimus, cēperint, cēperis, cēperīs, cēperit, cēperitis, cēperītis, cēperō, cēpērunt, cēpī, cēpīmus, cēpisse, cēpissēm, cēpissēmus, cēpissēt, cēpissēs, cēpisset, cēpissētis, cēpistī, cēpistis, cēpit},*

{*cupe, cupere, cuperem, cuperēminī, cuperēmur, cuperēmus, cuperent, cuperentur, cuperer, cuperēre, cuperēris, cuperēs, cuperet, cuperētis, cuperētur, cuperis, cupī, cupiam, cupiāminī, cupiāmur, cupiāmus, cupiant, cupiantur, cupiar, cupiāre, cupiāris, cupiās, cupiat, cupiātis, cupiātur, cupiēbam, cupiēbāminī, cupiēbāmur, cupiēbāmus, cupiēbant, cupiēbantur, cupiēbar, cupiēbāre, cupiēbāris, cupiēbās, cupiēbat, cupiēbātis, cupiēbātūr, cupiēminī, cupiēmur, cupiēmus, cupiēndī, cupiēndō, cupiēndum, cupiēndus, cupiēns, cupient, cupientur, cupiēre, cupiēris, cupiēs, cupiet, cupiētis, cupiētur, cupiēndī, cupimur, cupimus, cupiō, cupior, cupis, cupit, cupīta, cupītā, cupītae, cupītam, cupītārum, cupītās, cupite, cupītī, cupitis, cupītīs, cupitō, cupītō, cupitor, cupītōrum, cupītōs, cupitōte, cupītū, cupītūm, cupitūr, cupītūs, cupiunt, cupiuntō, cupiuntor, cupiuntur, cupīveram, cupīverāmus, cupīverant, cupīverās, cupīverat, cupīverātis, cupīvēre, cupīverim, cupīverimus, cupīverītis, cupīverītis, cupīverītis, cupīverītis, cupīverō, cupīverūnt, cupīvī, cupīvīmus, cupīvīsse, cupīvīssem, cupīvīssēm, cupīvīssēt, cupīvīssēs, cupīvīstī, cupīvīstis, cupīvīt},*

{*eam, eamus, eant, eas, eat, eatis, eatur, eo, ēs, ēsa, ēsā, ēsae, ēsam, ēsī, ēsō, ēsū, eundi, eundo, eundum, eundus, eunt, eunto, i, ibam, ibamus, ibant, ibas, ibat, ibatis, ibatur, ibimus, ibis, ibit, ibitis, ibitur, ibo, ibunt, iens, ieram, ieramus, ierant, ieras, ierat, ieratis, iere, ierim, ierimus, ierint, ieris, ierit, ieritis, iero, ierunt, ii, iimus, iit, imus, ire, irem, iremus, irent, ires, iret, iretis, iretūr, iri, is, isse, issem, issemus, issent, isses, isset, issetis, isti, istis, it, ite, itis, ito, itote, itu, itum, itur, iturus, itus, ivi, ivisti, ivit},*

{*edam, edāminī, edāmur, edāmus, edant, edantur, edar, edāre, edās, edat, edātis, edātūr, ede, edēbam, edēbāminī, edēbāmur, edēbāmus, edēbant, edēbantur, edēbar, edēbāre, edēbās, edēbat, edēbātis, edēbātūr, edēminī, edēmur, edēmus, edendī, edendō, edendum, edendus, edēns, edent, edentur, ēderam, ēderāmus, ēderant, ēderās, ēderat, ēderātis, edere, edēre, ēdēre, ederem, ederēminī, ederēmur, ederēmus, ederent, ederentur, ederer, ederēre, ederēris, ederēs, ederet, ederētis, ederētūr, ēderim, ēderīmus, ēderint, ederis, edēris, ēderīs, ēderit, ēderitis, ēderō, ēderunt, edēs, edet, edētis, edētūr, edī, ēdī, edim, ediminī, edimur, edimus, edīmus, ēdimus, edint, edis, edīs, ēdissem, ēdissēm, ēdissēt, ēdissētis, ēdistī, ēdistis, edit, ēedit, edite, editis, editō, editor, editōte, editur, edō, edor, edunt, eduntō, eduntor, eduntur, ēsārum, ēsās, ēsē, ēsī, ēsō, ēsōrum, ēsōs, esse, ēsse, ēssem, ēsēm, ēsēt, ēsētis, ēst, ēste, ēstis, ēstō, ēstōte, ēstur, ēsum, ēsūrus, ēsus},*

{*fer, feram, ferāminī, ferāmur, ferāmus, ferant, ferantur, ferar, ferāre, ferāris, ferās, ferat, ferātis, ferātūr, ferēbam, ferēbāminī, ferēbāmur, ferēbāmus, ferēbant, ferēbāre, ferēbāris, ferēbās, ferēbat, ferēbātis, ferēbātūr, ferēminī, ferēmur, ferēmus, ferēndī, ferēndō, ferēndum, ferēndus, ferēns, ferēnt, ferēntur, ferērē, ferēris, ferēs, feret, ferētis, ferētūr, ferēmī, ferēmīs, ferēmīs, ferō, feror, ferre, ferrem, ferrēmī, ferrēmīs, ferrent, ferrentur, ferrērē, ferrēris, ferrēs, ferret, ferrētis, ferrētūr, ferrī, ferrīs, fers, fert, ferte, fertis, fertō, fertor, fertōte, fertur, ferunt, feruntō, feruntor, feruntur, lāta, lātā, lātāe, lātārum, lātās, lāte, lātī, lātīs, lātō, lātōrum, lātōs, lātū, lātūm, lātūrūs, lātūs, tuleram, tulerāmus, tulerant, tulerās, tulerat, tulerātīs, tulēre, tulerim, tulerimus, tulerītis, tulerint, tuleris, tulerīs, tulerit, tuleritis, tulerītis, tulerō, tulērunt, tulī, tulimus, tulisse, tulissem, tulissēm, tulissēt, tulissētis, tulistī, tulistis, tulit},*

{*mālam, mālēbam, mālēbāmus, mālēbāt, mālēbātis, mālēmus, mālēs, mālet, mālētis, mālim, mālīmus, mālīnt, mālīs, mālit, mālītis, mālle, māllēm, māllēmus, māllēnt, māllēs, māllēt, māllētis, mālō, mālueram, māluerāmus, māluerant, māluerās, māluerat, māluerātis, māluēre, māluerīm, māluerimus, māluerīmus, māluerint, mālueris, māluerīs, māluerit, mālueritis, māluerītis, māluerō, māluērunt, māluī, māluīmus, māluisse, māluissēm, māluissēmus, māluissēnt, māluissēs, māluissēt, māluissētis, māluistī, māluistis, māluit, mālumus, mālunt, māvīs, māvult, māvultis},*

{morere, morerēminī, morerēmur, morerentur, morerer, morerēre, morerēris, morerētur, moreris, morī, moriāminī, moriāmur, moriantur, moriar, moriāre, moriāris, moriātur, moriēbāminī, moriēbāmur, moriēbantur, moriēbar, moriēbāre, moriēbāris, moriēbātur, moriēminī, moriēmur, moriendī, moriendō, moriendum, moriendus, moriēns, morientur, moriēre, moriēris, moriētur, morimīnī, morimur, morior, moritor, moritur, moritūrus, moriuntor, moriuntur, mortua, mortuā, mortuae, mortuam, mortuārum, mortuās, mortue, mortuī, mortuīs, mortuō, mortuōrum, mortuōs, mortuū, mortuum, mortuus},

{nōlam, nōlēbam, nōlēbāmus, nōlēbānt, nōlēbās, nōlēbat, nōlēbātis, nōlēmus, nōlēns, nōlēnt, nōlēs, nōlēt, nōlētis, nōlī, nōlim, nōlīmus, nōlīnt, nōlīs, nōlīt, nōlīte, nōlītō, nōlītōte, nōlle, nōllem, nōllēmus, nōllent, nōlēs, nōllet, nōlētis, nōlō, nōlueram, nōluerāmus, nōluerant, nōluerās, nōluerat, nōluerātis, nōluēre, nōluerim, nōluerimus, nōluerīmus, nōluerint, nōlueris, nōluerīs, nōluerit, nōlueritis, nōluerītis, nōluerō, nōluerunt, nōluī, nōluīmus, nōluīsse, nōluīssem, nōluīssēmus, nōluīscent, nōluīssēs, nōluīsset, nōluīssētis, nōluīstī, nōluīstis, nōluīt, nōlūmus, nōlūnt, nōlūntō},

{passa, passā, passae, passam, passārum, passās, passe, passī, passīs, passō, passōrum, passōs, passū, passum, passūrus, passus, patere, paterēminī, paterēmur, paterentur, paterēre, paterēris, paterētur, pateris, patī, patiāminī, patiāmur, patientur, patiar, patiāre, patiāris, patiātur, patiēbāminī, patiēbāmur, patiēbantur, patiēbar, patiēbāre, patiēbāris, patiēbātur, patiēminī, patiēmur, patientī, patiēdō, patiendum, patientus, patientur, patiēre, patiēris, patiētur, patimīnī, patimur, patior, patitor, patitūr, patiuntor, patiuntur},

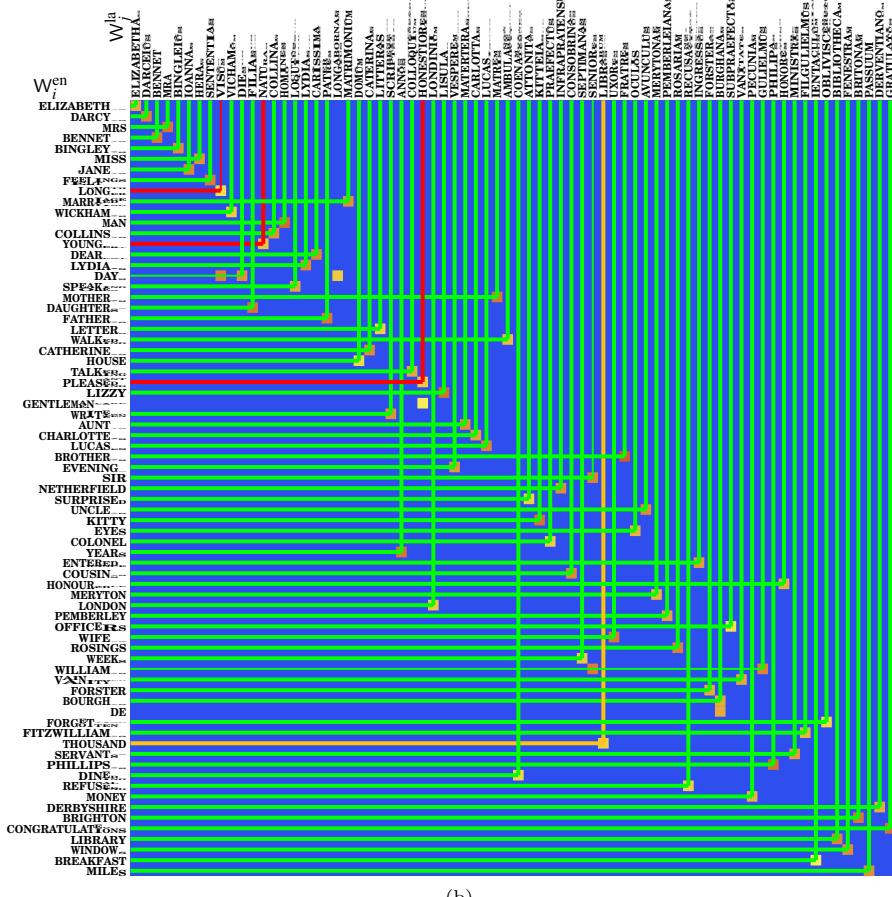
{rape, rapere, raperem, raperēminī, raperēmur, raperēmus, raperent, raperer, raperēre, raperēris, raperēs, raperet, raperētis, raperētur, raperis, rapī, rapiām, rapiāminī, rapiāmur, rapiāmus, rapiant, rapiantur, rapiar, rapiāre, rapiāris, rapiās, rapiat, rapiātis, rapiātur, rapiēbam, rapiēbāminī, rapiēbāmur, rapiēbāmus, rapiēbānt, rapiēbāntur, rapiēbar, rapiēbāre, rapiēbāris, rapiēbās, rapiēbat, rapiēbātis, rapiēbātur, rapiēminī, rapiēmur, rapiēmus, rapiendī, rapiendō, rapiendum, rapiendus, rapiēns, rapiēnt, rapiēntur, rapiēre, rapiēris, rapiēs, rapiēt, rapiētis, rapiētur, rapimīnī, rapimur, rapiō, rapiōr, rapiōs, rapiōt, rapiēt, rapiētis, rapiōtō, rapitor, rapitōte, rapitūr, rapiunt, rapiuntō, rapiuntor, rapiuntur, raptā, raptās, raptāe, raptām, raptārūm, raptās, rapte, raptī, raptīs, raptō, raptōrūm, raptōs, raptū, raptūm, raptūrūs, raptūs, rapueram, rapuerāmus, rapuerant, rapuerās, rapuerat, rapuerātis, rapuēre, rapuerim, rapuerimus, rapuerīmus, rapuerint, rapueris, rapuerīs, rapuerit, rapueritis, rapuerītis, rapuerō, rapuerunt, rapuī, rapuīmus, rapuisse, rapuissem, rapuissēmus, rapuissent, rapuissēs, rapuisset, rapuissētis, rapuistī, rapuistis, rapuit},

{velim, velīmus, velint, velīs, velit, velītis, velle, vellem, vellēmus, vellent, vellēs, vellet, vellētis, vīs, volam, volēbam, volēbāmus, volēbānt, volēbās, volēbat, volēbātis, volēmus, volēns, volent, volēs, volet, volētis, volō, volt, voltis, volueram, voluerāmus, voluerant, voluerās, voluerat, voluerātis, voluēre, voluerim, voluerimus, voluerīmus, voluerint, volueris, voluerīs, voluerit, volueritis, voluerītis, voluerō, voluērunt, voluī, voluīmus, voluisse, voluissem, voluissēmus, voluissent, voluissēs, voluisset, voluissētis, voluistī, voluistis, voluit, volumus, volunt, vult, vultis}.

*Example 6.29.4.* In Fig. S9, we further apply the aforementioned word clustering algorithm to topic extraction and machine translation (see Table S1 for text source). Since vowel lengths are usually not marked in written Latin, conflations of certain concepts are unavoidable in our algorithm. For example, we have *mēnsis* “month” (nominative/genitive/vocative singular) vs. *mēnsīs* “tables” (dative/ablative plural), *miserīs* “poor” (masculine/feminine/neuter ablative plural) vs. *mīseris* “you will have sent” (second person singular, future perfect).

**ELIZABETHA.** HR REAT CERTIOREM VEL...  
**SORORP.** BENNET MRA BINGLEI<sup>2</sup> NOVE  
**RESPONSUS ANIMO** RESPONDIT COMITATU<sup>3</sup> AMORE<sup>4</sup> VISSU<sup>5</sup> VICHAM<sup>6</sup> DIE<sup>7</sup> FILIA<sup>8</sup> NATURA<sup>9</sup> DOMINA<sup>10</sup> AMIC<sup>11</sup>  
**VENI<sup>12</sup>** COLLINA<sup>13</sup> TEMPORA<sup>14</sup> HOMINA<sup>15</sup> MENTI<sup>16</sup> LOQU<sup>17</sup> VILLAM FAMILIA<sup>18</sup> CARISSIMA SCIRE<sup>19</sup> SATIS<sup>20</sup>  
**PATRE<sup>21</sup>** LONGA BORNARI VITAM<sup>22</sup> EXSPECTA<sup>23</sup> AUDI<sup>24</sup> DUBIT<sup>25</sup> GESE<sup>26</sup> ANIM<sup>27</sup> OP<sup>28</sup> MATRIMONIUM MISER<sup>29</sup> FACILIT<sup>30</sup> DOM<sup>31</sup>  
**CATERINA<sup>32</sup>** LITTERAS<sup>33</sup> ACCEPT<sup>34</sup> PRIM<sup>35</sup> PUELLA<sup>36</sup> PET<sup>37</sup> GRATIAS<sup>38</sup> SPERO<sup>39</sup> DUCERE<sup>40</sup> BENE<sup>41</sup> FEMIN<sup>42</sup> VERB<sup>43</sup> IUDIC<sup>44</sup>  
**COETUS** CONSUTUDINE<sup>45</sup> LIBER<sup>46</sup> CANT<sup>47</sup> CAMERA<sup>48</sup> SUABE<sup>49</sup> SEDE<sup>50</sup> LAR<sup>51</sup> VADE<sup>52</sup> COEN<sup>53</sup> SERMONEA<sup>54</sup> APART<sup>55</sup> MAGN<sup>56</sup> IUVEN<sup>57</sup> DISC<sup>58</sup> VIT<sup>59</sup>  
**INTRAPRATENS** CONSUETUDINE<sup>60</sup> STAV<sup>61</sup> MATER LAUD<sup>62</sup> MALE<sup>63</sup> MAXIMA<sup>64</sup> ACCIPERE<sup>65</sup> ROGAR<sup>66</sup> CURA<sup>67</sup> LOC<sup>68</sup> MANE<sup>69</sup> RATIONE<sup>70</sup> BREV<sup>71</sup> PARTEM<sup>72</sup> PLAC<sup>73</sup>  
**SIMILITUDIN** CANT<sup>74</sup> LIBRA<sup>75</sup> GARDNER<sup>76</sup> UXOR<sup>77</sup> OCULAS FRATR<sup>78</sup> CONV<sup>79</sup> HONESTIO<sup>80</sup> CAP<sup>81</sup> FELICITAT<sup>82</sup>  
**IGNOR<sup>83</sup>** BEAT<sup>84</sup> MIN<sup>85</sup> NUP<sup>86</sup> S<sup>87</sup> OPIN<sup>88</sup> HORAS<sup>89</sup> MUND<sup>90</sup> TIM<sup>91</sup> MOLES<sup>92</sup> DULCISSIMA<sup>93</sup> FESTUS<sup>94</sup> COEM<sup>95</sup> FESTUS<sup>96</sup> AD<sup>97</sup> FIN<sup>98</sup> TAC<sup>99</sup> ASCE<sup>100</sup>  
**MUTAT<sup>101</sup>** MON<sup>102</sup> PRO<sup>103</sup> MARIT<sup>104</sup> CINT<sup>105</sup> FORT<sup>106</sup> MIR<sup>107</sup> PEMBERLEIAN<sup>108</sup> EID<sup>109</sup> HOSPIT<sup>110</sup> RHAEDA<sup>111</sup> OFFIC<sup>112</sup> CITA<sup>113</sup> LEV<sup>114</sup> EXPLIC<sup>115</sup> SOLICIT<sup>116</sup> COGN<sup>117</sup> REVEN<sup>118</sup>  
**PERTURB<sup>119</sup>** EPISTOL<sup>120</sup> VOLUNT<sup>121</sup> ELEGANT<sup>122</sup> MUL<sup>123</sup> ROSARIAM COHORTE<sup>124</sup> JUCUND<sup>125</sup> RECUSA<sup>126</sup> SUMMA<sup>127</sup> OSTEND<sup>128</sup> LAP<sup>129</sup> AMPL<sup>130</sup> SUB<sup>131</sup> COND<sup>132</sup> MEMORIA<sup>133</sup> TRANQUILL<sup>134</sup>  
**AV<sup>135</sup>** FORTUNA<sup>136</sup> FASTID<sup>137</sup> AESTIMA<sup>138</sup> CUPID<sup>139</sup> MOD<sup>140</sup> FILI<sup>141</sup> AUDI<sup>142</sup> VIAM<sup>143</sup> GARDINER<sup>144</sup> INDIC<sup>145</sup> SUBRIDES<sup>146</sup> SUSPIC<sup>147</sup> COGIT<sup>148</sup> BENIGN<sup>149</sup> DILIGENT<sup>150</sup> TEN<sup>151</sup> MAIOR<sup>152</sup> NOM<sup>153</sup>  
**CASE<sup>154</sup>** ATTIN<sup>155</sup> FRATER<sup>156</sup> CONSPIC<sup>157</sup> CO<sup>158</sup> COGNOSC<sup>159</sup> VESTIB<sup>160</sup> OFTEN<sup>161</sup> INVIT<sup>162</sup> FEST<sup>163</sup> VERT<sup>164</sup> QUESA<sup>165</sup> OBSERV<sup>166</sup> AVUNCUL<sup>167</sup> LEG<sup>168</sup> AVUST<sup>169</sup>  
**CELLA OH VAN<sup>170</sup>** CONTING<sup>171</sup> SUPERBIA<sup>172</sup> NESC<sup>173</sup> DIMID<sup>174</sup> CRAST<sup>175</sup> FAVER<sup>176</sup> CONV<sup>177</sup> IUST<sup>178</sup> VERT<sup>179</sup> OBSERV<sup>180</sup> AVUNCUL<sup>181</sup> LEG<sup>182</sup> AVUST<sup>183</sup>  
**PRUDENT<sup>184</sup>** CONST<sup>185</sup> ASPIRC<sup>186</sup> APPAR<sup>187</sup> DOM<sup>188</sup> PECUNIA<sup>189</sup> GUILIELM<sup>190</sup> DIFFIC<sup>191</sup> PRO<sup>192</sup> SIGN<sup>193</sup> CONFIRM<sup>194</sup> COMFOR<sup>195</sup> PHILIP<sup>196</sup> JUNG<sup>197</sup> CAN<sup>198</sup>  
**INTELLIG<sup>199</sup>** EXCV<sup>200</sup> HONOR<sup>201</sup> COMMEND<sup>202</sup> LINGU<sup>203</sup> COLLOC<sup>204</sup> CONNUB<sup>205</sup> REFR<sup>206</sup> DU<sup>207</sup> SING<sup>208</sup> PRAT<sup>209</sup> REPET<sup>210</sup> CONV<sup>211</sup> PHILIP<sup>212</sup> JUNG<sup>213</sup> CAN<sup>214</sup>  
**ADMISSION<sup>215</sup>** SPATI<sup>216</sup> FRATER<sup>217</sup> CONSPIC<sup>218</sup> COGNOSC<sup>219</sup> VESTIB<sup>220</sup> OFTEN<sup>221</sup> INVIT<sup>222</sup> FEST<sup>223</sup> VERT<sup>224</sup> QUESA<sup>225</sup> OBSERV<sup>226</sup> AVUNCUL<sup>227</sup> LEG<sup>228</sup> AVUST<sup>229</sup>  
**FRIEND<sup>230</sup>** TRO<sup>231</sup> POD<sup>232</sup> STOL<sup>233</sup> CON<sup>234</sup> IN<sup>235</sup> CON<sup>236</sup> CON<sup>237</sup> RIDE<sup>238</sup> PER<sup>239</sup> OCCUP<sup>240</sup> PRAS<sup>241</sup> ADM<sup>242</sup> AFFECT<sup>243</sup> PERT<sup>244</sup> ABEG<sup>245</sup> CLARE<sup>246</sup> LUD<sup>247</sup> ABEG<sup>248</sup>  
**INTUIT<sup>249</sup>** HER<sup>250</sup> BEI CON<sup>251</sup> PRIV<sup>252</sup> OBLIV<sup>253</sup> ER<sup>254</sup> AEGROY<sup>255</sup> OL<sup>256</sup> L<sup>257</sup> INSOL<sup>258</sup> D<sup>259</sup> PRIV<sup>260</sup> NAR<sup>261</sup> DEPINGER<sup>262</sup> CON<sup>263</sup> PANT<sup>264</sup> VIGINTI<sup>265</sup> PRAS<sup>266</sup> NEG<sup>267</sup> LENT<sup>268</sup> APPROP<sup>269</sup>  
**ART<sup>270</sup>** INTEN<sup>271</sup> HUNSFORDIAN<sup>272</sup> BELLS<sup>273</sup> VIBR<sup>274</sup> QUOD<sup>275</sup> BIBLIOTHECA<sup>276</sup> DIVIT<sup>277</sup> CAUS<sup>278</sup> ARGU<sup>279</sup> PROF<sup>280</sup> AEST<sup>281</sup> REFL<sup>282</sup> L<sup>283</sup> GAV<sup>284</sup> FEN<sup>285</sup> BRITON<sup>286</sup> RITE<sup>287</sup> FACIL<sup>288</sup> SALVER<sup>289</sup> QUATTUOR<sup>290</sup> INFELIC<sup>291</sup>  
**GEORGIA<sup>292</sup>** AGNO<sup>293</sup> SERVAT<sup>294</sup> MUS<sup>295</sup> ULIV<sup>296</sup> LABOR<sup>297</sup> GRATUL<sup>298</sup> EXPON<sup>299</sup> TENT<sup>300</sup> EXP<sup>301</sup> REQU<sup>302</sup> PRAE<sup>303</sup> SEPT<sup>304</sup> PERN<sup>305</sup> DERVENTIANO<sup>306</sup> AFFL<sup>307</sup> ACRI<sup>308</sup> AFFR<sup>309</sup> HOD<sup>310</sup> POSTUL<sup>311</sup> ISS<sup>312</sup> CONNECT<sup>313</sup> HIS<sup>314</sup>

(a)



(b)

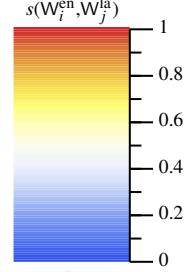


Fig. S9. Text mining in Latin. (a) Statistically identified topics ( $n_{ii} \geq 20$ ) in a Latin version of *Pride and Prejudice*, with the same color encoding scheme as Fig. S3. (b) Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{la}})$  between selected topics in English and Latin versions of *Pride and Prejudice*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) ranking of the similarity score for the optimal node. Green (resp. red) cross-hair indicates a correct (resp. incorrect) match. Amber cross-hair marks a link between distinct concepts that share the same hypernyms (such as proper names, kinship terms, or emotions).

## 7 Approximate word clustering in selected Slavic languages

In this section, we present the approximate clustering algorithms for two representative Slavic languages: Polish and Russian.

Despite differences in the writing systems (Polish uses the Latin script while Russian employs the Cyrillic alphabet), both these languages share some morphological features in common:

- Polish nouns decline in seven cases, Russian six (in singular and plural forms).
- Adjectives are declined in three genders and in all the cases available to nouns.
- The same verb in imperfective and perfective aspects often differ by a prefix or an alternation of consonant/vowel. A limited amount of verbs have etymologically unrelated stems in their imperfective and perfective aspects.
- Both Polish and Russian have a handful of “verbs of motion” that contrast multidirectional and unidirectional forms.

These characterisations also apply to other Slavic languages (Croatian, Czech, Slovak and Slovene written in the Latin script; Belarusian, Bulgarian, Macedonian, Serbian and Ukrainian written in the Cyrillic script), more or less. In addition, the two modern Baltic languages, Latvian and Lithuanian, are similar to Slavic languages at the morphological level.

Some common words in Polish and Russian have very short (even vowelless) stems, and the short stems themselves may allow vowel and/or consonant alternations. At times, words with very different meanings have nearly identical stems. Thus it is highly challenging to design approximate clustering algorithms for these languages, while avoiding confusions of unrelated words. In the methods developed in this section, we make a compromise between a workable algorithm and error-free clustering.

### 7.1 Modified Porter stemming algorithm for Polish

**Definition 7.1** (Polish stop words). If a word belongs to the following list<sup>87</sup>:

*a, aby, ach,acz,aczkolwiek, aj, albo, albowiem, ale, ależ, ani, aż, bardziej, bardzo, bądź, bądźcie, bądźmy, bez, będą, będąc, będąca, będące, będący, będąd, będązie, będąziemy, będąziesz, bo, bowiem, by, bycie, być, byli, byliby, bylibyście, bylibyśmy, byliście, byliśmy, był, byla, byłaby, byłabym, byłabyś, byłam, byłaś, byłby, byłbym, byłbyś, byłem, byłeś, było, byłoby, były, byłyby, byłybyście, byłybyśmy, byłyście, byłyśmy, bym, bynajmniej, byś, byście, byśmy, bywszy, cali, cała, cały, chociaż, choćby, chyba, ci, ciebie, cię, co, cokolwiek, coraz, coś, cóż, czasami, czasem, czego, czegokolwiek, czegoś, czemu, czemukolwiek, czemuś, często, czy, czyjś, czyli, czym, czymkolwiek, czymś, czyż, czyżby, czyżbym, czyżbyś, czyżbyście, czyżbyśmy, daleko, dla, dlaczego, dlatego, do, dobrze, dokąd, dokądś, doń, dopóki, doprawdy, dosyé, dość, dużo, dwa, dwaj, dwie, dwoje, dzięki, dzisiaj, dziś, gdy, gdyby, gdybym, gdybyś, gdybyście, gdybyśmy, gdyż, gdzie, gdziekolwiek, gdzieś, go, i, ich, ile, ileś, iloma, ilu, im, inna, inną, inne, innego, innej, innemu, inni, inny, innych, innym, innymi, iż, ja, jacy, jacyś, jak, jaka, jakaś, jaką, jakąś, jakby, jaki, jakich, jakichś, jakie, jakiego, jakiegoś, jakiej, jakiejś, jakiemu, jakiemuś, jakieś, jakim, jakimi, jakimiś, jakimś, jakiś, jakiz, jakkolwiek, jako, jakoś, ja, je, jeden, jedna, jednak, jednakże, jedną, jedne, jednego, jednej, jednemu, jedni, jedno, jednych, jednym, jednymi, jedyna, jedną, jedyną, jedynego, jedyną, jedynemu, jedyny, jedynie, jedynych, jedynym, jedynymi, jego, jej, jemu, jest, jestem, jesteś, jesteście, jesteśmy, jeszcze, jeśli, jeżeli, już, każda, każdą, każdej, każdejgo, każdej, każdemu, każdy, każdych, każdym, każdymi, kiedy, kiedyś, kilka, kim, kimkolwiek, kimś, kogo, kogokolwiek, kogoś, koło, komu, komukolwiek, komuś, kto, ktokolwiek, ktoś, która, którą, które, którego, której, któremu, który, których, którym, którymi, któryś, którzy, ku, lat, lecz, lub, ma, macie, mają, mając, mająca, mające, mający, mało, mam, mamy, masz, mą, me, meg, mej, memu, mi, miał, miałI, miała, miałaby, miałabym, miałabyś, miałam, miałas, miałby, miałbym, miałbyś, miałem, miałes, miało, miałoby, miały, miałyby, miałybyście, miałybyśmy, miałyście, miałyśmy, miano, mieć, miej, miejsce, miejmy, mieli, mieliby, mielibyście, mielibyśmy, mieliście, mieliśmy, mię, między, mimo, mną, mnie, mogą, mogąc, mogąca, mogące, mogący, mogę, mogli, mogliby, moglibyście, moglibyśmy, mogliście, mogliśmy, mogła, mogłaby, moglabym, moglabyś, mogłam, mogłaś, mogłem, mogłeś, mogło, mogłoby, mogłobym, mogłobyś, mogłom, mogłoś, mogły, mogłyby, mogłybyście, mogłybyśmy, mogłyście, mogłyśmy, moi, moich, moim, moimi, moja, moją, moje, mojego, mojej, mojemu, może, możecie, możemy, możesz, możliwe, można, móc, mógl, móglby, móglbym, móglbyś, mój, mu, musi, musiął, musiala, musialaby, musialabym, musialabyś, musialam, musialaś, musialby, musialbym, musialbys, musialeł, musialo, musialoby, musialobym, musialobyś, musialom, musialoś, musialy, musialyby, musialybyście, musialybyśmy, musialyście, musialyśmy, musicie, musicieć, musiel, musieliby, musielibyście, musielibyśmy, musielicie, musieliszy, musimy, musisz, musząc, musząca, muszące, muszący, muszę, my, mych, mym, mymi, na, nad, najbardziej, nam, nami, nań, naokoło, naprzeciw, nas,*

<sup>87</sup>Our list of Polish stop words is based on <https://pl.wikipedia.org/wiki/Wikipedia:Stopwords>, with extensive additions to roughly match their counterparts in English. In particular, we have included all the inflected forms of *co* “what”, *każdy* “every”, *kto* “who”, *nikt* “nobody”, *nic* “nothing”, *wszyscy*, *wszystek*, *wszystkie*, *wszystko* “all (in various genders)” and all the personal pronouns.

*nasi, nasz, nasza, naszą, nasze, naszego, naszej, naszemu, naszych, naszymi, natomiast, natychmiast, nawet, nią, nic, nich, niczego, niczemu, niczym, nie, niech, niechaj, nieco, niego, niej, niektóre, niektórych, niektórym, niektórymi, niektórzy, niemu, nigdy, nikim, nikogo, nikomu, nikt, nim, nimi, niż, no, o, oba, obaj, obie, obiema, oboje, obojga, obojgiem, obojgu, obok, oboma, obu, oby, obym, obyś, obyscie, obyśmy, od, okolo, on, ona, one, oni, ono, oprócz, oraz, oto, owa, ową, owe, owego, owej, owemu, owi, owo, owszem, owych, owym, owymi, ów, pan, pana, pani, po, pod, podczas, pomiędzy, pomimo, ponad, ponieważ, poprzez, potem, powinien, powinienniem, powinieneś, powinna, powinnam, powinnaś, powinni, powinniście, powinniśmy, powinno, powinny, powinnyście, powinnyśmy, poza, później, prawie, przecież, przeciw, przeciwko, przed, przede, przedtem, przez, przy, razem, roku, również, sam, sama, samą, same, samego, samej, samemu, sami, samo, samych, samym, samymi, są, siebie, się, skąd, skądś, skoro, sobą, sobie, spomiędzy, sponad, sposób, spośród, spoza, swa, swą, swe, swego,owej, swemu, swoi, swoich, swoim, swoimi, swoja, swoją, swoje, swojego, swojej, swojemu, swój, swych, swym, swymi, ta, tacy, tak, taka, taką, taki, takich, takie, takiego, takiej, takiemu, takim, takimi, także, tam, tamci, tamta, tamtą, tamte, tamtego, tamtej, tamtemu, tamten, tamto, tamtych, tamtym, tamtymi, tą, te, tego, tej, temu, ten, teraz, też, tę, to, tobą, tobie, toteż, trzeba, tu, tutaj, twa, twą, twe, twego, twej, twemu, twoi, twoich, twoim, twoimi, twoja, twoje, twojego, twojej, twojemu, twój, twych, twym, twymi, ty, tych, tyle, tylko, tyloma, tylu, tym, tymi, tyś, u, w, wam, wami, was, wasi, wasz, waszą, wasze, waszego, waszej, waszemu, waszych, waszym, waszymi, wbrew, we, według, wiele, wieloma, wielu, więc, więcej, wkrótce, właśnie, wobec, wszak, wszakże, wszyscy, wszystek, wszystka, wszystkich, wszystkie, wszystkiego, wszystkiej, wszystkiemu, wszystkim, wszystkimi, wszystko, wśród, wtedy, wy, z, za, zapewne, zaraz, zaś, zawsze, zbyt, ze, zł, znowu, znów, został, zza, ż, żaden, żadna, żadne, żadnych, że, żeby,*

then we consider it a Polish stop word. All the Polish stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list.  $\square$

In addition, we define string patterns **eatPolish**, **givePolish**, **MissPolish**, **MrsPolish** through inflected forms of four Polish words:

*jada, jadać, jadaj, jadają, jadając, jadającą, jadające, jadający, jadali, jadał, jadała, jadałe, jadało, jadały, jadana, jadane, jadani, jadanie, jadano, jadany, jadasz, jadł, jadła, jadle, jadło, jadły, je, jedli, jedz, jedzą, jedząc, jedząca, jedzące, jedzący, jedzeni, jedzenie, jedzona, jedzone, jedzono, jedzony, jesz, jeść — “eat”;*

*da, dacie, dać, dadzą, daj, dają, dając, dającą, dające, dający, dajcie, daje, dajecie, dajemy, dajesz, daję, dajmy, dali, daliby, dalibyście, dalibyśmy, daliście, daliśmy, dal, dała, dalaby, dalabym, dalabyś, dalaś, dalby, dalbym, dalbyś, dałem, dałeś, dalo, daloby, daly, dalyby, dalybyście, dalybysmy, dalyście, dalyśmy, dam, damy, dana, dane, dani, danie, dano, dany, dasz, dawać, dawaj, dawajcie, dawajmy, dawali, dawaliby, dawalibyście, dawalibyśmy, dawaliście, dawaliśmy, dawał, dawała, dawałaby, dawałabym, dawałabyś, dawałam, dawałaś, dawały, dawałbym, dawałyś, dawałem, dawałeś, dawało, dawało1, dawałoby, dawały, dawałyby, dawałybyście, dawałybysmy, dawałyście, dawałyśmy, dawana, dawane, dawani, dawanie, dawano, dawany, dawszy — “give”;*

*panna, pannach, pannami, panną, pannę, pannie, panno, pannom, panny — “Miss”;*

*paniach, paniami, panią, panie, paniom, pań — “Mrs”.*

In order to stay consistent with their counterparts in English, we are not going to treat these string patterns as stop words in Polish. Instead, we will exploit them in clustering of Polish content words.

### 7.1.1 Effective spelling and essential root

It is assumed that all Polish words are converted to lowercase<sup>88</sup> before going through any of the procedures below.

**Definition 7.2** (Polish Vowels and Verb Prefixes). Hereafter in §7.1, the symbol **V** stands for any member from the list of Polish vowels  $\{a, q, e, \dot{e}, i, o, \acute{o}, u, y\}$ . In line with the multiplicity notations introduced in Definition 3.3, the symbol **V<sub>m</sub>** stands for a text string formed by consecutive appearance of one or more (not necessarily identical or distinct) members from the set of Polish vowels.

Dual to the notations above, the symbol **C** stands for any character that does not belong to the list  $\{a, q, e, \dot{e}, i, o, \acute{o}, u, y\}$ , and **C<sub>m0</sub>** stands for a text string formed by consecutive appearance of zero or more (not necessarily identical or distinct) characters that do not belong to the same list.

The symbol **P** = ( $\emptyset | do | na | nad | o | ob | obe | po | pod | pode | prze | przed | przy | roz | s | sj | u | w | we | wy | z | za | ze$ ) represents the possible prefix of a Polish verb.  $\square$

<sup>88</sup>As of v11.0, *Mathematica* does not produce conversions  $A \rightarrow q$ ,  $\acute{E} \rightarrow \dot{e}$ ,  $\acute{N} \rightarrow \acute{n}$ ,  $\acute{S} \rightarrow \acute{s}$ ,  $\acute{Z} \rightarrow \acute{z}$  upon invoking the *ToLowerCase* operation. Therefore, case conversions for these special Polish letters need to be hand-coded if one implements our algorithm in *Mathematica*.

**Algorithm 7.3** (Polish effective spelling). For a Polish word  $\hat{\sigma}$ , its effective spelling  $\text{EffSpell}(\hat{\sigma})$  is constructed in sequential steps:

(1) Replace

<b>Prze(cz k)</b>	<b>(<math>\emptyset co</math>)dzienn~</b>	<b>(<math>\emptyset naj</math>)czarn~</b>	<b>(<math>\emptyset naj</math>)koniecz~</b>
<b>Přeček</b>	<b>dni</b>	<b>βłak</b>	<b>ñσecz</b>
$(\emptyset naj)wyżs(\emptyset z)~$ ηwykη	$(\emptyset naj)zim~$ ωiñt	$(\emptyset o)żeni~$ małżeństwo	$(\emptyset o po)(ślub żen)i~$ małżeństwi
$(\emptyset wy)marz~$ δριμ	$(mężat zamężn żonat)~$ małżeństw	$(ojcz tat)ul(\emptyset e)k~$ ojc	$(s ś)piesz$ χup
bogat~ ρbkit	bram~ βrau	dlu(g ż) λoy	dziewię~ 9e
kieli(ch sz)~ καλχ	ksi(qżecz łedz eg) książ	marc~ μ3ap	martw μɔrt
panowie panach	plan πλaν	pocz πoξ	poje~ je
rozpacz~ δεσπρ	rób~ zrob	skr~ σkp	sp(o ó)C~ sCp
świat~ ωświłδ	tajemnicz~ μyσticz	trz(y e(ch j m ma))~ 3e	stan(o ó)w κοñtw
um( $\emptyset a ie$ )r( $\emptyset z$ )X~ μɔrt	upły~ ωply	v 5	wiejsk~ ρiρλ
wynal(a e)(z ż)~ βynaido	wysoc(k k)~ ηwykη	wzmian~ umeñtn	x 10
znal(a e)(z ż)~ znajdo	zys(k zcz)~ γař	zaw( $\emptyset a ie$ )r( $\emptyset z$ )X~ iñκλ	stos stoσ
<b>MrsPolish</b>		<b>(<math>\emptyset z</math>)eatPolish(<math>\emptyset by</math>)(<math>\emptyset ś</math>)(<math>\emptyset cie m my</math>)</b>	<b>(iśc szli)</b>
μipσ		eat	idz
<u>czu(j l t t ~ć)</u> φιλυć	<u>de</u> δee	<u>l</u> 50	<u>lady</u> damskich
<u>rzec</u> rzekl	<u>ta(cie tach)</u> ojcu	<u>tat(V m w)m</u> ojcu	<u>~Cže</u> C

(2) Replace

<b>Pmiesza~</b>	<b>Ptrz</b>	<b>(do na o od po roz wy za)(łoz łóż)~</b>
<b>mixa</b>	<b>Pτρ</b>	<b>klad</b>
X $\epsilon$ ((poj przyj)(q e m))~ σX	brać~ wziąć	chud(l l n)~ woτ
licz( $\emptyset b$ )~ ñiuma	list~ λiστ	łaṁx~ μlaμx
opin~ ωpin	os(o ó)b~ oσoβ	pan(ow uj)~ κotρλ
przednio(c t)X~ θiγ	przekon~ βliφ	qm(a o ó)w δeñη
uszy $X^{\epsilon}$ (c ć j k l ł t ~ć) szyX	will~ βiλλ	rze(cz k ł k)~ rzekl
przyjaci(e o ó)(l l)( $\emptyset c k$ )X πprzyjak	zach~ Ξ	zna~ ζna
	zobacz~ widzi	życz~ ωiσt
		(iście iść szli szlobym szłom) szedlem
		~ $X^{\epsilon}(k l)ach$ ∅
		~kich Xa
		~oce ki
		oc

(3) Do dom~ → δoμ, pie(c|cz|k) → βak, **givePolish**( $\emptyset|by$ )( $\emptyset|ś$ )( $\emptyset|cie|m|my$ ) → dar; ~kacie → ka.

(4) Replace

$(\mathbf{P} \emptyset naj)l\sim$	$(\emptyset naj)gors(i z)\sim$	$(\emptyset naj)gorzej\sim$	$(\emptyset naj)lepiej\sim$	$(\emptyset naj)leps(i z)\sim$
$\lambda$	$\acute{z}l$	$\acute{z}le$	$\delta obrze$	$\delta obr$
$(\emptyset naj)mniejs(i z)\sim$	$(\emptyset naj)więks(i z)\sim$	$(\emptyset za)ta(nieć ńc)(\emptyset z)\sim$	$(mów powie)\sim$	$(my umy)\sim$
$mał$	$duż$	$tańc$	$πowie$	$uμy$
$(wiad wiedz)\sim$	$X_j^{\epsilon}(s z)t$	$bra(c č t)\sim$	$chc\sim$	$cioc\sim$
$wie$	$X_jl$	$βbra$	$choc$	$tsioc$
$ksi(q ę)dz\sim$	$mam\sim$	$mat(c ek k)\sim$	$mył\sim$	$naj(\hat{x})_m X_1^{\epsilon}(ej sz)X$
$księż$	$μam$	$μam$	$vy$	$(\hat{x})_m X_1$
$orzel(\emptyset ek k)\sim$	$pi(q ę ęc ęc)\sim$	$pies\sim$	$piln\sim$	$tydzień\sim$
$orł$	$5a$	$psa$	$sn$	$tygodni$
$(ojciec ojcze)$	$kwiecień$	$oj$	$pić$	$tatu(s ~s)$
$ojcach$	$kwietni$	$yoy$	$pici$	$uch((\emptyset V)(\emptyset ch m mi))$
				$uszach$
				$wielce$
				$duż$

(5) Do  $ch \rightarrow s\chi$ ,  $cz \rightarrow č$ ,  $l \rightarrow l$ ,  $ó \rightarrow o$ ,  $sz \rightarrow š$ ,  $\sim det \rightarrow \emptyset$ ,  $\sim i \rightarrow j$ .

(6) Do  $wi(ad|edz) \rightarrow wie$ ,  $Ci \rightarrow Cji$ ,  $X^{\epsilon}(\mathbf{P}mow) \rightarrow X_{maw}$ , and call the result  $\hat{\sigma}'$ .

(7) Break down  $\hat{\sigma}' = \hat{\sigma}_1 \hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}'^{[\min\{\lambda(\hat{\sigma}'), \ell(\hat{\sigma}')\}]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where  $\lambda(\hat{\sigma}') - 1$  equals either the last position occupied by the string pattern  $(\mathbf{P}|bez(\emptyset|e)|njie|pa|st(\emptyset|r))\hat{x}\sim$  in the non-void  $\hat{\sigma}'$ , or  $-1$  if  $\hat{\sigma}' = \emptyset$ .

(8) On  $\hat{\sigma}_2$ , perform the following substitutions in a sequel:

(8.1) Replace

$X^{\epsilon}(g k l)j$	$ck$	$dzk$	$ńc$	$ow(\emptyset č ji(č ec) nji(\emptyset c k) sk)$	$(a(c st t) ljiw (\emptyset aw nji)č(\emptyset yk ~ycy) njisk (a(l r)n(\emptyset ji)$
$X$	$c$	$d$	$ń$	$\emptyset$	
				$(ač(\emptyset ce k) a(\emptyset ńs)k ar(ce k z) by dl icie l i(st ści w) kacj s(k tw) y(ciel st ści w))$	
				$\emptyset$	
$\sim a(ček cy rek)$				$\sim X^{\epsilon}(nt st)(ce ek k)\hat{x}_{m_0}$	$\sim ij$
$\emptyset$				$X$	$i$
					$\sim mjiq$
					$mjenj$

(8.2) Do  $kV \rightarrow V$ ,  $\sim(ce|k) \rightarrow \emptyset$ ,  $\sim((a|i|y)s\chi|e|ego|em|emu|im|mj|ow(\emptyset|j)V_{m_0}|om|um|ym|Vč) \rightarrow \emptyset$ .

The result after these two steps of operations is called  $\hat{\sigma}'_2$ .

(9) Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ .

(10) Do  $č \rightarrow cj$ ,  $ń \rightarrow nj$ ,  $ś \rightarrow sj$ ,  $ź \rightarrow zj$ .

**Definition 7.4** (Polish protected range). Let  $\hat{\sigma}$  be a text string derived from a Polish word, its protected range  $\text{ProtRg}(\hat{\sigma}) = \min\{\lambda_1(\hat{\sigma}), \lambda_2(\hat{\sigma})\}$  is determined by two non-negative integers  $\lambda_1(\hat{\sigma})$  and  $\lambda_2(\hat{\sigma})$  specified through the following procedures:

- Look for the string pattern  $(na|u|we|wy)_{m_0} C_{m_0} V(j|s)_{m_0} \sim$  in the string  $\hat{\sigma}$ ;
- If the string pattern above is found, the last position occupied by such a string defines  $\lambda_1(\hat{\sigma})$ ; otherwise, set  $\lambda_1(\hat{\sigma}) = 0$ ;
- Look for the pattern  $(na|po|prze|przy|roz|s|w|wy|z)V_{m_0} CV_{m_0} C \sim$  in the string  $\hat{\sigma}$ ;
- If the pattern above is found, the last position occupied by such a pattern defines  $\lambda_2(\hat{\sigma})$ ; otherwise, set  $\lambda_2(\hat{\sigma}) = \ell(\hat{\sigma})$ .  $\square$

**Algorithm 7.5** (Polish essential root). Let  $\hat{\sigma}$  be the effective spelling of a Polish word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:

(1) Break down  $\hat{\sigma} = \hat{\sigma}_1 \hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}^{[\text{ProtRg}(\hat{\sigma})]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma})$  is equal to the protected range of  $\hat{\sigma}$ .

(2) On  $\hat{\sigma}_2$ , perform the following substitutions in a sequel:

(2.1) Do  $\sim V_m \rightarrow \emptyset$ .

(2.2)  $\text{Do } \sim \mathbf{V}_m(m|n(\emptyset|j)|w\check{s}) \rightarrow \mathbf{V}_m, \sim \mathbf{V}_{m_0}\check{s} \rightarrow \mathbf{V}_{m_0}, \sim(qc|ql|\check{ec}|l\check{et}) \rightarrow \emptyset.$

(2.3)  $\text{Do } \sim \mathbf{V}_m \rightarrow \emptyset.$

The result after these three steps of operations is called  $\hat{\sigma}'_2$ .

(3) *Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ .*

In Polish (as well as Russian and many other Slavic languages), the imperfect and perfect aspects of the same verb may differ by a prefix. To improve the performance of verb clustering, we need the following extensions to Definition 7.4 and Algorithm 7.5.

**Definition 7.6** (Polish verbal protected range). Let  $\hat{\sigma}$  be a text string derived from a Polish word, its verbal protected range  $\text{VbProtRg}(\hat{\sigma})$  is an integer determined as follows:

- Try to find the string pattern  $(do|na|o|odz|po|prze|we|wy)_{m_0} \mathbf{C}_{m_0} \mathbf{V}(j|l)_{m_0} \sim$  in the string  $\hat{\sigma}$ ;
- If the string pattern above is found, the last position occupied by such a string defines  $\text{VbProtRg}(\hat{\sigma})$ ; otherwise, set  $\text{VbProtRg}(\hat{\sigma}) = 0$ .  $\square$

**Algorithm 7.7** (Polish verbal essential root). *Let  $\hat{\sigma}$  be the token string associated with a Polish word, then its corresponding verbal essential root  $\text{VbEssRoot}(\hat{\sigma})$  is constructed in the following steps:*

- (1) *Break down  $\hat{\sigma} = \hat{\sigma}_1 \hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}^{[\text{VbProtRg}(\hat{\sigma})]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{VbProtRg}(\hat{\sigma})$  is equal to the verbal protected range of  $\hat{\sigma}$ .*
- (2) *Do  $\sim(\mathbf{V}_{m_0}(l)(\mathbf{V}_{m_0}(\emptyset|c|s|m)(\emptyset|j)(\emptyset|m)))_{m_0} \rightarrow \emptyset$  on  $\hat{\sigma}_2$ , and call the result  $\hat{\sigma}'_2$ .*
- (3) *Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ .*
- (4) *Do  $\sim \hat{\chi}^\epsilon(d|m|n|r|s|w|z)j \rightarrow \hat{\chi}, \sim c(\emptyset|j) \rightarrow t, \sim \check{z} \rightarrow z$ .*
- (5) *Replace*

$(bjierz bjior br bra bran wez wz wzji)$	$\sim dn$	$l(a al ej)$	$pji(\emptyset j t)$	$\sim \check{z}dz$
$bral$	$d$	$lan$	$pjl$	$zjdz$

(6) *Replace*

$(\emptyset do na o odz po roz wy za)kla(\emptyset d dd dz)$	$klad$		
$(\emptyset w)(i id idz s\chi odz \check{s}jla \check{s}jlo \check{s}jly \check{s}lji)$	$s\chi odz$	$(\emptyset w)(jad jedz jes\chi)$	$jes\chi$
$\underline{x_1^\epsilon(do ob odz po przy we wy)(j jd jdz \check{s} s\chi odz \check{s}ed)}$	$\underline{x_1' s\chi odz}$	$\underline{x_1^\epsilon(do ob odz po przy we wy)(jad jedz jes\chi jezdz)}$	$\underline{x_1' jes\chi}$
$\underline{x_2^\epsilon(odz odze po przy we wy)(j jd jdz \check{s} s\chi odz \check{s}ed)}$	$\underline{x_2' s\chi odz}$	$obejrz$	$odze(d jl l s\chi odz)$
		$oglqd$	$odzs\chi odz$

where  $\mathbf{X}'_1$  results from doing  $\text{przy} \rightarrow \pi rzy$  on  $\mathbf{X}_1$  and  $\mathbf{X}'_2$  results from doing  $\text{odze} \rightarrow \text{odz}$ ,  $\text{przy} \rightarrow \pi rzy$  on  $\mathbf{X}_2$ .

### 7.1.2 Admissible mutation and approximate clustering

In Polish (as well as Russian and many other Slavic languages), fleeting vowels may appear or disappear in inflected forms of the same word. We need to heuristically detect and remove these fleeting vowels to achieve better clustering results. Meanwhile, in Polish (as well as Russian and many other Slavic languages), one encounters vowel alternations in verb conjugations, as in Spanish or the Germanic languages.

**Algorithm 7.8** (Polish vowel blotting). *For a token string  $\hat{\sigma}$ , its blotted form  $\text{BlotFltV}(\hat{\sigma})$  is constructed as follows:*

- (1)  $\text{Do } \sim \mathbf{C}(e|jie|jio)^{\hat{\chi}^\epsilon}(c|k|l|n|r|s|j|w) \rightarrow \mathbf{C}\hat{\chi}.$
- (2)  $\text{Do } \sim \hat{\chi}kier \rightarrow \hat{\chi}kr.$
- (3)  $\text{Do } q \rightarrow \check{e}, i(a|o) \rightarrow ie.$

In what follows, we will construct a bivariate Boolean-valued function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  on a “simple heredity test function” in Algorithm 7.9, and a set of “admissible suffix mismatch” rules in Algorithm 7.10.

**Algorithm 7.9** (Simple heredity test). Let  $\hat{\alpha}'$  be a string obtained from  $\hat{\alpha}$  by the following steps:

- (1) Break down  $\hat{\alpha} = \hat{\alpha}_1 \hat{\alpha}_2$  into the concatenation of two strings  $\hat{\alpha}_1 = \hat{\alpha}^{[\min\{3, \ell(\hat{\alpha})\}]}$  (see the notation in Definition 3.1) and  $\hat{\alpha}_2$ , where the length of the first string  $\ell(\hat{\alpha}_1) = \min\{3, \ell(\hat{\alpha})\}$  is equal to 3 or the length of  $\hat{\alpha}$ , whichever is shorter.
- (2) On  $\hat{\alpha}_2$ , perform the following substitutions in a sequel:
  - (2.1) Do  $\sim \mathbf{C}_m \rightarrow \mathbf{X}'$ , where  $\mathbf{X}'$  is obtained from  $\mathbf{C}_m$  by deleting all the occurrences of the letter ‘‘j’’.
  - (2.2) Do  $\sim(n|(\mathbf{V}_{m_0}(l|n|\check{s}|t|w)_{m_0}))_m \rightarrow \emptyset$ .

The result after these two steps of operations is called  $\hat{\alpha}'_2$ .

- (3) Construct  $\hat{\alpha}' = \text{BlotFltV}(\hat{\alpha}_1 \hat{\alpha}'_2)$ .

Define  $\hat{\beta}'$  similarly. The Boolean-valued function  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if the lowercase form of  $\hat{\alpha}$  contains at least one instance of {q, a, e, ę, i, l, n, o, ó, r, u, y} **AND** at least one of the following six conditions holds:<sup>89</sup>

- (i)  $\hat{\alpha} = \hat{\beta}$ ;
- (ii)  $\hat{\beta} = \hat{\alpha}in$ ;
- (iii)  $\hat{\beta} = \hat{\alpha}n$ ;
- (iv)  $\hat{\beta} = \hat{\alpha}z$ ;
- (v)  $\min\{\ell(\hat{\alpha}'), \ell(\hat{\beta}')\} \geq 3$  **AND** ( $\hat{\alpha} = \hat{\beta}'$  **OR**  $\hat{\alpha}' = \hat{\beta}$  **OR**  $\hat{\alpha}' = \hat{\beta}'$  **OR**  $\hat{\beta}' = \hat{\alpha}'z$ ).
- (vi)  $\ell(\hat{\beta}) > \ell(\hat{\alpha}) \geq \frac{\ell(\hat{\beta})}{3}$  **AND**  $\hat{\beta} = \hat{\alpha}(a|e|o)_m(m|t)$ .

In what follows, we define  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{NW}^*(\hat{\alpha}, \hat{\beta})$  and  $\text{SuffixSW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{RootSW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{SW}^*(\hat{\alpha}, \hat{\beta})$  as what is done in the Danish case (Algorithm 5.7), which was also applied to two other Germanic languages treated in §5. Particular to the Polish case, we further define  $\text{SuffixNW}'(\hat{\alpha}, \hat{\beta})$  by performing the substitutions  $(e|ie)\sim \rightarrow \emptyset$ ,  $\sim j \rightarrow \emptyset$ ,  $\sim jc \rightarrow c$  to both components of  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta})$ , in a single sweep. The notation  $\text{SuffixSW}'(\hat{\alpha}, \hat{\beta})$  is defined similarly. If  $\text{SuffixNW}'(\hat{\alpha}, \hat{\beta}) = [\hat{\tau}_1, \hat{\tau}_2]$ , then  $\text{SuffixNW}'(\hat{\beta}, \hat{\alpha}) = [\hat{\tau}_2, \hat{\tau}_1]$ .

**Algorithm 7.10** (Admissible suffix mismatch and vowel alternation). For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function

$$\text{AdmMut}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta}))$$

returns **TRUE** if  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \emptyset$  **OR**  $[(a, e)][(a, o)][(q, ę)][(e, o)]$  **AND**  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$  contains at least one instance of {q, a, e, ę, i, l, n, o, ó, r, u, y} **AND** at least one of the following three conditions holds:

- (i)  $\text{SuffixNW}'(\hat{\alpha}, \hat{\beta}) = [\emptyset, \emptyset][[\emptyset, b]][[\emptyset, e]][[\emptyset, ie]][[c, k]][[c, t]][[\check{c}, k]][[dz, g]][[jc, t]][[s, \check{s}]][[s\chi, s]][[s\chi, \check{s}]][[z, \check{z}]][[zd, \check{zd}z]][[zg, \check{zd}\check{z}]]$ ;
- (ii)  $\text{SuffixNW}'(\hat{\beta}, \hat{\alpha}) = [\emptyset, \emptyset][[\emptyset, b]][[\emptyset, e]][[\emptyset, ie]][[c, k]][[c, t]][[\check{c}, k]][[dz, g]][[jc, t]][[s, \check{s}]][[s\chi, s]][[s\chi, \check{s}]][[z, \check{z}]][[zd, \check{zd}z]][[zg, \check{zd}\check{z}]]$ ;
- (iii)  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = [\emptyset|(\emptyset|j|l)\mathbf{V}_{m_0}(\emptyset|m|sj)(\emptyset|m|cj), \emptyset|(\emptyset|j|l)\mathbf{V}_{m_0}(\emptyset|m|sj)(\emptyset|m|cj)]$ .

Similarly, one can evaluate another Boolean-valued function

$$\text{AdmMut}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta}))$$

by trading all the occurrences of NW in the statements above with SW.

**Algorithm 7.11** (Heredity test function). For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , define  $\hat{\alpha}'$  and  $\hat{\beta}'$  as in Algorithm 7.9. The Boolean-valued function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if at least one of the following three conditions holds:

- (i)  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta}) = \text{TRUE}$ ;
- (ii)  $\ell(\text{RootNW}(\hat{\alpha}, \hat{\beta})) \geq \frac{\max\{\ell(\hat{\alpha}'), \ell(\hat{\beta}')\}}{2}$  **AND**  $\text{AdmMut}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta})) = \text{TRUE}$ ;
- (iii)  $\ell(\text{RootSW}(\hat{\alpha}, \hat{\beta})) \geq \frac{\max\{\ell(\hat{\alpha}'), \ell(\hat{\beta}')\}}{2}$  **AND**  $\text{AdmMut}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta})) = \text{TRUE}$ .

<sup>89</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical **OR**.

**Algorithm 7.12** (Polish verb aspect test). Let  $\text{NW}(\hat{\alpha}, \hat{\beta})$  be the result of performing Needleman–Wunsch alignment on strings  $\hat{\alpha}$  and  $\hat{\beta}$ . Let  $\mathbf{X}^*$  be an arbitrary non-empty string. The Boolean-valued function  $\text{VbAspTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if at least one of the following three conditions holds:

- (i)  $\text{NW}(\hat{\alpha}, \hat{\beta}) = \mathbf{X}^*$  (in other words,  $\hat{\alpha} = \hat{\beta}$ );
- (ii)  $\text{NW}(\hat{\alpha}, \hat{\beta}) = [\mathbf{P}, \mathbf{P}] \mathbf{X}^*$  (in other words,  $\hat{\alpha}$  agrees with  $\hat{\beta}$  up to a pair of word initial mismatching strings, both of which are possible prefixes of Polish verbs);
- (iii)  $\text{NW}(\hat{\alpha}, \hat{\beta}) = [\mathbf{P}, \mathbf{P}] \mathbf{X}^*([a, e][[a, o][[q, e][[e, o]]] \mathbf{X}^*$ , where the two occurrences of  $\mathbf{X}^*$  may or may not represent the same non-empty string.

Note that Smith–Waterman alignment is not used in this test.

The first two stages in the following approximate clustering algorithm for Polish words are very similar to the English counterpart (Algorithm 4.15), while the third stage is tailored for the verb aspects in Polish.

**Algorithm 7.13** (Approximate clustering of Polish words). The approximate clustering of a list of Polish words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  is completed in three stages:

- (1) We sort the list  $\{(\hat{\alpha}_1, \text{EffSpell}(\hat{\alpha}_1)), \dots, (\hat{\alpha}_N, \text{EffSpell}(\hat{\alpha}_N))\}$  alphabetically according to the second component (effective spelling) of each entry. If two consecutive neighbors in the alphabetized list  $\{(\hat{\alpha}_{(1)}, \text{EffSpell}(\hat{\alpha}_{(1)})), \dots, (\hat{\alpha}_{(N)}, \text{EffSpell}(\hat{\alpha}_{(N)}))\}$  satisfies  $\text{HrdTest}(\text{EffSpell}(\hat{\alpha}_{(n)}), \text{EffSpell}(\hat{\alpha}_{(n+1)})) = \text{FALSE}$ , where  $n \in \mathbb{Z} \cap [1, N]$ , then add a demarcation line between these two entries. In this way, the alphabetized list is divided into separate groups of words tagged with their effective spellings:  $\{g_1 = \{(\hat{\alpha}_{(1,1)}, \text{EffSpell}(\hat{\alpha}_{(1,1)})), \dots, (\hat{\alpha}_{(1,n_1)}, \text{EffSpell}(\hat{\alpha}_{(1,n_1)}))\}, \dots, g_M = \{(\hat{\alpha}_{(M,1)}, \text{EffSpell}(\hat{\alpha}_{(M,1)})), \dots, (\hat{\alpha}_{(M,n_M)}, \text{EffSpell}(\hat{\alpha}_{(M,n_M)}))\}\}$ , where each sublist also preserves alphabetic order. In particular, we have  $\hat{\alpha}_{(1,1)} = \hat{\alpha}_{(1)}$  and  $\hat{\alpha}_{(M,n_M)} = \hat{\alpha}_{(N)}$ .
- (2) For each group of words  $g_m = \{(\hat{\alpha}_{(m,1)}, \text{EffSpell}(\hat{\alpha}_{(m,1)})), \dots, (\hat{\alpha}_{(m,n_m)}, \text{EffSpell}(\hat{\alpha}_{(m,n_m)}))\}$  where  $m \in \mathbb{Z} \cap [1, M]$ , we augment it into a tagged entry  $G_m = (g_m, \text{EffSpell}(\hat{\alpha}_{(m,1)}), \text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})), \text{BlotFltV}(\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))))$ . The list  $\{G_1, \dots, G_M\}$  is sorted alphabetically, with respect to  $\text{BlotFltV}(\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})))$  (with highest priority),  $\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))$  (with medium priority), and  $\text{EffSpell}(\hat{\alpha}_{(m,1)})$  (with lowest priority). If two consecutive neighbors in the alphabetized list  $\{G_{(1)} = (g_{(1)}, \hat{\gamma}'_{(1)}, \hat{\gamma}''_{(1)}, \hat{\gamma}'''_{(1)}), \dots, G_{(M)} = (g_{(M)}, \hat{\gamma}'_{(M)}, \hat{\gamma}''_{(M)}, \hat{\gamma}'''_{(M)})\}$  satisfy

$$\text{SimpHrdTest}(\hat{\gamma}'''_{(m+1)}, \hat{\gamma}'''_{(m)}) = \text{FALSE}$$

AND

$$\text{SimpHrdTest}(\hat{\gamma}'''_{(m)}, \hat{\gamma}'''_{(m+1)}) = \text{FALSE}$$

AND

$$\text{HrdTest}(\hat{\gamma}''_{(m+1)}, \hat{\gamma}''_{(m)}) = \text{FALSE}$$

AND

$$\text{HrdTest}(\hat{\gamma}''_{(m)}, \hat{\gamma}''_{(m+1)}) = \text{FALSE}$$

AND

$$\text{HrdTest}(\hat{\gamma}'_{(m)}, \hat{\gamma}'_{(m+1)}) = \text{FALSE}$$

where  $m \in \mathbb{Z} \cap [1, M]$ , then add a demarcation line between these two entries. In this manner, the alphabetized list  $\{G_{(1)}, \dots, G_{(M)}\}$  is divided into separate groups of tagged clusters  $\{\Gamma_1 = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, \Gamma_K = \{G_{(K,1)}, \dots, G_{(K,m_K)}\}\}$ . Generate a list of word clusters  $\{\check{\Gamma}_1 = (\check{G}_{(1,1)}, \dots), \dots, \check{\Gamma}_K = (\check{G}_{(K,1)}, \dots)\}$  by discarding all the tags (effective spellings, essential roots, vowel blotted forms) from each  $\Gamma_k$ , where  $k \in \mathbb{Z} \cap [1, K]$ .

- (3) For each word cluster  $\check{\Gamma}_k = (\check{G}_{(k,1)}, \dots)$ , we augment it into a tagged entry

$$\mathcal{G}_k \equiv (\check{\Gamma}_k, \check{\Gamma}'_k, \check{\Gamma}''_k) := (\check{\Gamma}_k, \text{EssRoot}(\text{EffSpell}(\check{G}_{(k,1)})), \text{VbEssRoot}(\text{EssRoot}(\text{EffSpell}(\check{G}_{(k,1)})))).$$

Recall that  $\hat{\sigma}^{-1}$  is the reverse of the string  $\hat{\sigma}$  (Definition 3.1). The list  $\{\mathcal{G}_1 = (\check{\Gamma}_1, \check{\Gamma}'_1, \check{\Gamma}''_1), \dots, \mathcal{G}_K = (\check{\Gamma}_K, \check{\Gamma}'_K, \check{\Gamma}''_K)\}$  is sorted alphabetically, with respect to  $(\check{\Gamma}''_k)^{-1}$  (with higher priority) and  $\check{\Gamma}'_k$  (with lower priority). If two consecutive neighbors in the alphabetized list  $\{\mathcal{G}_{(1)} = (\check{\Gamma}_{(1)}, \check{\Gamma}'_{(1)}, \check{\Gamma}''_{(1)}), \dots, \mathcal{G}_{(K)} = (\check{\Gamma}_{(K)}, \check{\Gamma}'_{(K)}, \check{\Gamma}''_{(K)})\}$  satisfy

$$\text{VbAspTest}(\hat{\Gamma}''_{(k)}, \hat{\Gamma}''_{(k+1)}) = \text{FALSE}$$

AND

$$\text{VbAspTest}(\hat{\Gamma}'_{(k)}, \hat{\Gamma}'_{(k+1)}) = \text{FALSE}$$

where  $k \in \mathbb{Z} \cap [1, K]$ , then add a demarcation line between these two entries. In this manner, the alphabetized list  $\{\mathcal{G}_{(1)}, \dots, \mathcal{G}_{(K)}\}$  is divided into separate groups of tagged clusters  $\{\tilde{\Gamma}_1 = \{\mathcal{G}_{(1,1)}, \dots, \mathcal{G}_{(1,n_1)}\}, \dots, \tilde{\Gamma}_J = \{\mathcal{G}_{(J,1)}, \dots, \mathcal{G}_{(J,n_J)}\}\}$ , from which the output list of word clusters is generated, after removal of all the tags from each  $\tilde{\Gamma}_j$ , where  $j \in \mathbb{Z} \cap [1, J]$ .

*Example 7.13.1.* We refer to the Polish version of Wiktionary ([https://pl.wiktionary.org/wiki/Kategoria:Gramatyka\\_j%C4%99zyka\\_polskiego](https://pl.wiktionary.org/wiki/Kategoria:Gramatyka_j%C4%99zyka_polskiego)), for the divisions of declension and conjugation classes, of nouns and verbs respectively.

Here are some nouns that display the 1st masculine declension patterns:

*drani, drania, draniach, draniami, dranie, draniem, draniom, draniowi, draniów, draniu, drań* — “rascal”;

*goryl, goryla, gorylach, gorylami, goryle, gorylem, goryli, gorylom, gorylowi, gorylu* — “gorilla”;

*liści, liścia, liściach, liście, liściem, liściom, liściowi, liściu, liść, liścimi* — “leaf”.

2nd masculine declension:

*owoc, owocach, owocami, owoce, owocem, owocom, owocowi, owoców, owocu* — “fruit”;

*stróż, stróża, stróżach, stróżami, stróżem, stróżom, stróżowi, stróżowie, stróżu, stróży* — “guard”;

*tasiemca, tasiemcach, tasiemcamy, tasiemce, tasiemcem, tasiemcom, tasiemcowi, tasiemców, tasiemcu, tasiemiec* — “tapeworm”.

3rd masculine declension:

*chirurdzy, chirurg, chirurga, chirurgach, chirurgami, chirurgiem, chirurgom, chirurgowi, chirurgów, chirurgu* — “surgeon”;

*kożuch, kożucha, kożuchach, kożuchami, kożuchem, kożuchom, kożuchowi, kożuchów, kożuchu, kożuchy* — “sheepskin”;

*strażacy, strażak, strażaka, strażakach, strażakami, strażakiem, strażakom, strażakowi, strażaków, strażaku* — “fireman”.

4th masculine declension:

*chleb, chleba, chlebach, chlebami, chlebem, chlebie, chlebom, chlebowi, chlebów, chleby* — “bread”;

*doktor, doktora, doktorach, doktorami, doktorem, doktorom, doktorowi, doktorowie, doktorów, doktorze, doktorzy* — “doctor”;

*drozd, drozda, drozdach, drozdami, drozdem, drozdom, drozdowi, drozdów, drozdy, droździe* — “thrush”.

5th masculine declension:

*Marsjan, Marsjanach, Marsjanami, Marsjanie, Marsjanin, Marsjanina, Marsjaninem, Marsjaninie, Marsjaninowi, Marsjanom, Marsjany* — “Martian”.

1st neuter declension:

*pola, polach, polami, pole, polem, polom, polu, pól* — “field”;

*studia, studiach, studiami, studiem, studio, studiom, studiów, studiu* — “studio”.

2nd neuter declension:

*pluc, pluca, plucach, plucami, plucem, pluco, plucom, plucu* — “lung”.

3rd neuter declension:

*cuda, cudach, cudami, cudem, cudo, cudom, cudów, cudu, cudzie* — “wonder”;

*palcie, palt, palta, paltach, paltami, paltem, palto, paltom, paltu* — “overcoat”.

4th neuter declension:

*dziewczętę, dziewczę, dziewczęcia, dziewczęciem, dziewczęciu, dziewczęta, dziewczętach, dziewczętami, dziewczętom — “girl”.*

5th neuter declension:

*ramienia, ramieniem, ramieniu, ramię, ramion, ramiona, ramionach, ramionami, ramionom — “shoulder”.*

6th neuter declension:

*muzea, muzeach, muzeami, muzeom, muzeów, muzeum — “museum”.*

1st feminine declension:

*pani, paniach, paniami, panią, panie, paniom, pań — “lady”;*

*świni, świnia, świniach, świniami, świnią, świnie, świnę, świnio, świniom, świn — “pig”.*

2nd feminine declension:

*prac, praca, pracach, pracami, pracą, prace, prace, praco, pracom, pracy — “job”;*

*sadza, sadzach, sadzami, sadzą, sadze, sadzę, sadzo, sadzom, sadzy — “soot”.*

3rd feminine declension:

*much, mucha, muchach, muchami, muchą, muchę, mucho, muchom, muchy, musze — “bowtie”;*

*plotce, plotek, plotka, plotkach, plotkami, plotką, plotkę, plotki, plotko, plotkom — “rumour”.*

4th feminine declension:

*form, forma, formach, formami, formą, formę, formie, formo, formom, formy — “shape”.*

5th feminine declension:

*brew, brwi, brwiach, brwiami, brwią, brwiom — “eyebrow”;*

*kolei, kolej, kolejach, kolejami, koleją, kolej, kolejom — “rail”;*

*kości, kościach, kością, kościom, kość, kośćmi — “bone”.*

6th feminine declension:

*klacz, klaczach, klaczami, klaczą, klacze, klaczom, klaczy — “mare”;*

*mysz, myszach, myszami, myszą, myszom, myszy — “mouse”.*

As we send these nouns into our clustering algorithm, we receive the following results:

{*brew, brwi, brwiach, brwiami, brwią, brwiom*}, {*chirurdzy, chirurg, chirurga, chirurgach, chirurgami, chirurgiem, chirurgom, chirurgowi, chirurgów, chirurgu*}, {*chleb, chleba, chlebach, chlebami, chlebem, chlebie, chlebom, chlebowi, chlebów, chleby*}, {*cuda, cudach, cudami, cudem, cudo, cudom, cudów, cudu, cudzie*}, {*doktor, doktora, doktorach, doktorami, doktorem, doktorom, doktorowi, doktorowie, doktorów, doktorze, doktorzy*}, {*drani, drania, draniach, draniami, dranie, draniem, draniom, draniowi, draniów, draniu, drań*}, {*drozd, drozda, drozdach, drozdam, drozdem, drozdom, drozdow, drozdów, drozdy, droździe*}, {*dziewczętę, dziewczę, dziewczęcia, dziewczęciem, dziewczęciu, dziewczęta, dziewczętach, dziewczętami, dziewczętom*}, {*form, forma, formach, formami, formą, formę, formie, formo, formom, formy*}, {*goryl, goryla, gorylach, gorylami, goryle, gorylem, goryli, goryлом, goryлови, goryлу*}, {*klacz, klaczach, klaczami, klaczą, klacze, klaczom, klaczy*}, {*kolei, kolej, kolejach, kolejami, koleją, kolej, kolejom*}, {*kości, kościach, kością, kościom, kość, kośćmi*}, {*kożuch, kożucha, kożuchach, kożuchami, kożuchem, kożuchom, kożuchowi, kożuchów, kożuchu, kożuchy*}, {*liści, liścia, liściach, liście, liściem, liściom, liściowi, liściu, liść, liścimi*}, {*Marsjan, Marsjanach, Marsjanami, Marsjanie, Marsjanin, Marsjanina, Marsjaninem, Marsjaninie, Marsjanowi, Marsjanom, Marsjany*}, {*much, mucha, muchach, muchami, muchą, muchę, mucho, muchom, many, musze*}, {*muzea, muzeach, muzeami, muzeom, muzeów, muzeum*}, {*mysz, myszach, myszami, myszą, myszom, myszy*}, {*owoc, owocach, owocami, owoce, owocom, owocom*,

*owocowi, owoców, owocu}, {palcie, palt, palta, paltach, paltami, paltem, palto, paltom, paltu}, {pani, paniach, paniami, panią, panie, paniom, pań}, {plotce, plotek, plotka, plotkach, plotkami, plotką, plotkę, plotki, plotko, plotkom}, {pluc, pluca, plucach, plucami, plucem, pluco, plucom, plucu}, {pola, polach, polami, pole, polem, polom, polu, pól}, {prac, praca, pracach, pracami, pracą, prace, pracę, praco, pracom, pracy}, {ramienia, ramieniem, ramieniu, ramię, ramion, ramiona, ramionach, ramionami, ramionom}, {sadza, sadzach, sadzami, sadzą, sadze, sadzę, sadzo, sadzom, sadzy}, {strażacy, strażak, strażaka, strażakach, strażakami, strażakiem, strażakom, strażakowi, strażaków, strażaku}, {stróż, stróża, stróżach, stróżami, stróżem, stróżom, stróżowi, stróżowie, stróżu, stróży}, {studia, studiach, studiami, studiem, studio, studiom, studiów, studiu}, {świni, świnia, świnach, świniami, świną, świnie, świnę, świnio, świnom, świń}, {tasiemca, tasiemcach, tasiemcami, tasiemce, tasiemcem, tasiemcom, tasiemcowi, tasiemcu, tasiemiec}.*

We pick the following Polish verbs, followed by their English translations (enclosed in quotation marks) and conjugation classes (enclosed in parentheses):

*czyta, czytacie, czytać, czytaj, czytając, czytającą, czytającą, czytającą, czytajcie, czytajmy, czytali, czytaliby, czytalibyście, czytalibyśmy, czytaliście, czytaliśmy, czytał, czytala, czytałaby, czytałabym, czytałabyś, czytałam, czytałaś, czytaliby, czytaliby, czytałyś, czytałem, czytałeś, czytało, czytałoby, czytały, czytałyby, czytałybyście, czytałybyśmy, czytałyście, czytałyśmy, czytam, czytamy, czytana, czytane, czytani, czytanie, czytano, czytany, czytasz, przeczyta, przeczytacie, przeczytać, przeczytaj, przeczytaję, przeczytajcie, przeczytajmy, przeczytali, przeczytaliby, przeczytalibyście, przeczytalibyśmy, przeczytaliście, przeczytaliśmy, przeczytał, przeczytała, przeczytały, przeczytałabym, przeczytałybyś, przeczytałam, przeczytałaś, przeczytałyby, przeczytałbym, przeczytałyś, przeczytałem, przeczytałeś, przeczytało, przeczytałyby, przeczytały, przeczytałyby, przeczytałybyście, przeczytałybyśmy, przeczytałyście, przeczytałyśmy, przeczytam, przeczytamy, przeczytana, przeczytane, przeczytani, przeczytanie, przeczytano, przeczytany, przeczytasz, przeczytawszy — “read” (1);*

*umiał, umiała, umialaby, umialabym, umialabyś, umiałam, umiałaś, umialby, umialbym, umialbyś, umiałem, umialeś, umiało, umialoby, umiali, umialyby, umialbyście, umialbyśmy, umialyście, umialyśmy, umie, umiecie, umieć, umiej, umieją, umiejąc, umiejąca, umiejace, umiejący, umiejcie, umiejmy, umielii, umieliby, umielibyście, umielibyśmy, umielisicie, umielismy, umiem, umiemy, umiesz — “can” (2);*

*taniały, taniał, taniała, taniałaby, taniałabym, taniałamy, taniałaś, taniałby, taniałbym, taniałbyś, taniałe, tanialem, taniałeś, taniało, taniałoby, taniałobym, taniałobyś, taniałom, taniałoś, taniały, taniałyby, taniałybyście, taniałybyśmy, taniałyście, taniałyśmy, taniano, tanieć, taniej, tanieją, taniejąc, taniejąca, taniejace, taniejący, taniejcie, tanieje, taniejecie, taniejemy, taniejesz, tanieję, taniejmy, tanieli, tanieliby, tanielibyście, tanielibyśmy, tanielisicie, tanielismy, tanienie — “cheapen” (3);*

*malować, malowali, malowaliby, malowalibyście, malowalibyśmy, malowaliście, malowaliśmy, malował, malowała, malowałyby, malowałabym, malowałybyś, malowałam, malowałaś, malowałyby, malowałbym, malowałybyś, malowałem, malowałeś, malowało, malowałyby, malowały, malowałyby, malowałybyś, malowałyście, malowałyśmy, malowana, malowane, malowani, malowanie, malowan, malowany, maluj, malują, malując, malującą, malującę, malując, maluje, malujecie, malujemy, malujesz, maluję, malujmy, pomalować, pomalowali, pomalowaliby, pomalowalibyście, pomalowalibyśmy, pomalowaliście, pomalowaliśmy, pomalował, pomalowała, pomalowałyby, pomalowałabym, pomalowałybyś, pomalowałam, pomalowałaś, pomalowały, pomalowałybym, pomalowałybyś, pomalowałem, pomalowałeś, pomalowało, pomalowałyby, pomalowałybyś, pomalowałybyście, pomalowałybyśmy, pomalowałyście, pomalowałyśmy, pomalowana, pomalowane, pomalowani, pomalowanie, pomalowan, pomalowany, pomalowawszy, pomaluj, pomalując, pomalujcie, pomalujecie, pomalujemy, pomalujesz, pomaluję, pomalujmy — “portray” (4);*

*ciągną, ciągnąc, ciągnąca, ciągnące, ciągnący, ciągnącą, ciągnął, ciągnąłby, ciągnąłbym, ciągnąłbyś, ciągnąłem, ciągnąłeś, ciągnę, ciągnęli, ciągnęliby, ciągnęlibyście, ciągnęlibyśmy, ciągnęliście, ciągnęliśmy, ciągnęła, ciągnęlaby, ciągnęlabym, ciągnęlabyś, ciągnęłam, ciągnęłaś, ciągnęło, ciągnęloby, ciągnęły, ciągnęłyby, ciągnęłybyście, ciągnęłybyśmy, ciągnęłyście, ciągnęłyśmy, ciągnie, ciągnicie, ciągnimy, ciągnieni, ciągniesz, ciągnięci, ciągnicie, ciągnięta, ciągnięte, ciągnięto, ciągnęty, ciągnij, ciągnicie, ciągnijmy, ciągniona, ciągnione, ciągniony — “pull” (5a);*

*płyna, płynąc, płynąca, płynące, płynący, płynnąć, płynał, płynąłby, płynąłbym, płynąłbyś, płynąłem, płynąłeś, płynę, płynęli, płynęliby, płynęlibyście, płynęlibyśmy, płynęliście, płynęliśmy, płynęła, płynęlaby, płynęlabym, płynęlabyś, płynęłam, płynęłaś, płynęło, płynęloby, płynęły, płynęłyby, płynęłybyście, płynęłybyśmy, płynęłyście, płynęłyśmy, płynie, płyniecie, płyniemy, płyniesz, płynięcie, płynięto, płynię, płynicie, płynny — “swim” (5b);*

*chudli, chudliby, chudlibyście, chudlibyśmy, chudliście, chudliśmy, chudł, chudła, chudlaby, chudlabym, chudlabyś, chudłam, chudłaś, chudlb, chudlbym, chudlbys, chudlem, chudleś, chudło, chudloby, chudły, chudłyby, chudłybyście, chudłybyśmy, chudłyście, chudłyśmy, chudnq, chudnqc, chudnaca, chudnace, chudnacy, chudnacq, chudnq, chudnie, chudniecie, chudniemy, chudniesz, chudnij, chudnijcie, chudnijmy — “lose weight” (5c);*

*robi, robią, robiąc, robiąca, robiące, robiący, robicie, robić, robieni, robienie, robię, robili, robiliby, robilibyście, robilibyśmy, robiliście, robiliśmy, robil, robila, robilaby, robilabym, robilabyś, robilam, robilaś, robilby, robilbym, robilbyś, robilem, robileś, robilo, robiloby, robilobym, robilobyś, robilom, robilos, robily, robilby, robilbyście, robilbyśmy, robilyscie, robilysmy, robimy, robiona, robione, robiono, robiony, robisz, zrobi, zrobia, zrobicie, zrobic, zrobieni, zrobienie, zrobie, zrobili, zrobiliby, zrobilibyście, zrobilibyśmy, zrobiliście, zrobiliśmy, zrobil, zrobila, zrobilaby, zrobilabym, zrobilabyś, zrobilam, zrobilaś, zrobilby, zrobilbym, zrobilbyś, zrobilem, zrobileś, zrobilo, zrobiloby, zrobilobym, zrobilobyś, zrobilom, zrobilos, zrobily, zrobilby, zrobilibyście, zrobilibyśmy, zrobimy, zrobiona, zrobione, zrobiono, zrobiony, zrobisz, zrobiszy — “do” (6a);*

*uwierz, uwierzą, uwierzcie, uwierzenie, uwierzę, uwierzmy, uwierzono, uwierzy, uwierzycie, uwierzyc, uwierzyli, uwierzylby, uwierzylbyście, uwierzylbyśmy, uwierzyliszy, uwierzyliszy, uwierzyl, uwierzyla, uwierzylaby, uwierzylabym, uwierzylabyś, uwierzylam, uwierzylaś, uwierzylby, uwierzylbym, uwierzylbyś, uwierzylem, uwierzyleś, uwierzyo, uwierzloby, uwierzlobym, uwierzlobys, uwierzylom, uwierzyoś, uwierzyl, uwierztyby, uwierztybyście, uwierztybyśmy, uwierztylysie, uwierztylysy, uwierzymy, uwierzysz, uwierzyszy, wierz, wierzq, wierzqc, wierzaca, wierzace, wierzacy, wiercie, wierzenie, wierz, wierzmy, wierzono, wierz, wierzycie, wierzyć, wierzli, wierzliby, wierzlibyście, wierzlibyśmy, wierzliście, wierzliisy, wierzyl, wierzyla, wierzlaby, wierzlabym, wierzlabys, wierzylam, wierzylaś, wierzylby, wierzlbym, wierzlbys, wierzylem, wierzyleś, wierzyo, wierzloby, wierzlobym, wierzlobys, wierzylom, wierzyoś, wierzyl, wierzlyby, wierzlybyście, wiezybysy, wierzlyscie, wierzlysy, wierzymy, wierzysz — “believe” (6b);*

*widzą, widząc, widząca, widzące, widzący, widzenie, widzę, widział, widziała, widzialaby, widzialabym, widzialabyś, widzialam, widzialaś, widzialby, widzialbym, widzialbyś, widzialem, widzialeś, widziale, widzialoby, widzialobym, widzialobyś, widzialom, widzialeś, widzialy, widzialyby, widzialybyście, widzialybyśmy, widzialyście, widzialyśmy, widziana, widziane, widziani, widziano, widziany, widzicie, widzieć, widzieli, widzieliby, widzielibyście, widzielibyśmy, widzieliście, widzieliśmy, widzimy, widzisz, widź, widżcie, widżmy, zobacz, zobaczą, zobaczcie, zobaczenie, zobaczę, zobaczmy, zobaczono, zobaczy, zobaczycie, zobaczyć, zobaczyli, zobaczyliby, zobaczylibyście, zobaczylibyśmy, zobaczyliście, zobaczyliśmy, zobaczył, zobaczyła, zobaczylaby, zobaczylabym, zobaczylabys, zobaczyłam, zobaczyłaś, zobaczyby, zobaczylbym, zobaczylbys, zobaczyłem, zobaczyłeś, zobaczylo, zobaczyloby, zobaczyły, zobaczylyby, zobaczylybyście, zobaczylybyśmy, zobaczylyscie, zobaczylysy, zobaczymy, zobaczysz, zobaczywszy — “see” (7a);*

*leż, leżał, leżała, leżałaby, leżałabym, leżałam, leżałas, leżałby, leżałbym, leżałem, leżaleś, leżalo, leżaloby, leżalobym, leżalobyś, leżalom, leżaloś, leżaly, leżalyby, leżalybyście, leżalybyśmy, leżalyście, leżalyśmy, leżano, leżq, leżąc, leżaca, leżace, leżacy, leżcie, leżec, leżeli, leżeliby, leżelibyście, leżelibyśmy, leżeliście, leżeliśmy, leżenie, leż, leżmy, leży, leżycie, leżymy, leżysz — “lie” (7b);*

*odczyta, odczytacie, odczytać, odczytaj, odczytają, odczytajcie, odczytajmy, odczytali, odczytaliby, odczytalibyście, odczytalibyśmy, odczytaliście, odczytaliśmy, odczytal, odczytała, odczytalaby, odczytalabym, odczytalabyś, odczytałam, odczytałaś, odczytalby, odczytalbym, odczytalbyś, odczytałem, odczytales, odczytał, odczytaloby, odczytalobym, odczytalobyś, odczytalom, odczytałos, odczytały, odczytałyby, odczytałybyście, odczytałybyśmy, odczytałyście, odczytałyśmy, odczytam, odczytamy, odczytana, odczytane, odczytani, odczytanie, odczytano, odczytany, odczytasz, odczytawszy, odczytuj, odczytują, odczytując, odczytującą, odczytujące, odczytujący, odczytujcie, odczytuje, odczytujecie, odczytujemy, odczytujesz, odczytuję, odczytujmy, odczytywać, odczytywali, odczytywaliby, odczytywalibyście, odczytywalibyśmy, odczytywaliście, odczytywaliśmy, odczytywał, odczytywała, odczytywały, odczytywałabym, odczytywałaby, odczytywałabym, odczytywałabyś, odczytywałam, odczytywałaś, odczytywałby, odczytywałbym, odczytywałbyś, odczytywałem, odczytywałeś, odczytywało, odczytywałoby, odczytywałobym, odczytywałobyś, odczytywałom, odczytywałoś, odczytywały, odczytywałyby, odczytywałybyście, odczytywałybyśmy, odczytywałyście, odczytywałyśmy, odczytywana, odczytywane, odczytywani, odczytywanie, odczytywano, odczytywany — “decipher” (8a);*

*zyska, zyskacie, zyskać, zyskaj, zyskają, zyskajcie, zyskajmy, zyskali, zyskaliby, zyskalibyście, zyskalibyśmy, zyskaliście, zyskaliśmy, zyskał, zyskała, zyskalaby, zyskalabym, zyskalabyś, zyskałam, zyskałaś, zyskaliby, zyskalibyś, zyskałam, zyskałes, zyskał, zyskałoby, zyskałobym, zyskałobys, zyskałom, zyskałos, zyskały, zyskałyby, zyskałybyście, zyskałybyśmy, zyskałyście, zyskałyśmy, zyskam, zyskamy, zyskana, zyskani, zyskanie,*

*zyskano, zyskany, zyskasz, zyskawszy, zyskiwać, zyskiwali, zyskiwaliby, zyskiwalibyście, zyskiwalibyśmy, zyskiwaliście, zyskiwaliśmy, zyskiwał, zyskiwała, zyskiwałaby, zyskiwałabym, zyskiwałabyś, zyskiwałam, zyskiwałaś, zyskiwałby, zyskiwałbym, zyskiwałbyś, zyskiwałem, zyskiwałeś, zyskiwało, zyskiwałoby, zyskiwały, zyskiwałyby, zyskiwałybyście, zyskiwałybyśmy, zyskiwałyście, zyskiwałyśmy, zyskiwana, zyskiwane, zyskiwani, zyskiwanie, zyskiwano, zyskiwany, zyskuj, zyskują, zyskując, zyskującą, zyskującce, zyskującący, zyskujcie, zyskuje, zyskujecie, zyskujemy, zyskujesz, zyskuję, zyskujmy* — “gain” (8a);

*mazać, mazali, mazaliby, mazalibyście, mazalibyśmy, mazaliście, mazaliśmy, mazal, mazała, mazałaby, mazałabym, mazałabyś, mazalam, mazałaś, mazały, mazalbym, mazałyś, mazalem, mazałeś, mazało, mazaloby, mazały, mazałyby, mazałybyście, mazałybyśmy, mazałyście, mazałyśmy, mazana, mazane, mazani, mazanie, mzano, mazany, maż, mażq, mażąc, mażąca, mażące, mażący, mażcie, maże, mażecie, mażemy, mażesz, mażę, mażmy* — “smear” (9);

*pici, picie, pić, pij, piją, pijąc, pijąca, pijące, pijący, pijcie, pije, pijecie, pijemy, pijesz, piję, pijmy, pili, piliby, pilibyście, pilibyśmy, piliście, piliśmy, pił, pila, piłaby, piłabym, piłabyś, piłam, piłaś, piłby, piłbym, piłbyś, piłem, piłeś, piło, piłoby, piłobym, piłobyś, piłom, piłoś, piły, piłyby, piłybyście, piłyśmy, piłyście, piłyśmy, pita, pite, pito, pity* — “drink” (10a);

*lać, lali, laliby, lalibyście, lalibyśmy, laliście, laliśmy, lat, lala, lataby, latabym, latabyś, lalam, lalaś, latby, latbym, latbyś, lałem, lałeś, lało, lałoby, laly, lalyby, lalybyście, lałybyśmy, lałyście, lałyśmy, lana, lane, lani, lanie, lano, lany, lej, leją, lejąc, lejącą, lejące, lejący, lejcie, leje, lejecie, lejemy, lejesz, leję, lejmy* — “pour” (10b);

*bierz, bierzcie, bierze, bierzecie, bierzemy, bierzesz, bierzmy, biorą, biorąc, biorąca, biorące, biorący, biorę, brać, brali, braliby, bralibyście, bralibyśmy, braliście, braliśmy, brał, brała, brałaby, brałabym, brałabyś, brałam, brałaś, brałby, brałbym, brałbyś, brałem, brałeś, brało, brałoby, brały, brałyby, brałybyście, brałybyśmy, brałyście, brałyśmy, brana, brane, brani, branie, brany, wezmą, wezmę, wezmiesz, wziąć, wziął, wziąłby, wziąłbym, wziąłbyś, wziąłem, wziąłeś, wziąwszy, wzięci, wzięcie, wzięli, wzięliby, wzięlibyście, wzięlibyśmy, wzięliście, wzięliśmy, wzięla, wzięlaby, wzięlabym, wzięlabyś, wzięłam, wzięłaś, wzięło, wzięloby, wzięły, wzięłyby, wzięłybyście, wzięłybyśmy, wzięłyście, wzięłyśmy, wzięta, wzięte, wzięto, wzięty* — “take” (10c);

*wiezie, wieziecie, wieziemy, wiezieni, wiezienie, wieziesz, wieziona, wiezione, wieziono, wieziony, wiozą, wioząc, wioząca, wiożące, wiożący, wiozę, wiozła, wiozłaby, wiozłabym, wiozłabyś, wiozłam, wiozłaś, wiozo, wiozłyby, wiozłyby, wiozłybyście, wiozłybyśmy, wiozłyście, wiozłyśmy, wióz, wiózby, wiózlbym, wiózbyś, wiózlem, wiózleś* — “carry” (11).

These verbs are clustered by our algorithm as follows:

{*bierz, bierzcie, bierze, bierzecie, bierzemy, bierzesz, bierzmy, biorą, biorąc, biorąca, biorące, biorący, biorę, brać, brali, braliby, bralibyście, bralibyśmy, braliście, braliśmy, brał, brała, brałaby, brałabym, brałabyś, brałam, brałaś, brałby, brałbym, brałbyś, brałem, brałeś, brało, brałoby, brały, brałyby, brałybyście, brałybyśmy, brałyście, brałyśmy, brana, brane, brani, branie, brany, wezmą, wezmę, wezmiesz, wziąć, wziął, wziąłby, wziąłbym, wziąłbyś, wziąłem, wziąłeś, wziąwszy, wzięci, wzięcie, wzięli, wzięliby, wzięlibyście, wzięlibyśmy, wzięliście, wzięliśmy, wzięla, wzięlaby, wzięlabym, wzięlabyś, wzięłam, wzięłaś, wzięło, wzięloby, wzięły, wzięłyby, wzięłybyście, wzięłybyśmy, wzięłyście, wzięłyśmy, wzięta, wzięte, wzięto, wzięty*},

{*chudli, chudliby, chudlibyście, chudlibyśmy, chudliście, chudliły, chudła, chudlaby, chudlabym, chudlabyś, chudlam, chudłaś, chudły, chudbym, chudbyś, chudlem, chudłeś, chudło, chudloby, chudły, chudbyby, chudłybyście, chudłybyśmy, chudłyście, chudłyśmy, chudnq, chudnac, chudnaca, chudnace, chudnacy, chudnac, chudnq, chudnie, chudniecie, chudniemy, chudniesz, chudnij, chudnijcie, chudnijmy*},

{*ciągną, ciągnąć, ciągnąca, ciągnące, ciągnący, ciągnąć, ciągnął, ciągnąłby, ciągnąłbym, ciągnąłbyś, ciągnąłtem, ciągnąłs, ciągnę, ciągneli, ciągneliwy, ciągneliwyście, ciągneliwyśmy, ciągneliście, ciągneliśmy, ciągnela, ciągnelab, ciągnelabym, ciągnelabys, ciągnelam, ciągnelaś, ciągnęlo, ciągneloby, ciągnęły, ciągnęłyby, ciągnęłybyście, ciągnęłybyśmy, ciągnęłyście, ciągnęłyśmy, ciągnie, ciągnicie, ciągniemy, ciągnieni, ciągniesz, ciągnęci, ciągnęcie, ciągniąta, ciągnięte, ciągnięto, ciągnięty, ciągnij, ciągnijcie, ciągnijmy, ciągniona, ciągnione, ciągniony*},

{*czyta, czytacie, czytać, czytaj, czytają, czytając, czytającą, czytające, czytający, czytajcie, czytajmy, czytali, czytaliby, czytalibyście, czytalibyśmy, czytaliscie, czytaliśmy, czytal, czytala, czytalaby, czytalabym, czytalabyś*,

czytałam, czytałaś, czytałby, czytałbym, czytałbyś, czytałem, czytałeś, czytało, czytałyby, czytały, czytałyby, czytałyście, czytałybyśmy, czytałyście, czytałyśmy, czytam, czytamy, czytana, czytane, czytani, czytanie, czytano, czytany, czytasz, przeczyta, przeczytacie, przeczytać, przeczytaj, przeczytają, przeczytajcie, przeczytajmy, przeczytali, przeczytaliby, przeczytalibyście, przeczytalibyśmy, przeczytaliście, przeczytaliśmy, przeczytal, przeczytala, przeczytalały, przeczytalałbym, przeczytalałyś, przeczytalaś, przeczytalały, przeczytalałbym, przeczytalałyś, przeczytalałem, przeczytalałeś, przeczytalały, przeczytalały, przeczytalałyby, przeczytalałybyście, przeczytalałybyśmy, przeczytalałyście, przeczytalałyśmy, przeczytala, przeczytana, przeczytane, przeczytani, przeczytanie, przeczytano, przeczytany, przeczytawszy},

{lać, lali, laliby, lalibyście, lalibyśmy, laliście, laliśmy, lał, lala, lałaby, lałabym, lałabyś, lałam, lałaś, lałby, lałbym, lałbyś, lałem, lałeś, lało, lałoby, lały, lałyby, lałybyście, lałybyśmy, lałyście, lałyśmy, lana, lane, lani, lanie, lano, lany, lej, leją, lejąc, lejącą, lejące, lejący, lejecie, leje, lejecie, lejemy, lejesz, leję, lejmy},

{leż, leżał, leżała, leżałaby, leżałabym, leżałabyś, leżałam, leżałaś, leżałby, leżałbym, leżałbyś, leżałem, leżałob, leżałobym, leżałobyś, leżałom, leżałos, leżałły, leżałby, leżałbyście, leżałbyśmy, leżałyscie, leżałysmy, leżano, leżą, leżąc, leżąca, leżace, leżący, leżcie, leżeć, leżeli, leżeliby, leżelibyście, leżelibyśmy, leżeliście, leżeliśmy, leżenie, leżę, leżmy, leży, leżcie, leżmy, leżysz},

{malować, malowali, malowaliby, malowalibyście, malowalibyśmy, malowaliście, malowaliśmy, malował, malowała, malowałaby, malowałabym, malowałabyś, malowałam, malowałaś, malowałby, malowałbym, malowałyś, malowalem, malowaleś, malowało, malowaloby, malowały, malowałyby, malowałybyście, malowałybyśmy, malowałyście, malowałyśmy, malowana, malowane, malowani, malowanie, malowan, malowany, maluj, malując, malującą, malującce, malujący, malujcie, maluje, malujecie, malujemy, malujesz, maluję, malujmy},

{mazać, mazali, mazaliby, mazalibyście, mazalibyśmy, mazaliście, mazaliśmy, mazał, mazała, mazałaby, mazałabym, mazałabyś, mazałam, mazałaś, mazałby, mazałbym, mazałbyś, mazałem, mazałes, mazało, mazałoby, mazały, mazałyby, mazałybyście, mazałybyśmy, mazałyście, mazałyśmy, mazana, mazane, mazani, mazanie, mazano, mazany, maż, mażą, mażąc, mażąca, mażace, mażacy, mażcie, mażecie, mażemy, mażesz, mażę, mażmy},

{odczyta, odczytacie, odczytać, odczytaj, odczytają, odczytacie, odczytajmy, odczytali, odczytaliby, odczytalibyście, odczytalibyśmy, odczytaliście, odczytaliśmy, odczytał, odczytała, odczytalaby, odczytalabym, odczytalabyś, odczytałam, odczytałaś, odczytały, odczytaliby, odczytałem, odczytałeś, odczytało, odczytaloby, odczytałoby, odczytałom, odczytałoś, odczytały, odczytałyby, odczytałybyście, odczytałybyśmy, odczytałyście, odczytałyśmy, odczytam, odczytamy, odczytana, odczytane, odczytani, odczytan, odczytano, odczytany, odczytasz, odczytawszy, odczytuj, odczytują, odczytując, odczytującą, odczytującce, odczytującący, odczytujcie, odczytuję, odczytujecie, odczytujemy, odczytujesz, odczytuję, odczytujmy, odczytywać, odczytywali, odczytywaliby, odczytywalibyście, odczytywalibyśmy, odczytywaliście, odczytywaliśmy, odczytywał, odczytywała, odczytywałaby, odczytywałabym, odczytywałabyś, odczytywałam, odczytywałaś, odczytywałby, odczytywałbym, odczytywałbyś, odczytywałem, odczytywałeś, odczytywało, odczytywałoby, odczytywałbym, odczytywałobyś, odczytywałom, odczytywałoś, odczytywały, odczytywałyby, odczytywałybyście, odczytywałybyśmy, odczytywałyście, odczytywałyśmy, odczytywana, odczytywane, odczytywani, odczytywanie, odczytywano, odczytywany},

{pici, picie, pić, pij, piją, pijąc, pijąca, pijące, pijący, pijcie, pije, pijecie, pijemy, pijesz, piję, pijmy, pili, piliby, pilibyście, pilibyśmy, piliście, piliśmy, pil, pila, pilaby, pilabym, pilabyś, pilam, pilaś, pilby, pilbym, pilbyś, pilem, pileś, pilo, piloby, pilobym, pilobyś, pilom, piłoś, pily, piloby, pilobyście, pilobyśmy, pilobyście, pita, pite, pito, pity},

{pływą, płynąć, płynąca, płynące, płynący, płynąć, płynał, płynąłby, płynąłbym, płynąłbyś, płynąłem, płynąłeś, płynę, płyneli, płyneliby, płynelibyście, płynelibyśmy, płyneliście, płyneliśmy, płynęła, płyneliby, płyneliby, płynelibyś, płynelabym, płynelam, płynelaś, płynęło, płyneloby, płynęły, płynęłyby, płynęłybyście, płynęłybyśmy, płynęłyście, płynęły, płynie, płynecie, płyniem, płyniesz, płynięcie, płynięto, płyn, płynie, płynny},

{pomalować, pomalowali, pomalowaliby, pomalowalibyście, pomalowalibyśmy, pomalowaliście, pomalowaliśmy, pomalował, pomalowała, pomalowały, pomalowały, pomalowałyś, pomalowałam, pomalowałaś, pomalowały, pomalowaliby, pomalowalbym, pomalowalbyś, pomalowalem, pomalowałeś, pomalowało, pomalowały, pomalowały, pomalowałyby, pomalowałybyście, pomalowałybyśmy, pomalowałyście, pomalowałyśmy, pomalowana, pomalowane, pomalowani, pomalowanie, pomalowan, pomalowany, pomalowawszy, pomaluj, pomalują, pomalujcie, pomalujecie, pomalujemy, pomalujesz, pomaluję, pomalujmy},

{robi, robią, robiąc, robiąca, robiące, robiący, robicie, robić, robieni, robienie, robię, robili, robiliby, robilibyście, robilibyśmy, robiliście, robiliśmy, robil, robila, robilaby, robilabym, robilabyś, robilam, robilaś, robilby, robilbym, robilbyś, robilem, robileś, robilo, robiloby, robilobym, robilobyś, robilom, robilos, robily, robilby, robilbyście, robilbyśmy, robilbyście, robilbymy, robimy, robiona, robione, robiony, robisz, zrobi, zrobia, zrobicie, zrobić, zrobieni, zrobienie, zrobię, zrobili, zrobiliby, zrobilibyście, zrobilibyśmy, zrobiliście, zrobiliśmy, zrobil, zrobila, zrobilaby, zrobilabym, zrobilabyś, zrobilam, zrobilaś, zrobilby, zrobilbym, zrobilbyś, zrobilem, zrobileś, zrobilo, zrobiloby, zrobilobym, zrobilobyś, zrobilom, zrobilos, zrobily, zrobilby, zrobilbyście, zrobilbyśmy, zrobilbyście, zrobiliśmy, zrobimy, zrobiona, zrobione, zrobiony, zrobisz, zrobiszy},

{taniały, taniał, taniała, taniałaby, taniałabym, taniałabyś, taniałam, taniałaś, taniałby, taniałbym, taniałbyś, taniale, tanialem, taniałeś, taniało, taniałoby, taniałobym, taniałobyś, taniałom, taniałoś, taniały, taniałyby, taniałybyście, taniałybyśmy, taniałyście, taniałyśmy, taniano, tanieć, taniej, tanieją, taniejąc, taniejąca, taniejące, taniejący, taniejcie, tanieje, taniejecie, taniejemy, taniejesz, tanieję, taniejmy, tanieli, tanieliby, tanielibyście, tanielibyśmy, tanieliście, tanieliśmy, tanienie},

{umiał, umiała, umiałaby, umiałabym, umiałabyś, umiałam, umiałaś, umialby, umialbym, umialbyś, umiałem, umialeś, umiało, umialoby, umiąły, umiąłyby, umiąłybyście, umiąłybyśmy, umiąłyście, umiąłyśmy, umie, umiecie, umieć, umiej, umieją, umiejąc, umiejąca, umiejące, umiejacy, umiejcie, umiejmy, umieliby, umielibyście, umielibyśmy, umielisię, umieliszy, umiem, umiem, umiesz},

{uwierz, uwierzą, uwierzcie, uwierzenie, uwierzę, uwierzmy, uwierzono, uwierzy, uwierzycie, uwierzyć, uwierzyli, uwierzyliby, uwierzylibyście, uwierzylibyśmy, uwierzyliście, uwierzyliśmy, uwierzył, uwierzyła, uwierzylaby, uwierzylabym, uwierzylabys, uwierzylam, uwierzyłaś, uwierzyliby, uwierzylibym, uwierzylibyś, uwierzyliem, uwierzyłeś, uwierzyło, uwierzyłoby, uwierzyłobym, uwierzyłobyś, uwierzyłom, uwierzyłoś, uwierzyły, uwierzyłyby, uwierzyłybyście, uwierzyłybyśmy, uwierzyłyście, uwierzyłyśmy, uwierzymy, uwierysz, uwierzywszy, wierz, wierzq, wierzqc, wierzqca, wierzqe, wierzqcy, wierzcie, wierzenie, wierzę, wierzmy, wierzono, wierz, wierzcie, wierzyć, wierzli, wierzliby, wierzlibyście, wierzlibyśmy, wierzliście, wierzliśmy, wierzyl, wierzyla, wierzlaby, wierzlabym, wierzlabys, wierzlam, wierzlaś, wierzliby, wierzlibym, wierzlibyś, wierzlem, wierzleś, wierzyo, wierzloby, wierzlobym, wierzlobys, wierzlom, wierzyoś, wierzły, wierzłyby, wierzłybyście, wiezybyśmy, wierzłyście, wierzłyśmy, wierzmy, wierzysz},

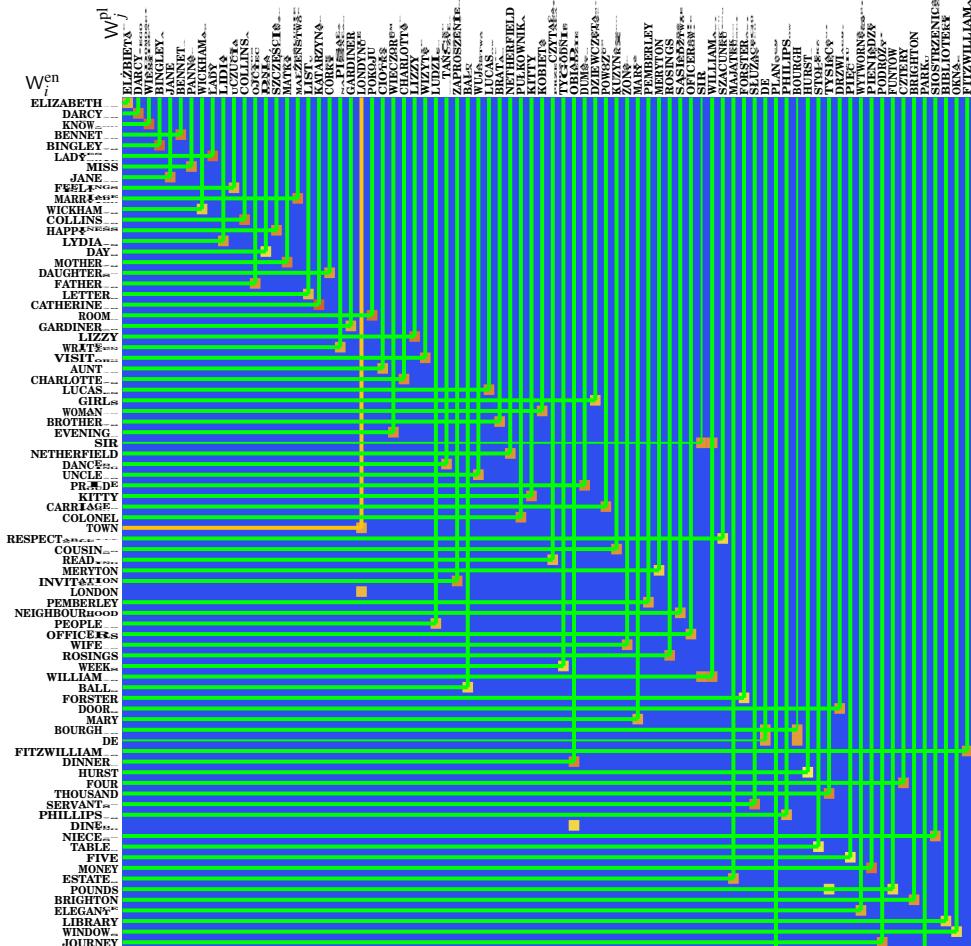
{widzą, widząc, widząca, widzące, widzący, widzenie, widzę, widzi, widział, widziała, widziałaby, widziałabym, widziałabyś, widziałam, widziałaś, widziałby, widziałbym, widziałbyś, widziałem, widzialeś, widziale, widzialoby, widzialobym, widzialobyś, widzialom, widzialeś, widzali, widzialo, widzialoby, widzialobym, widzialobyś, widzialom, widzialeś, widzali, widzialyby, widzialybyście, widzialybyśmy, widzialyśmy, widziana, widziane, widziani, widziano, widziany, widzicie, widzieć, widzieli, widzieliby, widzielibyście, widzielibyśmy, widzieliście, widzieliśmy, widzimi, widzisz, widź, widżcie, widżmy, zobacz, zobaczą, zobaczenie, zobaczę, zobacczmy, zobacczono, zobaczy, zobaczyście, zobaczyć, zobaczyli, zobaczyliby, zobaczylibyście, zobaczylibyśmy, zobaczyliście, zobaczyśmy, zobaczył, zobaczyła, zobaczyliby, zobaczylabym, zobaczylabyś, zobaczyłam, zobaczyłaś, zobaczyły, zobaczyłbym, zobaczyłbyś, zobaczyłem, zobaczyłeś, zobaczyło, zobaczyły, zobaczyły, zobaczyłyby, zobaczyłybyście, zobaczyłybyśmy, zobaczyłyście, zobaczyśmy, zobaczymy, zobacczysz, zobaczywszy},

{wiezie, wieziecie, wieziemy, wiezieni, wiezienie, wieziesz, wieziona, wiezione, wieziono, wiezony, wiozq, wiozqc, wioząca, wiozące, wiozący, wiozę, wiozla, wiozlaby, wiozlabym, wiozlabyś, wiozlam, wiozlaś, wiozlo, wiozloby, wiozły, wiozłyby, wiozłybycie, wiozłybyśmy, wiozłyście, wiozłyśmy, wiózl, wiózlb, wiózlbym, wiózlbys, wiózlem, wiózleś},

{zyska, zyskacie, zyskać, zyskaj, zyskajq, zyskajcie, zyskajmy, zyskali, zyskaliby, zyskalibyście, zyskalibyśmy, zyskaliście, zyskaliśmy, zyskał, zyskała, zyskały, zyskałabym, zyskałyby, zyskałam, zyskałaś, zyskałyby, zyskałbym, zyskałyś, zyskałem, zyskałoś, zyskały, zyskałyby, zyskałybyście, zyskałybyśmy, zyskałyście, zyskałyśmy, zyskam, zyskamy, zyskana, zyskane, zyskani, zyskanie, zyskano, zyskany, zyskas, zyskawszy, zyskiwać, zyskiwali, zyskiwaliby, zyskiwalibyście, zyskiwalibyś, zyskiwaliście, zyskiwaliśmy, zyskiwał, zyskiwała, zyskiwałaby, zyskiwałabym, zyskiwałabyś, zyskiwałam, zyskiwałaś, zyskiwałby, zyskiwałbym, zyskiwałbyś, zyskiwalem, zyskiwaleś, zyskiwalo, zyskiwałoby, zyskiwały, zyskiwałyby, zyskiwałybyście, zyskiwałybyśmy, zyskiwałyście, zyskiwałyśmy, zyskiwana, zyskiwane, zyskiwani, zyskiwanie, zyskiwano, zyskiwany, zyskuj, zyskujq, zyskując, zyskującqa, zyskującce, zyskujączy, zyskujcie, zyskuje, zyskujecie, zyskujemy, zyskujesz, zyskuję, zyskujmy}.

*Example 7.13.2.* In Fig. S10, we further apply the aforementioned word clustering algorithm to topic extraction and machine translation (see Table S1 for text sources).

ELŻBIETA MÓWIEĆ DARCY EGO PEWNE WIEZ BINGLEY MYSŁ  
 JANE SPRAW BENNET PANNE MŁODE NADZIEJE WIELKI SIOSTR CHCI WICKHAM  
 PANNE WAPNI ŁADY DOMU CAL ZOSTAĆ DAŁ LIDIA UCZUCIA CHWILA MILĘ DOBRE CZŁOWIEK  
 COLLINS OJCIEC SŁOŃ PRZYJACIŁ STARA CZAS WYZE SZCZEŚCIE MATKA PRZYSZŁO  
 WYRAZ MAŁZENSTW PRZYPUSZCZAM LIST WZGLED DROG PIĘK PODOBNE KONIEC SIESTA  
 PRAGNIE BEZPOMODŁ ODPARŁ OPOWIĘDZIŁ RADOŚĆ DŁUGA PRZYJEMNOŚĆ PANIA PYTAĆ ZACHOWAĆ  
 RODZINĘ USTĘP STWIERDZIŁ ZUPELNIE SZYBKO PRZYECHALA PIRAKAN ZYCIE POSTANOWIŁA DALEJ OBAWA  
 LEPKĘJ WRESZCIE MYSŁISZCZ ZDANY UWAGA WYDAJE POKOJU PIERWSZA MAŁA CIAGN MINIEST NASTĘPNEGO RÓWNIE SPOTKAĆ BOZIĘJ TRUDNO CIOŁA  
 ZAJĘTE ZDUMIENIE WIECZORE CHARLOTTA PODOBNA ZAMAJĄCA WODA LIZZY ROZSĄD WYJECHAL UPRZEJMOS DODAĆ RODZINY LONGBOURN SLYSM  
 WIZYTY LUDZI STOSUNKI ZAPROSZENIE DRUGI OSTATNIE BAL TROCHE WUJ LUCAS DZIAŁOŚĆ POTRAFIĆ RZEWY  
 WRÄZENIĘ POMYSŁY SPOKOJNY PIEKNY OGÓŁE OCHŁA BRATA OCZYWIŚCIE NETHERFIELD WIDOWIĘ ZŁOŻYĆ USMIECHEN  
 PRZYKRO ZAWALONA KOLPĘZKI SZCZEGÓLSKA UZNAŁ GŁEBOKA SZCZEGÓLSKA STAŁA KOBIEC ZWRÓCIA NOWSZA  
 HOZIOPRY SWIADKAMI RZECZ CZYTAĆ TYCHODZIĘSIĘ LIPĘ MINNYS OSWIADCZYŁ SIESTREZENIE PASTOR OBIAŁA  
 SKÓBLIZNA WYGŁADAŁA WŁASNOŚĆ GRACJA CZESŁAW ZWYKŁE DUMA DZIEWCZĘŁA WSTARZĘ RZOMIAWAĆ PROST OGRONIE  
 POWIĘZIĘ ZŁOŁĘTRZ YCZĘSŁAW KUZYNA ŹONY MARA RANKA PEMBERLEY NAŁĘŻY WYPADKĘ BRC KAZA WROCŁAW ZALEDWY ROZMOWY  
 WYBRZESZCZ POSTEP PRZESŁAŁ WŁASCIWY USTAWA PAR CIESZES WIERZYŁ SADZIE SPÓRZESZ KROTKA POWAŻNI POJECIE MILCZEŃ GOSCI GŁOSI DAWN POCZĘTNI ZYWA PRZYNAFEL MERTYON  
 ZAROWSKO STRON ROSINGS KŁĘBIĘ WYJAŚNIEĆ UCZCZEN WYDAWAŁ PODZIĘKOWA UDAN MOLAK ZAI LASKAWA GODZINĘ PRZYPADKĘ NIECHCE STANIE ZNAŁ WATPAVOSZ OPOWIĘSĆ OCZ  
 PRZEDSTAWIĆ OPONOWAŁ LUBIE SLUCHAĆ RESZTA WŁOŻTAWIAŁ SASHIBROWA CIEKAWOŚĆ HEDON WZROKOWO OFICER MOKSA SIR 920H92 DMIAN DUEDE  
 OTRYZMATEK POZCZAT WILLIAM LADNE PRZYZNAYŁ LATWO PRZERZECZŁ TŁUMACZYŁ OGROŃCZ MIESZAKAŁ SŁACZNIĘ  
 WSPOMNIĘĆ OPIS WŁOŻNIKI NIEPORÓW ZAŁAWA ŚJADŁA OPINIA MAJĄTEK ZDROWIE URODZIŁA BACZNA SZKATKA CHARTER WNBISKU USPOSOBIENIE ŚLUB HRASTWA FORSTER  
 SŁUŻĄCZKI WŁOŻNIKI MEZCZYZNA POMOC WŁOŻNIKI UWIERZIŁ WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI  
 POWIĘZIĘ BŁOGI WYJAZD ŚWIECZ DE WŁOŻNIKI SPACER POGODĘ SPŁEZIŁ ALEXANDR WŁOŻNIKI  
 DREWŁI TWARY WŁOŻNIKI DZWIAŁ PRZYCZYN KONIECZNOŚĆ ZEŁWY MYLKĘ WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI  
 GLUPIAKI GŁADMĘ WŁOŻNIKI UPŁYWAŁ WŁOŻNIKI ZŁOŻAŁA PRZYJAZD PÓBYTU PLEBANTY ZNA BOZ WŁOŻNIKI WŁOŻNIKI  
 PRÓŻNOŚĆ POMAG JAZU WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI  
 SALONI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI  
 KAROLINA DZWIAŁ PODRÓŻY MOTORYKI GŁOŚNIKI FUNKTOW DUCHY STRACIŁ PRZYWIĘZIŁA  
 TONI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI  
 WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI  
 ZIMĘ WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI  
 WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI  
 KUPIĘ ĆZARZAJĄCY PRZYPOMINA KILU ZŁE DORTANIA WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI  
 WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI WŁOŻNIKI  
 (a)



(b)

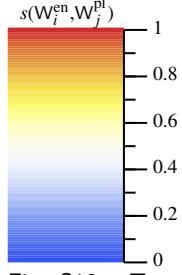


Fig. S10. Text mining in Polish. (a) Statistically identified topics ( $n_{ii} \geq 20$ ) in a Polish version of *Pride and Prejudice*, with the same color encoding scheme as Fig. S3. (b) Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{pl}})$  between selected topics in English and Polish versions of *Pride and Prejudice*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) ranking of the similarity score for the optimal node. Green (resp. amber) cross-hair indicates an exact (resp. a close but non-exact) match.



## 7.2 Modified Porter stemming algorithm for Russian

Our analysis mainly centers on Russian documents in the modern (post-1918) orthography.<sup>90</sup> We also assume that the Russian documents under investigation do not carry stress marks (which only appear in pedagogical contexts) on individual words. We remove the diacritical mark from ё in the effective spelling.

In modern Russian, the following present tense conjugations for the Russian verb “to be”

*еси, есмы, есмъ, есте, есть, суть*

are no longer used as their English counterparts, but they are found in works of 19th-century Russian authors (like Dostoevsky). The word *есть* is still used in modern constructions of possessions (such as *У меня есть ...* “I have ...”). We are going to include all the inflected forms of the Russian verb “to be” in our list of stop words.

**Definition 7.14** (Russian stop words). If a word belongs to the following list<sup>91</sup>:

*а, б, без, более, больше, будем, будет, будеме, будешь, будто, буду, будут, будучи, будь, будьте, бы, быв, быбав, быбавши, быбавший, быбавем, быбаает, быбаете, быбаешь, быбай, быбайте, бывал, бывали, бывало, быбать, быбаю, быбают, быбающий, быбая, быбшая, быбшего, быбшее, быбшей, быбшем, быбшему, быбшею, быбши, быбшие, быбший, быбшим, быбших, быбшиую, быбыл, быбыла, быбыли, быбыло, быбыть, в, вам, вами, вас, ваш, ваша, ваше, вашего, вашей, вашем, вашему, вашею, ваши, вашим, вашими, ваших, вашу, вдоль, вдруг, ведь, весь, весьма, взад, вместе, вместо, вне, внутри, во, вовне, вовсе, возле, вокруг, вон, во-преки, вот, впереди, впрочем, все, всегда, всего, всей, всем, всеми, всему, всех, всею, всё, вскоре, вслед, всю, вся, всяк, всяка, всякая, всяки, всякий, всяким, всяками, всяких, всяко, всякое, всякой, всяком, всяку, всяку, вы, где, да, давай, давайте, даже, дай, дайте, делав, делавша, делавшего, делавшее, делавшей, делавшем, делавшему, делавшему, делавшею, делавши, делавши, делавшой, делавшим, делавши-ми, делавших, делавшую, делаем, делаема, делаемо, делаемого, делаемое, делаемой, делаемом, делаемому, делаемою, делаемую, делаемы, делаемы, делаемый, делаемым, делаемыми, делаемых, дела-ет, делаете, делаешь, делай, делайте, делал, делала, делали, делало, делан, делана, деланная, деланного, деланное, деланной, деланном, деланному, деланную, деланные, деланный, деланным, делан-ными, деланных, делано, деланы, делать, делаю, делают, делающа, делающего, делающе, делающей, делающем, делающему, делающею, делающие, делающий, делающим, делающими, делающих, делающю, дела-я, для, до, должен, должна, должно, должны, другая, другие, другим, другими, других, другого, другое, другой, другом, другому, другою, другую, е, его, едва, ее, её, ей, еле, ему, если, если, есмы, есмъ, есте, есть, еще, ещё, ю, ж, же, за, затем, зато, зачем, здесь, и, ибо, из, изо, или, им, именно, ими, иначе, иная, иногда, иного, иное, иной, ином, иному, иною, иную, иные, иным, иными, иных, их, к, ка, каждая, каждого, каждое, каждой, каждом, каждому, каждою, каждую, каждые, каждый, каждым, каждыми, каждых, как, какая, какие, каким, какими, каких, каков, какова, каково, какового, каково, како-вой, каковом, каковому, каковою, каковую, каковы, каковые, каковым, каковыми, каковых, какого, какое, какой, каком, какому, какою, какую, кем, ко, когда, кого, кое, коею, кое, коею, кои, коим, коими, коих, кой, ком, кому, конечно, которая, которого, которое, которой, которому, кото-рою, которую, которые, который, которым, которых, кою, коя, кроме, кто, куда, ли, либо, лишь, любая, любого, любое, любой, любом, любому, любюю, любую, любые, любым, любыми, любых, м, между, мене, меня, мне, мног, многа, многая, многи, многие, многий, многим, многими, многих, много, много, много, многое, многой, многом, многому, многую, мной, мною, мог, моги, могите, могла, мог-ли, могло, могу, могут, могущая, могущего, могущее, могущей, могущем, могущему, могущею, могущие, могущий, могущим, могущими, могущих, могущую, могшая, могшего, могшее, могшей, могшем, могшему, могшею, могшие, могший, могшим, могшии, могших, могшиую, мое, моего, моей, моем, моему, моё, можем, может, можете, можешь, можно, мои, моим, моими, моих, мой, мою, моя, мы, на, навстречу, над, надо, назад, наиболее, наконец, нам, нами, напролет, напролёт, напротив, нас, насколько, настолько, наш, наша, наше, нашего, нашей, нашем, нашему, нашую, наши, нашим, нашими, наших, нашу, не, негде, него, нее, неё,ней, некем, некогда, некого, некому, некоторая, некоторого, некоторое, некоторой, некотором, некоторому, некоторую, некоторую, некоторую, некоторые, некоторый, некоторым, некоторыми, некоторых, некуда, нельзя, нем, немног, немнога, немногая, немноги, немногие, немногий, немногим, немногими, немно-гих, немного, немногого, немногое, немногой, немногом, немногому, немногую, немногую, нему, передко, нескольким, несколькими, нескольких, несколько, несмотря, нет, неужели, нечего, нечем, нечemu, нею, нём,*

<sup>90</sup>When 19th-century Russian literature is printed and read today, the spellings therein are thoroughly modernized.

<sup>91</sup>Our list of Russian stop words is based on <http://snowball.tartarus.org/algorithms/russian/stop.txt>, with extensive modifications to roughly match their counterparts in English. In particular, we have included all the inflected forms of *быть* “to be”, *каждый* “every”, *какой* “what kind”, *кто* “who”, *чей* “whose”, *что* “what”, and all the possessive pronouns declined in six cases and three genders. Some archaic, poetic or dialectal forms are also included.

ни, нибудь, нигде, никак, никакая, никакие, никаким, никакими, никаких, никакого, никакое, никакой, никаком, никакому, никакою, никакую, никем, никогда, никого, никому, никто, никуда, ним, ними, ниоткуда, нисколько, них, ничего, ничей, ничем, ничему, ничто, ничьего, ничьей, ничьему, ничьею, ничьё, ничьём, ничьи, ничьим, ничьими, ничых, ничью, ничья, но, ну, ныне, нэи, о, об, оба, обе, обеим, обеими, обеих, обо, обоего, обоим, обоими, обоих, один, одна, однако, одни, одним, одними, одних, одно, одного, одной, одном, одному, одною, одну, около, он, она, они, оно, опять, от, откуда, отнюдь, ото, отсюда, оттого, оттуда, отчего, очень, перед, по, под, подле, пожалуй, позади, пока, поныне, поскольку, после, посреди, потом, потому, почем, почему, почём, почти, поэтому, прежде, при, притом, про, против, пускай, пустъ, раз, разве, с, сам, сама, самая, сами, самим, самими, самих, само, самого, самое, самоё, самой, самом, самому, самою, саму, самую, самые, самый, самым, самыми, самых, свое, своего, своей, своем, своему, своею, своё, свои, своим, своих, свой, свою, своя, сделав, сделавшая, сделавшего, сделавшее, сделавшей, сделавшем, сделавшему, сделавшему, сделавши, сделавши, сделавший, сделавшим, сделавшими, сделавших, сделавшую, сделаем, сделает, сделаете, сделашь, сделай, сделайте, сделал, сделала, сделали, сделало, сделан, сделана, сделанная, сделанного, сделанное, сделанной, сделанном, сделанному, сделанною, сделанную, сделанные, сделанный, сделанным, сделанными, сделаных, сделано, сделаны, сделать, сделаю, сделаю, себе, себя, сего, сей, сейчас, сем, сему, сею, сё, сём, сие, сии, сим, сими, сих, сию, сия, сквозь, сколь, скольким, сколькими, скольких, сколько, скоро, слишком, словно, снова, со, собой, собою, совсем, согласно, спустя, среди, столь, стольким, столькими, стольких, столько, суть, сущая, сущего, сущее, сущей, сущем, сущему, сущему, сущие, сущий, сущим, сущими, сущих, сущую, сюда, та, так, такая, также, таки, такие, таким, такими, таких, таков, такова, таково, таковы, такого, такое, такой, таком, такому, такою, такую, там, твое, твоего, твоему, твою, твоё, твои, твоими, твоих, твой, твою, твоя, те, тебе, тебя, тем, теми, теперь, тех, то, тобой, тобою, тогда, того, тоже, той, только, том, тому, том, точно, тою, три, ту, туда, тут, ты, у, уж, уже, хоть, хотя, часто, чего, чей, чем, чему, через, чём, что, чтоб, чтобы, чуть, чье, чьего, чьей, чем, чьему, чьею, чьё, чьём, чьи, чым, чьими, чых, чью, чья, эи, эта, эти, этим, этими, этих, это, этого, этой, этом, этому, этот, эту, эти, я,

then we consider it a Russian stop word. All the Russian stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list.  $\square$

The following Russian verb has highly irregular conjugations:

дав, давав, дававшая, дававшего, дававшее, дававшей, дававшем, дававшему, дававшему, дававши, дававши, дававший, дававшим, дававшими, дававших, дававшую, даваем, даваема, даваемо, даваемо, даваемо, даваемое, даваемой, даваемом, даваемому, даваемою, даваемую, даваемы, даваемые, даваемый, даваемым, даваемыми, даваемых, давай, давайте, давал, давала, давали, давало, давать, давая, давшая, давшего, давшее, давшей, давшем, давшему, давшему, давши, давшие, давший, давшим, давшими, давших, давшую, дадим, дадите, дадут, даем, дает, даете, даешь, даём, даёт, даёте, даёшь, дай, дайте, дал, дала, дали, дало, дам, дан, дана, данная, данного, данное, данной, данном, данному, данною, данную, данные, данный, данным, данными, данных, дано, даны, даст, дать, даши, даю, дают, дающая, дающег, дающе, дающей, дающем, дающему, дающею, дающие, дающий, дающим, дающими, дающих, дающую — “give”.

We are not going to include these conjugated forms as stop words, so as to be consistent with other languages under our consideration. Nevertheless, we will define a string pattern **giveRussian** using the list above, to facilitate the clustering of Russian content words.

For sorting purposes (*Mathematica* v11.0 does not seem to support multi-level sorting of Russian words), we will need to transliterate the Russian effective spelling and essential root into the Latin script. The transliteration is not meant to be phonetically accurate, but only serves as a convenient one-to-one mapping. The “alphabetic order” in our clustering algorithm for Russian words will refer to the Latin transliterations, rather than the native order according to the Russian alphabet.

**Algorithm 7.15** (Transliteration of Russian Letters). *A text string  $\hat{\sigma}$  derived from a Russian word is transliterated into Latin-based script through the following replacements:*

а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ѣ	ѧ
a	b	v	g	d	e	zh	z	i	y	k	l	m	n	o	r	s	t	u	f	h	c	ch	sh	th	y	ye	ya	

The result of such a transliteration is denoted by RuLat( $\hat{\sigma}$ ).

### 7.2.1 Effective spelling and essential root

It is assumed that all Russian words are converted to lowercase<sup>92</sup> before going through any of the procedures below.

**Definition 7.16** (Russian Effective Vowels and Verb Prefixes). Hereafter in §7.2, the symbol  $\mathbf{V}^*$  stands for any member from the list  $\{a, e, u, o, y, ы, ь, ю, я\}$ , the so-called Russian effective vowels.<sup>93</sup> In line with the multiplicity notations introduced in Definition 3.3, the symbol  $\mathbf{V}_m^*$  stands for a text string formed by consecutive appearance of one or more (not necessarily identical or distinct) members from the set of Russian vowel extensions.

Dual to the notations above, the symbol  $\mathbf{C}^*$  stands for any character that does not belong to the list  $\{a, e, u, o, y, ы, ь, ю, я\}$ , and  $\mathbf{C}_{m_0}^*$  stands for a text string formed by consecutive appearance of zero or more (not necessarily identical or distinct) characters that do not belong to the same list.

The symbol  $\mathbf{P} = (\emptyset|ө|өз|өө|өөз|өс|өы|ðо|ðа|ðи|ðи|на|о|оð|оðо|ом|омо|непе|но|под|подо|нри|раз|разо|рас|с|ко|у)$  represents the possible prefix of a Russian verb.

We further define  $\mathbf{P}^* = (\emptyset|ө|өз|өө|өөз|өс|өы|из|ис|иб|об|ом|омо|непе|но|под|подо|нри|раз|разо|рас|с|ко|у|П)$  to account for a possible insertion of the Russian hard sign  $\flat$  after a prefix ending in a consonant.  $\square$

**Algorithm 7.17** (Russian effective spelling). *For a Russian word  $\hat{\sigma}$ , its effective spelling EffSpell( $\hat{\sigma}$ ) is constructed in sequential steps:*

(1) If  $\hat{\sigma}$  contains  $(a|e|ë|u|o|y|ы|ь|э|ю|я)$ , then leave it as is; otherwise, replace it by  $\xi\hat{\sigma}\xi$ .

(2) Replace<sup>94</sup>

$(\emptyset no)ciуд~$ σύιτ	$(\emptyset y)c(ee еви ел яд)~$ σύιτ				$X^e(\emptyset no y)чт~$ Хчест			
$zрX^e(a e u o)\hat{\chi}~$ $\gamma\rho\hat{X}'$	жел~	завтрак~	зey(к ч)~	изум~	$испX^e(o ы)~$	каза~	камин~	кокет~
колен~	корач	коро(тк ток ч)	крам(к ок ч)	кри~	Xσξπ	кажущий	χεарθ	զօկտ
коλүен~	σηρт	σηρт	σηρт	κρρи	кро(m тк ток и)	мылδ		
кров(∅ e)л~	кроват~	крыл~	л(e ë)շ(∅ к ч)~		од(∅ и)հ~			
роωφ	ββεδδ	крыл	լայγτχ		I			
остан(а о)в~	отча~	раv(∅ e)н	разговар~	разум~	са(ди дя жу)~	солни~		
остаñов	отча	ραβн	говор	раçум	σыти	σοլε		
ств(ова у)Х	стен~	сторон~	m(e ë ь).м~	тарел~	յш(a e ek u κ и)Х~	խօ(մ ч)		
сть	вагл	σайδн	термк	тарел	յար	χօտ		
холл~	холм~	чис(∅ e)л	денни	сент	стара(∅ я)	յх(a e o om y)		
χօլլ	хօլμ	ξισλ	δεñնи	σանкт	старый	յար		

where  $\hat{\chi}'$  results from doing  $\epsilon \rightarrow \beta$ ,  $m \rightarrow \mu$ ,  $n \rightarrow \tilde{n}$ ,  $u \rightarrow \tau\sigma$ ,  $ч \rightarrow \xi$  on  $\hat{\chi}$ .

(3) Replace

$(\emptyset no)стара~$ τρηа	$(eже каждо)дневнХ~$ день				$γրան(у ч)~$ гранич	$б(a o)լմ$	внимա~
воени~	вол	дамV~	женат~	завтра(∅ ии)~	кров~	марտ~	милион
войн	βвол	δдаmV	μар	τμօրω	κρօբ	μրատσ	πէմօն
пал(е ь)ц	план~	погод	предани~	предел~	пր~	βιրծ	πիկք
φιñγ	πλաñ	ωεθ	πրատ	λից	πρ	γարձ	μար
соба(к ч)	специал	странн	стрем	сударын	сумм~	տոն(∅ o)(κ ч иши)~	
δօγ	спеξլ	φρемδ	ριչ	μադամ	σүմի	θиñ	
үбө(∅ ж)д	угод~	улыц	xва(m ч)~	частX <sup>e</sup> (e и и ь я)~	յոշХ~	լեdi	па
βиñк	εñօδ	τρит	χβատ	чаστХ	յօշ	լլեδδ	πաաσ
							ալաօտ

<sup>92</sup> As of v11.0, *Mathematica* does not produce desired case conversions of Russian letters upon invoking the *ToLowerCase* or *ToUpperCase* operation. Therefore, case conversions for all the Russian letters need to be hand-coded if one implements our algorithm in *Mathematica*.

<sup>93</sup> We note that although ə is indeed a vowel letter, it never appears in the ending of a Russian word. The soft sign ь was historically a short vowel, and it still occurs often in modern Russian word endings.

<sup>94</sup> Words like *кровать* “bed” and *кровли* “roofs” appear to carry commonly seen verb suffixes -атъ and -ли, but are actually not related to *кровъ* “blood”. We have created exceptional effective spellings to separate these words (and their inflections) into three classes. However, we note that our algorithm does not distinguish *кровъ* “shelter” (and its inflections) from any of the inflected forms of *кровъ*. Furthermore, *кровъ* “blood”, *кровоточить* “to bleed” and *кровотечение* “bleeding” are split into three different clusters by our current method.

## (4) Replace

$(\emptyset вз) волн~$	$(\emptyset вы)m(e \ddot{e})(\kappa \chi)(\emptyset и)~$	$(\emptyset но) жен^{X^{\infty}}(и ю я)~$	$(\emptyset о\delta с) нов\hat{x}~$	$(\emptyset но) счаст~$
$\omega\beta$	$тек$	$ма\beta$	$\tilde{нов}\hat{x}$	$счаст$
$(\emptyset на по про) чит~$	$(\emptyset по) чувств~$	$(в воз по) люб~$	$(по про) яв~$	$P^* \text{giveRussian}(\emptyset съ ся)~$
$\check{чит}$	$\varphiул$	$люб$	$яв$	$P^* \rho_гы\betaр$
$бени бo(к ч)$	$брак(o у)\hat{x}$	$вез велик~$	$вид втия~$	$выжима~$
$\betae\tilde{v}\tilde{v}$	$тио\beta$	$реф\sigma\hat{x}$	$бо\betaи$	$выжима~$
$дал(ек ёк ьн)$	$д\beta(a е у(м м х))$	$девуш(\emptyset е)к~$	$девят~$	$ди(к ч)~$
$\varphiар$	$2e2$	$дев$	$9я$	$ди\betaи$
$дум$	$жизн~$	$заб замуж$	$зл~$	$зна(л н ю)~$
$\deltaум$	$жсит$	$\zetaаб$	$\zetaл$	$и\betaол~$
$ко\tilde{n}\zeta$	$лид~$	$мил~$	$мир~$	$моля~$
$ко\tilde{n}\zeta$	$лид$	$ми\beta$	$ми\beta$	$мягк~$
$отве(m ч)~$	$отлив~$			$n(e(в ви л m) o(em em eши ём ёт ёши ю ют ющ))$
$\alpha\beta\sigma\omega$	$отлив$			$песн$
$парк~$	$перв~$	$нес(\emptyset е)н~$	$ней$	$пи\betaьменн~$
$парк$	$1e$	$песн$	$\pi\xi$	$пи\betaат$
$поз^{X^{\infty}}(б о)~$	$поим~$	$покл~$	$поко$	$по\betaч(e \ddot{e})(\emptyset с)m~$
$3X$	$поня$	$покл$	$ру\beta$	$по\beta\beta\rho$
$npV_m^*m$	$np^{X^{\infty}}(и о)ли\beta~$	$приним~$	$прия\betaел~$	$прол(e \ddot{e})~$
$\pi pV_m^*\tau$	$prXли\beta$	$приня$	$прия\betaел$	$ле$
$разлу(к ч)~$	$свет~$	$сл(y ы)и$	$смел~$	$собы~$
$\sigma\betaе$	$свет$	$лы\beta$	$смел$	$с\betaи~$
$треб$	$трем~$	$трон~$	$уи~$	$цел$
$тре\beta$	$3и$	$тrog$	$Qu$	$человек$
$брат(\emptyset а е ом у ьев ья ьям ьями ьях)$				$лет(\emptyset a(\emptyset m mi x))$
	$\beta\betaрат\beta$			$год$
$мы\beta(V^* \emptyset)$	$окон$	$пи\betaем$	$полях$	$сем(i ~b)$
$мыть$	$окно$	$пи\betaьмо$	$поля$	$7e$
				$сн(ax об у ы)$
				$сон$

## (5) Replace

$(no c) лож~$	$P^* е\betaу~$	$P^*(бу\beta бы)$	$P^*(и\beta и\betam)~$	$P^* ле(m ч)~$	$брак~$	$звол$	$знатн~$	$зна(к ч)~$
$клад$	$P^* exa$	$P^* \betae$	$P^* и\beta$	$ф\betaет$	$ма\beta$	$а\beta\omega$	$\tilde{нов}\hat{x}$	$му\beta$
$мер(\overline{i})~$	$плать~$	$пож~$	$пре~$	$прив~$	$прил~$	$сем(e\beta и)~$	$чет$	$P^* е\betaе(m m te и\beta)$
$мmer$	$бплать$	$ж$	$пре$	$\beta$	$л$	$фа\beta$	$кел$	$P^* exa$

## (6) Replace

$(до\betaч до\betaч(\emptyset а е)к)$		$(\emptyset наи)(вы\betaок вы\betaси)~$		$(\emptyset наи)(худи\beta ху\betaж)~$
$до\betaчер$		$вы\betaок$		$плох$
$(\emptyset наи) млади~$	$(\emptyset наи) низи~$	$(\emptyset наи по) лучи~$	$(\emptyset с) дел~$	$(mem m\beta\betat)(o уши уше)к~$
$ммолод$	$ни\beta$	$хороши$	$бдел$	$тетя$
$C_m^* осст~$	$бал~$	$бл~$	$глян$	$де(и я)~$
$C_m^* озм$	$ббал$	$жжбл$	$гля\beta$	$бгдея$
$e(z\beta x)~$	$жб$	$забет$	$игр~$	$ис\betaм$
$с\betaе$	$ж$	$за$	$ы\betaгр$	$исзз$
$наи\hat{x}_m mX$	$нест$	$обр~$	$на~$	$пол\hat{v}^{\infty}(a е и о ы ы ю я)~$
$\hat{x}_m ш$	$нес$	$ччобр$	$ппа$	$пол\hat{x}$

слав~	слов~	собств~	спор~	стол~	тз	юноши~	дасм	де	дорого	отцах
зслав	цслав	зопств	цспор	цстол	ци	юност	дал	дпде	дорог	отец
<u>nan(<math>\mathbf{V}^* \emptyset (\emptyset \ddot{u} mi x)</math>)</u>				пони	(мать мамоч  $\emptyset a e$   $\kappa$ ~  $\mathbf{mam}(\mathbf{V}^* \emptyset (\emptyset \ddot{u} mi x))$ )				матер	
отец		ппони								

## (7) Replace

$(\emptyset no c co)жсм~$	$\mathbf{C}_{m_0}^* \text{ен}~$	да~	$e(e \ddot{u})к$	зв	особ~	певши	пон~	риск~	сл~
жсал	$\mathbf{C}_{m_0}^* \text{гюн}$	чдача	ееч	зоб	ззоб	по	н	риотч	жсл
соб~	$\mathbf{X}_1(\emptyset вы no y)_M \mathbf{X}_2^\epsilon(o ы)~$	$\underline{\text{день}}$	$\underline{\text{лоб}}$	$\underline{\text{мягок}}$	$\underline{\text{ром}}$	$\underline{\text{училищ}}\mathbf{X}$	(живши ~живить)	~ести	~еу
зоб	$\mathbf{X}_1\text{мюж}\mathbf{X}_2$	дни	лба	мягч	рома	σкуλ	жил	ел	~ыца

and call the result  $\hat{\sigma}'$ .

- (8) Break down  $\hat{\sigma}' = \hat{\sigma}_1 \hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}'^{[\min\{\lambda(\hat{\sigma}'), \ell(\hat{\sigma}')\}]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where  $\lambda(\hat{\sigma}') - 1$  equals either the last position occupied by the string pattern  $(\mathbf{P}^*|\text{без}|\text{безд}|\text{бес}|\text{не}|\text{nna}|cmp)\hat{\chi}$  in the non-void  $\hat{\sigma}'$ , or  $-1$  if  $\hat{\sigma}' = \emptyset$ .

- (9) On  $\hat{\sigma}_2$ , perform the following substitutions in a sequel:

- (4.1) Do  $(ак|вич|лив|чек|чиш|чик|як|(мел|чик)\mathbf{X}) \rightarrow \emptyset$ ,  $(ар|ац|ник|ници|ств|яр) \rightarrow \emptyset$ ,  $\mathbf{X}^\epsilon(e|б)\beta \rightarrow \mathbf{X}$ ,  $(кац|(\emptyset|у)уп(о|е|у)) \rightarrow \emptyset$ ,  $\sim(сы|ся|ться) \rightarrow \emptyset$ ,  $\sim\text{мя} \rightarrow \text{мен}$ ,  $\sim\text{и}\mathbf{V}^*(e|u|\ddot{u}|m|x|ю|я) \rightarrow \text{и}\mathbf{u}$ .
- (4.2) Do  $\sim((a|о|я)m|(a|и|ы|я)x|\mathbf{V}^*е\mathbf{в}|(e|о)го|(e|о)м|у|\ddot{u}|(\ddot{u})_{m_0}me|m|о\mathbf{в}|\mathbf{V}^*m(\emptyset|б)|и\mathbf{в}) \rightarrow \emptyset$ ,  $(ac|a|uc|oc)m \rightarrow \emptyset$ ,  $\sim(eу|к) \rightarrow \emptyset$ ,  $(κ|у)\mathbf{V}^*(\emptyset|x) \rightarrow \mathbf{V}^*$ ,  $\sim(\mathbf{V}^*m|(\ddot{u})_{m_0}me) \rightarrow \emptyset$ .

The result after these two steps of operations is called  $\hat{\sigma}'_2$ .

- (10) Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ .

- (11) Do  $(\ddot{e}|\ddot{E}) \rightarrow e$ ,  $\text{ммальч}\mathbf{V}^*\mathbf{X}\sim \rightarrow \text{ммальч}$ .

**Definition 7.18** (Russian protected range). Let  $\hat{\sigma}$  be a text string derived from a Russian word, its protected range  $\text{ProtRg}(\hat{\sigma}) = \min\{\lambda_1(\hat{\sigma}), \lambda_2(\hat{\sigma})\}$  is determined by two non-negative integers  $\lambda_1(\hat{\sigma})$  and  $\lambda_2(\hat{\sigma})$  specified through the following procedures:

- Look for the string pattern  $(вы|на|у|у)_{m_0} \mathbf{C}_{m_0}^* \mathbf{V}^*(\ddot{u})_{m_0} \sim$  in the string  $\hat{\sigma}$ ;
- If the string pattern above is found, the last position occupied by such a string defines  $\lambda_1(\hat{\sigma})$ ; otherwise, set  $\lambda_1(\hat{\sigma}) = 0$ ;
- Look for the pattern  $(\emptyset|вз|вы|на|но|нро|пар|пач|с) \mathbf{V}_{m_0}^* \mathbf{C}^* \mathbf{V}_{m_0}^* \mathbf{C}^* \sim$  in the string  $\hat{\sigma}$ ;
- If the pattern above is found, the last position occupied by such a pattern defines  $\lambda_2(\hat{\sigma})$ ; otherwise, set  $\lambda_2(\hat{\sigma}) = \ell(\hat{\sigma})$ .  $\square$

**Algorithm 7.19** (Russian essential root). Let  $\hat{\sigma}$  be the effective spelling of a Russian word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:

- (1) Break down  $\hat{\sigma} = \hat{\sigma}_1 \hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}^{[\text{ProtRg}(\hat{\sigma})]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma})$  is equal to the protected range of  $\hat{\sigma}$ .

- (2) Do  $б \rightarrow e$  on  $\hat{\sigma}_1$ , and call the result  $\hat{\sigma}'_1$ .

- (3) On  $\hat{\sigma}_2$ , perform the following substitutions in a sequel:

- (3.1) Do  $\sim\mathbf{V}_m^* \rightarrow \emptyset$ ;
- (3.2) Do  $\sim\mathbf{V}_{m_0}^*(\text{и}\mathbf{u}|\ddot{u}(\emptyset|и|у)|м|нн)_m \rightarrow \mathbf{V}_{m_0}^*$ ;
- (3.3) Do  $\sim\mathbf{V}_m^* \rightarrow \emptyset$ .

The result after these three steps of operations is called  $\hat{\sigma}'_2$ .

- (4) Concatenate  $\hat{\sigma}'_1$  and  $\hat{\sigma}'_2$ . If  $\ell(\hat{\sigma}'_1 \hat{\sigma}'_2) = 3$ , do  $\sim(a|e|u|\ddot{u}|o|y|ы|б|и|я) \rightarrow \emptyset$  on  $\hat{\sigma}'_1 \hat{\sigma}'_2$ , and call the result  $\hat{\sigma}'$ ; otherwise, define  $\hat{\sigma}' = \hat{\sigma}'_1 \hat{\sigma}'_2$ .

- (5) Do  $\sim(\ddot{u}|б)\hat{\chi}(\mathbf{V}^*е)_{m_0} \rightarrow e\hat{\chi}$  on  $\hat{\sigma}'$ .

- (6) Replace

<u>дн~</u>	<u>((Ø no c co)жса(Ø в m) сж)</u>	<u>(реb ребен ребяят)</u>	<u>X<sup>ε</sup>(Ø в л n u)e</u>
<u>день</u>	<u>жсал</u>	<u>дем</u>	<u>Xu</u>
<u>краси</u>	<u>усл</u>	<u>чдачад</u>	<u>~жеb</u>
<u>красив</u>	<u>услти</u>	<u>чдач</u>	<u>еед</u>
			<u>~жисв</u>
			<u>жисл</u>
			<u>~ne (в л m)</u>
			<u>но</u>

To improve the performance of verb clustering in Russian, we need the following extension to Algorithm 7.19.

**Algorithm 7.20** (Russian verbal essential root). *Let  $\hat{\sigma}$  be the token string associated with a Russian word, then its corresponding verbal essential root  $VbEssRoot(\hat{\sigma})$  is constructed in the following steps:*

- (1) Break down  $\hat{\sigma} = \hat{\sigma}_1\hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}^{[\text{ProtRg}(\hat{\sigma})]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma})$  is equal to the protected range of  $\hat{\sigma}$ .
- (2) Do  $\sim(\mathbf{V}_{m_0}(e|l|H|\mathbf{V}_m m)_{m_0})_m \rightarrow \emptyset$  on  $\hat{\sigma}_2$ , and call the result  $\hat{\sigma}'_2$ .
- (3) Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ .
- (4) Do приод → прииод.
- (5) Do  $\mathbf{P}^*(e\bar{o}|e\bar{z}\bar{ж}) \sim \mathbf{P}^* ex$ , (бp|бз|боз)(Ø|еm|ьm) → бep, P<sup>\*</sup>(u\bar{d}|хoж|u|иueди) → P<sup>\*</sup>xod.

### 7.2.2 Admissible mutation and approximate clustering

In Russian (as well as Polish and many other Slavic languages), fleeting vowels may appear or disappear in inflected forms of the same word. We need to heuristically detect and remove these fleeting vowels to achieve better clustering results. Meanwhile, in Russian (as well as Polish and many other Slavic languages), one encounters vowel alternations in verb conjugations, as in Spanish or the Germanic languages.

**Algorithm 7.21** (Russian vowel blotting). *For a token string  $\hat{\sigma}$ , its blotted form  $\text{BlotFltV}(\hat{\sigma})$  is constructed as follows:*

- (1) Do  $\sim \mathbf{C}^*(e|o)\hat{\chi}^e(\kappa|H|l|u)((\mathbf{V}_{m_0}^*(e|l|H|m)_{m_0})_m|u) \rightarrow \mathbf{C}^*\hat{\chi}$ .
- (2) Do  $\sim \hat{\chi}^e(e|m)ep((\mathbf{V}_{m_0}^*(e|l|H|m)_{m_0})_m|u) \rightarrow \hat{\chi}p$ .
- (3) Do  $e \rightarrow u$ ,  $o \rightarrow a$ .

In what follows, we will construct a bivariate Boolean-valued function  $HrdTest(\hat{\alpha}, \hat{\beta})$  on a “simple heredity test function” in Algorithm 7.22, and a set of “admissible suffix mismatch” rules in Algorithm 7.23.

**Algorithm 7.22** (Simple heredity test). *Let  $\hat{\alpha}'$  be a string obtained from  $\hat{\alpha}$  by the following steps:*

- (1) Break down  $\hat{\alpha} = \hat{\alpha}_1\hat{\alpha}_2$  into the concatenation of two strings  $\hat{\alpha}_1 = \hat{\alpha}^{[\min\{3, \ell(\hat{\alpha})\}]}$  (see the notation in Definition 3.1) and  $\hat{\alpha}_2$ , where the length of the first string  $\ell(\hat{\alpha}_1) = \min\{3, \ell(\hat{\alpha})\}$  is equal to 3 or the length of  $\hat{\alpha}$ , whichever is shorter.
- (2) Do  $\sim((\mathbf{V}_{m_0}^*(e|l|H|m|u)_{m_0})_m|u)$  on  $\hat{\alpha}_2$ , and call the result  $\hat{\alpha}'_2$ .
- (3) Construct  $\hat{\alpha}' = \hat{\alpha}_1\hat{\alpha}'_2$ .

Define  $\hat{\beta}'$  similarly. The Boolean-valued function  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if the lowercase form of  $\hat{\alpha}$  contains at least one instance of {a, e, u, o, y, ы, б, ю, я, лб, pm, ch} **AND** at least one of the following five conditions holds:<sup>95</sup>

- (i)  $\hat{\alpha} = \hat{\beta}$ ;
- (ii)  $\hat{\beta} = \hat{\alpha}uh$ ;
- (iii)  $\hat{\beta} = \hat{\alpha}u$ ;
- (iv)  $\min\{\ell(\hat{\alpha}'), \ell(\hat{\beta}')\} \geq 3$  **AND** ( $\hat{\alpha} = \hat{\beta}'$  **OR**  $\hat{\alpha}' = \hat{\beta}$  **OR**  $\hat{\alpha}' = \hat{\beta}'$ ).
- (v)  $\ell(\hat{\beta}) > \ell(\hat{\alpha}) \geq \frac{\ell(\hat{\beta})}{3}$  **AND**  $\hat{\beta} = \hat{\alpha}(a|e|o|я)_m(e|m|m)$ .

In what follows, we define  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{RootNW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{NW}^*(\hat{\alpha}, \hat{\beta})$  and  $\text{SuffixSW}(\hat{\alpha}, \hat{\beta})$ ,  $\text{SuffixSW}^*(\hat{\alpha}, \hat{\beta})$ ,  $\text{SW}^*(\hat{\alpha}, \hat{\beta})$  as what is done in the Danish case (Algorithm 5.7), which was also applied to two other Germanic languages treated in §5. Particular to the Russian case, we further define  $\text{SuffixNW}'(\hat{\alpha}, \hat{\beta})$  by performing the substitutions  $\sim((\mathbf{V}_{m_0}^*(e|l|H|\mathbf{V}_m^* m)_{m_0})_m|u) \rightarrow \emptyset$  to both components of  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta})$ . The notation  $\text{SuffixSW}'(\hat{\alpha}, \hat{\beta})$  is defined similarly. If  $\text{SuffixNW}'(\hat{\alpha}, \hat{\beta}) = [\hat{\tau}_1, \hat{\tau}_2]$ , then  $\text{SuffixNW}'(\hat{\beta}, \hat{\alpha}) = [\hat{\tau}_2, \hat{\tau}_1]$ .

<sup>95</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical **OR**.

**Algorithm 7.23** (Admissible suffix mismatch and vowel alternation). *For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function*

$$\text{AdmMut}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta}))$$

*returns TRUE if  $\text{NW}(\hat{\alpha}, \hat{\beta}) = \emptyset$  or  $[\emptyset, u][[a, o]][[e, u]][[u, \emptyset]]$  AND ( $\text{RootNW}(\hat{\alpha}, \hat{\beta})$  contains at least one instance of {a, e, u, o, y, ы, ь, ю, я} OR  $\Omega(\text{RootNW}(\hat{\alpha}, \hat{\beta})) = (\kappa|h|l|y|p))$  AND at least one of the following three conditions holds:*

- (i)  $\text{SuffixNW}'(\hat{\alpha}, \hat{\beta}) = [\varepsilon, \text{ж}][[\varepsilon, \text{з}][[\partial, \text{ж}][[\text{ж}, \text{з}][[\text{ж}, \text{у}][[\kappa, \text{у}][[\text{с}, \text{и}][[\text{м}, \text{и}][[\text{x}, \text{и}][[\text{е}, \text{и}][[\text{o}, \text{и}][[\text{с}, \text{и}][[\text{м}, \text{и}][[\text{д}, \text{и}][[\text{с}, \text{и}][[\text{м}, \text{и}];$
- (ii)  $\text{SuffixNW}'(\hat{\beta}, \hat{\alpha}) = [\varepsilon, \text{ж}][[\varepsilon, \text{з}][[\partial, \text{ж}][[\text{ж}, \text{з}][[\text{ж}, \text{у}][[\kappa, \text{у}][[\text{с}, \text{и}][[\text{м}, \text{и}][[\text{x}, \text{и}][[\text{е}, \text{и}][[\text{o}, \text{и}][[\text{с}, \text{и}][[\text{м}, \text{и}][[\text{д}, \text{и}][[\text{с}, \text{и}][[\text{м}, \text{и}];$
- (iii)  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = [\emptyset | (\mathbf{V}_m(\emptyset | \text{и} | \text{м} | \text{и} | \text{у})), \emptyset | (\mathbf{V}_m(\emptyset | \text{и} | \text{м} | \text{и} | \text{у}))].$

Similarly, one can evaluate another Boolean-valued function

$$\text{AdmMut}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta}))$$

by trading all the occurrences of NW in the statements above with SW.

**Algorithm 7.24** (Heredity test function). *The structure of the Russian heredity test function HrdTest( $\hat{\alpha}, \hat{\beta}$ ) is identical to the Polish version (Algorithm 7.11), except that the functions SimpHrdTest, RootNW, SuffixNW, NW, RootSW, SuffixSW, SW and the strings  $\hat{\alpha}', \hat{\beta}'$  must follow the Russian rules stated above.*

**Algorithm 7.25** (Russian verb aspect test). *Let NW( $\hat{\alpha}, \hat{\beta}$ ) be the result of performing Needleman–Wunsch alignment on strings  $\hat{\alpha}$  and  $\hat{\beta}$ . Let  $\mathbf{X}^*$  be an arbitrary non-empty string. The Boolean-valued function VbAspTest( $\hat{\alpha}, \hat{\beta}$ ) returns TRUE if at least one of the following four conditions holds:*

- (i)  $\text{NW}(\hat{\alpha}, \hat{\beta}) = \mathbf{X}^*$  (in other words,  $\hat{\alpha} = \hat{\beta}$ );
- (ii)  $\text{NW}(\hat{\alpha}, \hat{\beta}) = [\mathbf{P}^*, \mathbf{P}^*]\mathbf{X}^*$  (in other words,  $\hat{\alpha}$  agrees with  $\hat{\beta}$  up to a pair of word initial mismatching strings, both of which are possible prefixes of Russian verbs);
- (iii)  $\text{NW}(\hat{\alpha}, \hat{\beta}) = [\mathbf{P}^*, \mathbf{P}^*]\mathbf{X}^*([\emptyset, u][[a, o]][[e, u]][[u, \emptyset]][[u, e]][[o, a]]\mathbf{X}^*$ , where the two occurrences of  $\mathbf{X}^*$  may or may not represent the same non-empty string);
- (iv)  $\text{NW}(\hat{\alpha}, \hat{\beta}) = [\mathbf{P}^*, \mathbf{P}^*]\mathbf{X}^*[\hat{\tau}_1, \hat{\tau}_2]$  so that either  $[\hat{\tau}_1', \hat{\tau}_2']$  or  $[\hat{\tau}_2', \hat{\tau}_1']$  matches  $[\varepsilon, \text{ж}][[\varepsilon, \text{з}][[\partial, \text{ж}][[\text{ж}, \text{з}][[\text{ж}, \text{у}][[\kappa, \text{у}][[\text{с}, \text{и}][[\text{м}, \text{и}][[\text{x}, \text{и}][[\text{е}, \text{и}][[\text{o}, \text{и}][[\text{с}, \text{и}][[\text{м}, \text{и}][[\text{д}, \text{и}][[\text{с}, \text{и}][[\text{м}, \text{и}][[\text{у}, \text{и}][[\text{е}, \text{и}][[\text{o}, \text{и}][[\text{с}, \text{и}][[\text{м}, \text{и}][[\text{д}, \text{и}][[\text{с}, \text{и}][[\text{м}, \text{и}];$  Here,  $\hat{\tau}_1'$  results from doing  $\sim((\mathbf{V}_{m_0}^*(\text{и} | \text{и} | \mathbf{V}_{m_0}^* \text{м})_{m_0} | \text{и})) \rightarrow \emptyset$  on  $\hat{\tau}_1$ , and  $\hat{\tau}_2'$  is similarly defined.

Note that Smith–Waterman alignment is not used in this test.

The following approximate clustering algorithm for Russian words are very similar to the Polish counterpart (Algorithm 7.13), except that Latin transliteration (Algorithm 7.15) is used before sorting the tokens. For clarity, we state the algorithm in full.

**Algorithm 7.26** (Approximate clustering of Russian words). *The approximate clustering of a list of English words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  is completed in three stages:*

- (1) We sort the list  $\{(\hat{\alpha}_1, \text{EffSpell}(\hat{\alpha}_1)), \dots, (\hat{\alpha}_N, \text{EffSpell}(\hat{\alpha}_N))\}$  alphabetically according to the Latin transliteration of second component (Algorithm 7.15 applied to the effective spelling) of each entry. If two consecutive neighbors in the alphabetized list  $\{(\hat{\alpha}_{(1)}, \text{EffSpell}(\hat{\alpha}_{(1)})), \dots, (\hat{\alpha}_{(N)}, \text{EffSpell}(\hat{\alpha}_{(N)}))\}$  satisfies  $\text{HrdTest}(\text{EffSpell}(\hat{\alpha}_{(n)}), \text{EffSpell}(\hat{\alpha}_{(n+1)})) = \text{FALSE}$ , where  $n \in \mathbb{Z} \cap [1, N]$ , then add a demarcation line between these two entries. In this way, the alphabetized list is divided into separate groups of words tagged with their effective spellings:  $\{g_1 = \{(\hat{\alpha}_{(1,1)}, \text{EffSpell}(\hat{\alpha}_{(1,1)})), \dots, (\hat{\alpha}_{(1,n_1)}, \text{EffSpell}(\hat{\alpha}_{(1,n_1)}))\}, \dots, g_M = \{(\hat{\alpha}_{(M,1)}, \text{EffSpell}(\hat{\alpha}_{(M,1)})), \dots, (\hat{\alpha}_{(M,n_M)}, \text{EffSpell}(\hat{\alpha}_{(M,n_M)}))\}\}$ , where each sublist also preserves alphabetic order. In particular, we have  $\hat{\alpha}_{(1,1)} = \hat{\alpha}_{(1)}$  and  $\hat{\alpha}_{(M,n_M)} = \hat{\alpha}_{(N)}$ .
- (2) For each group of words  $g_m = \{(\hat{\alpha}_{(m,1)}, \text{EffSpell}(\hat{\alpha}_{(m,1)})), \dots, (\hat{\alpha}_{(m,n_m)}, \text{EffSpell}(\hat{\alpha}_{(m,n_m)}))\}$  where  $m \in \mathbb{Z} \cap [1, M]$ , we augment it into a tagged entry  $G_m = (g_m, \text{EffSpell}(\hat{\alpha}_{(m,1)}), \text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})), \text{BlotFltV}(\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))))$ . The list  $\{G_1, \dots, G_M\}$  is sorted alphabetically, with respect to  $\text{RuLat}(\text{BlotFltV}(\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))))$  (with highest priority),  $\text{RuLat}(\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})))$  (with medium priority), and  $\text{RuLat}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))$  (with lowest priority). If two consecutive neighbors in the alphabetized list  $\{G_{(1)} = (g_{(1)}, \hat{\gamma}'_{(1)}, \hat{\gamma}''_{(1)}, \hat{\gamma}'''_{(1)}), \dots, G_{(M)} = (g_{(M)}, \hat{\gamma}'_{(M)}, \hat{\gamma}''_{(M)}, \hat{\gamma}'''_{(M)})\}$  satisfy

$$\text{SimpHrdTest}(\hat{\gamma}'''_{(m+1)}, \hat{\gamma}'''_{(m)}) = \text{FALSE}$$

**AND**

$$\text{SimpHrdTest}(\hat{\gamma}_{(m)}''', \hat{\gamma}_{(m+1)}''') = \text{FALSE}$$

**AND**

$$\text{HrdTest}(\hat{\gamma}_{(m+1)}'', \hat{\gamma}_{(m)}'') = \text{FALSE}$$

**AND**

$$\text{HrdTest}(\hat{\gamma}_{(m)}'', \hat{\gamma}_{(m+1)}'') = \text{FALSE}$$

**AND**

$$\text{HrdTest}(\hat{\gamma}_{(m)}', \hat{\gamma}_{(m+1)}') = \text{FALSE}$$

where  $m \in \mathbb{Z} \cap [1, M]$ , then add a demarcation line between these two entries. In this manner, the alphabetized list  $\{G_{(1)}, \dots, G_{(M)}\}$  is divided into separate groups of tagged clusters  $\{\Gamma_1 = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, \Gamma_K = \{G_{(K,1)}, \dots, G_{(K,m_K)}\}\}$ . Generate a list of word clusters  $\{\check{\Gamma}_1 = (\check{G}_{(1,1)}, \dots), \dots, \check{\Gamma}_K = (\check{G}_{(K,1)}, \dots)\}$  by discarding all the tags (effective spellings, essential roots, vowel blotted forms) from each  $\Gamma_k$ , where  $k \in \mathbb{Z} \cap [1, K]$ .

- (3) For each word cluster  $\check{\Gamma}_k = (\check{G}_{(k,1)}, \dots)$ , we augment it into a tagged entry

$$\mathcal{G}_k \equiv (\check{\Gamma}_k, \check{\Gamma}'_k, \check{\Gamma}''_k) := (\check{\Gamma}_k, \text{EssRoot}(\text{EffSpell}(\check{G}_{(k,1)})), \text{VbEssRoot}(\text{EssRoot}(\text{EffSpell}(\check{G}_{(k,1)})))).$$

Recall that  $\hat{\sigma}^{-1}$  is the reverse of the string  $\hat{\sigma}$  (Definition 3.1). The list  $\{\mathcal{G}_1 = (\check{\Gamma}_1, \check{\Gamma}'_1, \check{\Gamma}''_1), \dots, \mathcal{G}_K = (\check{\Gamma}_K, \check{\Gamma}'_K, \check{\Gamma}''_K)\}$  is sorted sorted alphabetically, with respect to  $(\text{RuLat}(\check{\Gamma}''_k))^{-1}$  (with higher priority) and  $\text{RuLat}(\check{\Gamma}'_k)$  (with lower priority). If two consecutive neighbors in the alphabetized list  $\{\mathcal{G}_{(1)} = (\check{\Gamma}_{(1)}, \check{\Gamma}'_{(1)}, \check{\Gamma}''_{(1)}), \dots, \mathcal{G}_{(K)} = (\check{\Gamma}_{(K)}, \check{\Gamma}'_{(K)}, \check{\Gamma}''_{(K)})\}$  satisfy

$$\text{VbAspTest}(\hat{\Gamma}''_{(k)}, \hat{\Gamma}''_{(k+1)}) = \text{FALSE}$$

**AND**

$$\text{VbAspTest}(\hat{\Gamma}'_{(k)}, \hat{\Gamma}'_{(k+1)}) = \text{FALSE}$$

where  $k \in \mathbb{Z} \cap [1, K]$ , then add a demarcation line between these two entries. In this manner, the alphabetized list  $\{\mathcal{G}_{(1)}, \dots, \mathcal{G}_{(K)}\}$  is divided into separate groups of tagged clusters  $\{\tilde{\Gamma}_1 = \{\mathcal{G}_{(1,1)}, \dots, \mathcal{G}_{(1,n_1)}\}, \dots, \tilde{\Gamma}_J = \{\mathcal{G}_{(J,1)}, \dots, \mathcal{G}_{(J,n_J)}\}\}$ , from which the output list of word clusters is generated, after removal of all the tags from each  $\tilde{\Gamma}_j$ , where  $j \in \mathbb{Z} \cap [1, J]$ .

*Example 7.26.1.* Most Russian nouns distinguish six cases (nominative, genitive, dative, accusative, instrumental, prepositional) in singular and plural. A limited number of nouns have additional forms in three more cases (vocative, locative and partitive).

The sample nouns below are chosen according to the following web pages:

[https://en.wiktionary.org/wiki/Appendix:Russian\\_nouns#Declension\\_paradigms](https://en.wiktionary.org/wiki/Appendix:Russian_nouns#Declension_paradigms)

[https://en.wiktionary.org/wiki/Template:ru-noun-table#Basic\\_examples](https://en.wiktionary.org/wiki/Template:ru-noun-table#Basic_examples)

апреле, апрелей, апрелем, апрели, апрель, апрело, апрея, апрелям, апрелями, апрелях — “April”;

блюд, блюда, блюдам, блюдами, блюдах, блюде, блюдо, блюдом, блюду — “dish”;

бог, бога, богам, богами, богах, боге, боги, богов, богом, богу, боже — “god”;

бояр, боярам, боярами, боярах, бояре, боярин, боярина, боярине, боярином, боярину — “boyar”;

вид, вида, видам, видами, видах, виде, видов, видом, виду, виды — “appearance”;

волчиц, волчица, волчицам, волчицами, волчицах, волчице, волчицей, волчицем, волчицы, волчицу — “maturewolf”;

воротец, воротца, воротцам, воротцами, воротцах — “wicket”;

времена, временам, временами, временах, временем, времени, времён, время — “time”;

газет, газета, газетам, газетами, газетах, газете, газетой, газетою, газету, газеты — “newspaper”;

герое, героев, героем, герои, герой, герою, героя, героям, героями, героях — “hero”;

гетеросексуальности, гетеросексуальность, гетеросексуальностью — “heterosexuality”;

голов, голова, головам, головами, головах, голове, головой, головою, голову, головы — “head”;  
дверей, двери, дверь, дверьми, дверью, дверям, дверями, дверях — “door”;  
дев, дева, девам, девами, девах, деве, дево, девой, девою, деву, девы — “maiden”;  
день, дне, дней, днём, дни, дню, дня, дням, днями, днях — “day”;  
детей, дети, детьми, детям, детях, ребёнка, ребёнке, ребёнком, ребёнку, ребёнок, ребят, ребята, ребятам, ребятами, ребятах — “child”;  
друг, друга, друге, другом, другу, другже, друзей, друзья, друзьям, друзьями, друзьях — “friend”;  
железа, железам, железами, железах, железе, железой, железою, железу, железы, желёз — “gland”;  
жилищ, жилища, жилищам, жилищами, жилищах, жилище, жилищем, жилищу — “dwelling”;  
житие, житием, житии, житий, житию, жития, житиям, житиями, житиях — “life”;  
журнал, журнала, журналам, журналами, журналах, журнале, журналов, журналом, журналу, журналы — “magazine”;  
искр, искра, искрам, искрами, искрах, искре, искрой, искрою, искру, искры — “spark”;  
историей, историою, истории, историй, историю, история, историям, историями, историях — “history”;  
камень, каменев, каменя, каменьям, каменьями, каменьях, камне, камней, камнем, камни, камню, камня, камням, камнями, камнях — “stone”;  
карандаш, карандаша, карандашам, карандашами, карандашах, карандаше, карандашей, карандаши, карандашом, карандашу — “pencil”;  
книг, книга, книгам, книгами, книгах, книге, книги, книгой, книгою, книгу — “book”;  
колен, колена, коленам, коленами, коленах, колене, колено, коленом, колену — “tribe”;  
колена, колене, коленей, колени, колено, коленом, колену, коленям, коленями, коленях — “knee”;  
колена, колене, колено, коленом, колену, коленев, коленя, коленям, коленями, коленях — “joint”;  
копеек, копейка, копейкам, копейками, копейках, копейке, копейки, копейкой, копейкою, копейку — “kopek”;  
лёд, льда, льдам, льдами, льдах, льде, льдов, льдом, льду, льды — “ice”;  
лоскут, лоскута, лоскутам, лоскутами, лоскутах, лоскуте, лоскутов, лоскутом, лоскуту, лоскуты, лоскутьев, лоскутья, лоскутьям, лоскутьями, лоскутьях — “tag”;  
лошадей, лошади, лошадь, лошадьми, лошадью, лошадям, лошадями, лошадях — “horse”;  
мальчат, мальчата, мальчатам, мальчатами, мальчатах, мальчонка, мальчонкам, мальчонками, мальчонках, мальчонке, мальчонки, мальчонков, мальчонком, мальчонку, мальчонок — “littleboy”;  
море, морей, морем, морю, моря, морям, морями, морях — “sea”;  
мост, моста, мостам, мостами, мостах, мосте, мостов, мостом, мосту, мосты — “bridge”;  
муж, мужса, мужсе, мужей, мужсем, мужсу, мужся, мужсья, мужсьям, мужсьями, мужсьях — “husband”;  
неделе, неделей, неделею, недели, недель, неделю, неделя, неделем, неделеми, неделех — “week”;  
нож, ножса, ножам, ножами, ножах, ноже, ножей, ножси, ножом, ножу — “knife”;  
няне, няней, нянею, няни, нянь, няню, няня, няням, нянями, нянях — “nurse”;

озерца, озерце, озерцо, озерцом, озерцу, озёрец, озёрца, озёрцам, озёрцами, озёрцах — “lake”;  
окна, окнам, окнами, окнах, окне, окно, окном, окну, окон — “window”;  
ореол, ореола, ореолам, ореолами, ореолах, ореоле, ореолов, ореолом, ореолу, ореолы — “halo”;  
отец, отца, отцам, отцами, отцах, отце, отцов, отцом, отцу, отцы — “father”;  
палец, пальца, пальцам, пальцами, пальцах, пальце, пальцев, пальцем, пальцу, пальцы — “finger”;  
писем, письма, письмам, письмами, письмах, письме, письмо, письмом, письму — “letter”;  
площадей, площади, площадь, площадью, площадям, площадями, площадях — “square”;  
повод, повода, поводе, поводом, поводу, поводьев, поводья, поводям, поводьями, поводьях — “rein”;  
поле, полей, полем, полю, поля, полям, полями, полях — “field”;  
пух, пуха, пухе, пухом, пуху — “feather”;  
работ, работа, работам, работами, работах, работе, работой, работею, работу, работы — “task”;  
рот, рта, ртам, ртами, ртах, рте, ртов, ртом, рту, рты — “mouth”;  
сапожек, сапожка, сапожкам, сапожками, сапожках, сапожке, сапожки, сапожков, сапожком, сапожку,  
сапожок — “boot”;  
сестёр, сестра, сестре, сестрой, сестрою, сестру, сестры, сёстрам, сёстрами, сёстрах, сёстры — “sis-  
ter”;  
словаре, словарей, словарём, словари, словарь, словарю, словаря, словарям, словарями, словарях — “dic-  
tionary”;  
случае, случаев, случаем, случаи, случай, слушаю, слушая, слушаем, слушаями, слушаях — “event”;  
сна, снам, снами, снах, сне, снов, сном, сну, сны, сон — “dream”;  
снег, снега, снегам, снегами, снегах, снеге, снегов, снегом, снегу — “snow”;  
стол, стола, столам, столами, столах, столе, столов, столом, столу, столы — “table”;  
сын, сына, сынам, сынами, сынах, сыне, сынов, сыновей, сыновья, сыновьям, сыновьями, сыновьях, сыном,  
сыну, сыны — “son”;  
товарищ, товарища, товарищам, товарищами, товарищах, товарице, товарищей, товарищем, това-  
рищи, товаришу — “comrade”;  
урок, урока, урокам, уроками, уроках, уроке, уроки, уроков, уроком, уроку — “lesson”;  
учащегося, учащемсяся, учащемуся, учащиеся, учащийся, учащимися, учащимся, учащихся — “pupil”;  
учение, учением, учении, учений, учению, учения, учениям, учениями, учениях — “learning”;  
училищ, училища, училищам, училищами, училищах, училище, училищем, училищу — “college”;  
учителе, учителей, учителем, учители, учитель, учителю, учителя, учителям, учителями, учителях —  
“teacher”;  
феномен, феномена, феноменам, феноменами, феноменах, феномене, феноменов, феноменом, феномену,  
феномены — “phenomenon”;  
хвала, хвалам, хвалами, хвалях, хвале, хвалой, хвалю, хвалу, хвалы — “praise”;  
чай, чае, чаёв, чаи, чай, чаю, чая, чаям, чаями, чаях — “tea”.

It is worth noting that, unlike Polish, the consonants of the Russian noun stem mostly remain stable in all declined forms. The sample nouns given above are clustered by our algorithm as follows:

- {апреле, апрелей, апрелем, апрели, апрель, апрело, апреля, апрелям, апрелями, апрелях},
- {блюд, блюда, блюдам, блюдами, блюдах, блюде, блюдо, блюдом, блюду},
- {бог, бога, богам, богами, богах, боге, боги, богов, богом, богу, боже},
- {бояр, боярам, боярами, боярах, бояре, боярин, боярина, боярине, боярином, боярину},
- {вид, вида, видам, видами, видах, виде, видов, видом, виду, виды},
- {волчищ, волчища, волчищам, волчищами, волчицах, волчище, волчицей, волчищем, волчищи, волчицу},
- {воротец, воротца, воротцам, воротцами, воротцах},
- {времена, временам, временами, временах, временем, временем, времён, время},
- {газет, газета, газетам, газетами, газетах, газете, газетой, газетою, газету, газеты},
- {герое, героев, героем, герои, герой, герою, героя, героям, героями, героях},
- {гетеросексуальности, гетеросексуальность, гетеросексуальностью},
- {голов, голова, головам, головами, головах, голове, головой, головою, голову, головы},
- {дверей, двери, дверь, дверьми, дверью, дверям, дверями, дверях},
- {дев, дева, девам, девами, девах, деве, дево, девой, девою, деву, девы},
- {день, дне, дней, днём, дни, дню, дня, дням, днями, днях},
- {детей, дети, детьми, детям, детях, ребёнка, ребёнке, ребёнком, ребёнку, ребёнок, ребят, ребята, ребятам, ребятами, ребятах},
- {друг, друга, друге, другом, другу, друже, друзей, друзья, друзьям, друзьями, друзьях},
- {железа, железам, железами, железах, железе, железой, железою, железу, железы, желёз},
- {жилищ, жилища, жилищам, жилищами, жилицах, жилище, жилищем, жилишу},
- {житие, житием, житии, житий, житию, жития, житиям, житиями, житиях},
- {журнал, журнала, журналам, журналами, журналах, журнале, журналов, журналом, журналу, журналы},
- {искр, искра, искрам, искрами, искрах, искре, искрой, искрою, искру, искры},
- {историей, историою, истории, историй, историю, история, историям, историями, историях},
- {камень, каменьев, каменья, каменьям, каменьями, каменьях, камне, камней, камнем, камни, камню, камня, камням, камнями, камнях},
- {карандаши, карандаша, карандашам, карандашами, карандашах, карандаше, карандашей, карандаши, карандашом, карандашу},
- {книг, книга, книгам, книгами, книгах, книге, книги, книгой, книгою, книгу},
- {колен, колена, коленам, коленами, коленах, колене, коленей, колени, колено, коленом, колену, коленьев, коленья, коленьям, коленьями, коленях, коленям, коленями, коленях},
- {копеек, копейка, копейкам, копейками, копейках, копейке, копейки, копейкой, копейкою, копейку},
- {лёд, льда, льдам, льдами, льдах, льде, льдов, льдом, льду, льды},

- {лоскут, лоскута, лоскутам, лоскутами, лоскутах, лоскуте, лоскутов, лоскутом, лоскуту, лоскуты, лоскутьев, лоскутья, лоскутьям, лоскутьями, лоскутьях},
- {лошадей, лошади, лошадь, лошадьми, лошадью, лошадям, лошадями, лошадях},
- {мальчат, мальчата, мальчатам, мальчатами, мальчатах, мальчонка, мальчонкам, мальчонками, мальчонках, мальчонке, мальчонки, мальчонков, мальчонком, мальчонку, мальчонок},
- {море, морей, морем, морю, моря, морям, морями, морях},
- {мост, моста, мостам, мостами, мостах, мосте, мостов, мостом, мосту, мосты},
- {муж, мужа, муже, мужей, мужем, мужу, мужья, мужьям, мужьями, мужьях},
- {неделе, неделей, неделено, недели, недел, неделю, неделя, неделям, неделами, неделях},
- {нож, ножа, ножам, ножами, ножах, ноже, ножей, ножи, ножом, ножу},
- {няне, няней, нянею, няни, нянь, няню, няня, няням, нянями, нянях},
- {озерца, озерце, озерцо, озерцом, озерцу, озёрец, озёрца, озёрцам, озёрцами, озёрцах},
- {окна, окнам, окнами, окнах, окне, окно, окном, окну, окон},
- {ореол, ореола, ореолам, ореолами, ореолах, ореоле, ореолов, ореолом, ореолу, ореолы},
- {отец, отца, отцам, отцами, отцах, отце, отцов, отцом, отцу, отцы},
- {пальца, пальцам, пальцами, пальцах, пальце, пальцев, пальцем, пальцу, пальцы},
- {писем, письма, письмам, письмами, письмах, письме, письмо, письмом, письму},
- {площадей, площади, площадь, площадью, площадям, площадями, площадях},
- {повод, повода, поводе, поводом, поводу, поводьев, поводья, поводьям, поводьями, поводьях},
- {поле, полей, полем, полю, поля, полям, полями, полях},
- {пух, пуха, пухе, пухом, пуху},
- {работ, работа, работам, работами, работах, работе, работой, работою, работу, работы},
- {рот, рта, ртам, ртами, ртах, рте, ртов, ртом, рту, рты},
- {сапожек, сапожка, сапожкам, сапожками, сапожках, сапожке, сапожки, сапожков, сапожком, сапожку, сапожок},
- {сестёр, сестра, сестре, сестрой, сестрою, сестру, сестры, сёстрам, сёстрами, сёстрах, сёстры},
- {словаре, словарей, словарём, словари, словарь, словарю, словаря, словарям, словарями, словарях},
- {слушае, слушаев, слушаем, случай, случаю, случая, случаям, случаями, случаях},
- {сна, снам, снами, снах, сне, снов, сном, сну, сны, сон},
- {снег, снега, снегам, снегами, снегах, снеге, снегов, снегом, снегу},
- {стол, стола, столам, столами, столах, столе, столов, столом, столу, столы},
- {сын, сына, сынам, сынами, сынах, сыне, сынов, сыновей, сыновья, сыновьям, сыновьями, сыновьях, сыном, сыну, сыны},
- {товарищ, товарища, товарищам, товарищами, товарищах, товарице, товарищей, товарищем, товарищи, товаришу},

{урок, урока, урокам, уроками, уроках, уроке, уроки, уроков, уроком, уроку},  
 {учащегося, учащемсяся, учащемуся, учащеся, учащийся, учащимися, учащимся, учащихся, учителе, учителей, учителем, учители, учитель, учителю, учителя, учителям, учителями, учителях},  
 {учение, учением, учении, учений, учению, учения, учениям, учениями, учениях},  
 {училищ, училища, училищам, училищами, училищах, училище, училищем, училищу},  
 {феномен, феномена, феноменам, феноменами, феноменах, феномене, феноменов, феноменом, феномену, феномены},  
 {хвала, хвалам, хвалами, хвалах, хвале, хвалой, хвалою, хвалу, хвалы},  
 {чай, чаём, чаёв, чай, чаю, чая, чаям, чаями, чаях},

As we send the following sample adjectives (including participles of certain verbs)

большая, большего, большее, большей, большем, большему, большею, большие, больший, большим, большими, больших, большого, большое, большой, большом, большому, большою, большую, велик, велика, велики, велико, наибольшая, наибольшего, наибольшее, наибольшей, наибольшем, наибольшему, наибольшую, наибольшие, наибольший, наибольшим, наибольшими, наибольших, наибольшую — “big”;

делавшая, делавшего, делавшее, делавшей, делавшем, делавшему, делавшую, делавшие, делавший, делавшим, делавшими, делавших, делавшую — “do (past active participles)”;

делающая, делающего, делающее, делающей, делающим, делающему, делающею, делающие, делающий, делающим, делающими, делающих, делающую — “do (present active participles)”;

дорог, дорога, дорогая, дороги, дорогие, дорогим, дорогими, дорогих, дорого, дорогого, дорогое, дорогой, дорогом, дорогому, дорогою, дорогую, дороже — “expensive”;

красив, красива, красивая, красивее, красивейший, красиво, красивого, красивое, красивой, красивом, красивому, красивою, красивую, красивы, красивые, красивый, красивым, красивыми, красивых, краше — “beautiful”;

крут, крута, крутая, крутейший, круто, крутого, крутое, крутой, крутом, крутому, крутою, крутую, круты, крутые, крутым, крутыми, крутых, круче — “steep”;

любим, любима, любимая, любимо, любимого, любимое, любимой, любимом, любимому, любимою, любимию, любими, любимые, любимый, любимым, любимыми, любимых — “love (present passive participles)”;

молод, молода, молодая, молодо, молодого, молодое, молодой, молодом, молодому, молодою, молодую, молоды, молодые, молодым, молодыми, молодых, моложе — “young”;

мягка, мягкая, мягки, мягкие, мягкий, мягким, мягкими, мягких, мягко, мягкого, мягкое, мягкой, мягким, мягкому, мягкою, мягкую, мягок, мягчайший, мягче, наимягчайший — “soft”;

решена, решено, решены, решён, решённая, решённого, решённое, решённой, решённом, решённому, решённою, решённую, решённые, решённый, решённым, решёнными, решённых — “decide (past passive perfect participles)”;

тих, тиха, тихая, тихи, тихие, тихий, тихим, тихими, тихих, тихо, тихого, тихое, тихой, тихом, тихому, тихою, тихую, тишайший, тишие — “silent”;

част, часта, частая, часто, частого, частое, частой, частом, частому, частою, частую, часты, частые, частый, частым, частыми, частых, чаще — “frequent”.

to our algorithm, we receive the following results:

{большая, большего, большее, большей, большем, большему, большею, большие, больший, большим, большими, больших, большого, большое, большой, большом, большому, большою, большую, велик, велика, велики, велико, наибольшая, наибольшего, наибольшее, наибольшей, наибольшем, наибольшему, наибольшую, наибольшие, наибольший, наибольшим, наибольшими, наибольших, наибольшую},

{делавшая, делавшего, делавшее, делавшей, делавшем, делавшему, делавшею, делавшие, делавший, делавшим, делавшими, делавших, делавшую, делающая, делающего, делающее, делающей, делающим, делающему, делающею, делающие, делающий, делающим, делающими, делающих, делающую},

{дорог, дорога, дорогая, дороги, дорогие, дорогим, дорогими, дорогих, дорого, дорогого, дороже, дорогой, дорогом, дорогому, дорогою, дорогую, дороже},

{красив, красива, красивая, красивее, красивейший, красиво, красивого, красивое, красивой, красивом, красивому, красивою, красивую, красивы, красивые, красивый, красивым, красивыми, красивых, красище},

{крут, крута, крутая, крутейший, круто, крутого, крутое, крутой, крутом, крутому, крутою, крутую, круты, крутые, крутым, крутыми, крутых, круче},

{любим, любима, любимая, любимо, любимого, любимое, любимой, любимом, любимому, любимою, любимую, любими, любимые, любимый, любимым, любимыми, любимых},

{молод, молода, молодая, молodo, молодого, молодое, молодой, молодом, молодому, молодою, молодую, молоды, молодые, молодым, молодыми, молодых, моложе},

{мягка, мягкая, мягки, мягкие, мягкий, мягким, мягкими, мягких, мягко, мягкого, мягкое, мягкой, мягким, мягкому, мягкою, мягкую, мягок, мягчайший, мягче, наимягчайший},

{решена, решено, решены, решён, решённая, решённого, решённое, решённой, решённом, решённому, решённою, решённую, решённые, решённый, решённым, решёнными, решённых},

{тих, тиха, тихая, тихи, тихие, тихий, тихим, тихими, тихих, тихо, тихого, тихое, тихой, тихом, тихому, тихою, тихую, тишайший, тишие},

{част, часта, частая, часто, частого, частое, частой, частом, частому, частою, частую, часты, частые, частый, частым, частыми, частых, чаще},

*Example 7.26.2.* Like Polish, a great majority of Russian verbs come in imperfective/perfective pairs. Consonant alternations still find their way into the Russian verb stems (more precisely, the end of stems) in conjugated forms. In the following, we list representative verbs from all the 16 Zaliznyak classes ([https://en.wiktionary.org/wiki/Appendix:Russian\\_verbs](https://en.wiktionary.org/wiki/Appendix:Russian_verbs)), together with their imperfective/perfective counterparts that are formed in various ways (we note that the two aspects of the same verb may not belong to the same Zaliznyak class). These verb conjugations may also involve consonant changes in the stem, or other types of irregularities. The Zaliznyak class number, which will be shown parenthetically after the English translation for each word family below, describes at least one aspect (imperfective or perfective) of the verb in question.

делав, делавши, делавший, делаем, делаемый, делает, делаете, делаешь, делай, делайте, делал, делала, делали, делало, деланный, делать, делаю, делают, делающий, делая, сделав, сделавши, сделавший, сделаем, сделает, сделаешь, сделай, сделайте, сделал, сделала, сделали, сделало, сделанный, сделать, сделаю, сделаю — “do” (1);

нарисовав, нарисовавши, нарисовавший, нарисовал, нарисовала, нарисовали, нарисовало, нарисованный, нарисовать, нарисуем, нарисует, нарисуете, нарисуешь, нарисуй, нарисуйте, нарисую, нарисуют, рисовав, рисовавши, рисовавший, рисовал, рисовала, рисовали, рисовало, рисованный, рисовать, рисуем, рисуемый, рисует, рисуете, рисуешь, рисуй, рисуйте, рисую, рисуют, рисующий, рисуя — “draw” (2a);

блевав, блевавши, блевавший, блевал, блевала, блевали, блевало, блеванём, блеванёт, блеванёте, блеванёшь, блевани, блеваните, блевану, блеванув, блеванувши, блеванувший, блеванул, блеванула, блеванули, блевануло, блеванут, блевануть, блевать, блёванный, блюём, блюёт, блюёте, блюёшь, блюй, блюйт, блюют, блюющий, блюя — “vomit” (2b);

гиб, гибла, гибли, гибло, гибнем, гибнет, гибнете, гибнешь, гибни, гибните, гибну, гибнув, гибнуши, гибнуший, гибнул, гибнут, гибнуть, гибущий, погиб, погибла, погибли, погибло, погибнем, погибнет, погибнете, погибнейшь, погибни, погибните, погибну, погибнут, погибнуть, погибши, погибший — “die” (3a);

искнём, рискнёт, рискнёте, рискнёшь, рискни, рискните, рискну, рискнув, рискнуши, рискнуший, рискнул, рискнула, рискнули, рискнуло, рискнут, рискнуть, рисковав, рисковавши, рисковавший, рисковал, рисковала, рисковали, рисковало, рисковать, рискуем, рискует, рискуете, рискуешь, рискуй, рискуйте, рискую, рискуют, рискующий, рискуя — “risk” (3b);

взглядывав, взглядавши, взглядавший, взглядываем, взглядываемый, взглядывает, взглядываете, взгля-  
дываешь, взглядывай, взглядывайте, взглядал, взглядовали, взглядовало, взглядывать, взгля-  
дываю, взглядывают, взглядывающий, взглядывая, взглянем, взглянет, взглянете, взглянешь, взгляни, взгля-  
ните, взгляну, взглянув, взглянуши, взглянувший, взглянул, взглянула, взглянули, взглянуло, взглянут, взгля-  
нуть, глядев, глядевши, глядевший, глядел, глядела, глядели, глядело, глядеть, гляди, глядим, глядит, гля-  
дите, глядишь, глядя, глядят, глядящий, гляжу — “glance” (3c);

жаленный, жалив, жаливши, жаливший, жалил, жалила, жалили, жалило, жалим, жалимый, жалит,  
жалите, жалить, жалишь, жаль, жальте, жалю, жала, жалят, жалищий, ужаленный, ужалив, ужалив-  
ши, ужаливший, ужалил, ужалила, ужалили, ужалило, ужалим, ужалит, ужалите, ужалить, ужалишь,  
ужаль, ужальте, ужалю, ужалят — “sting” (4a);

пощади, пощадив, пощадивши, пощадивший, пощадил, пощадила, пощадили, пощадило, пощадим, поща-  
дит, пощадите, пощадить, пощадишь, пощадят, пощажённый, пощажсу, щади, щадив, щадивши, ща-  
дивший, щадил, щадила, щадили, щадило, щадим, щадимый, щадит, щадите, щадить, щадишь, щадя,  
щадят, щадящий, щажённый, щажсу — “spare” (4b);

люби, любив, любивши, любивший, любил, любила, любили, любило, любим, любимый, любит, любите, лю-  
бить, любишь, любленный, люблю, любя, любят, любящий, полюби, полюбив, полюбивши, полюбивший,  
полюбил, полюбила, полюбили, полюбило, полюбим, полюбим, полюбите, полюбить, полюбиишь, полюб-  
ленный, полюблю, полюбят — “love” (4c);

слуша, слышав, слышавши, слышавший, слышал, слышала, слышали, слышало, слышанный, слышат, слы-  
шать, слышащий, слышим, слышимый, слышит, слышите, слышишь, слышу, услышав, услышавши, услы-  
шавший, услышал, услышала, услышали, услышало, услышанный, услышат, услышать, услышим, услы-  
шит, услышите, услышишь, услышу, услышь, услышьте — “hear” (5a);

бренча, бренчав, бренчавши, бренчавший, бренчал, бренчала, бренчали, бренчало, бренчат, бренчать, брен-  
чащий, бренчи, бренчим, бренчит, бренчите, бренчишь, бренчу — “jingle” (5b);

изгнав, изгнавши, изгнавший, изгнал, изгнала, изгнали, изгнало, изгнанный, изгнать, изгони, изгоним, изго-  
нит, изгоните, изгонишь, изгоню, изгоняв, изгонявш, изгоняший, изгоняем, изгоняемый, изгоняет, изго-  
няете, изгоняешь, изгоняй, изгоняйте, изгонял, изгоняла, изгоняли, изгоняло, изгонят, изгонять, изгоняю,  
изгоняют, изгоняющий, изгоняя — “banish” (5c);

веем, веемый, веет, веете, веешь, вей, вейте, вею, веют, веющий, вея, веяв, веявиши, веавший, веял, веяла,  
веали, веяло, веянный, веять — “flutter” (6a);

взыдав, взыдавши, взыдавший, взыдаем, взыдает, взыдывает, взыдаешь, взыдай, взыдайте, взыдал, взыдава-  
ли, взыдаво, взыдавать, взыдаю, взыдают, взыдающий, взыдавая, воззвав, воззвавши, воззвавший, воззвал,  
воззвала, воззвали, воззвало, воззванный, воззвать, воззовём, воззовёт, воззовёте, воззовёшь, воззови, воз-  
зовите, воззовут — “appeal” (6b);

взыскав, взыскавши, взыскавший, взыскал, взыскала, взыскали, взыскало, взысканный, взыскать, взыс-  
киав, взыскиавши, взыскиавший, взыскиаем, взыскиаемый, взыскиает, взыскиваете, взыскиваешь,  
взыскивай, взыскивайте, взыскивал, взыскивала, взыскивали, взыскивало, взыскивать, взыскиваю, взыски-  
вают, взыскивающий, взыскивая, взыщем, взыщет, взыщете, взыщешь, взыщи, взыщите, взыщу, взыщут  
— “recover” (6c);

влез, влезав, влезавши, влезавший, влезаем, влезает, влезаете, влезаешь, влезай, влезайт, влезал, влезала,  
влезали, влезало, влезать, влезаю, влезают, влезающий, влезая, влезем, влезет, влезете, влезешь, влезла,  
влезли, влезло, влезть, влезу, влезут, влезши, влезший, влезь, влезьте — “meddle” (7a);

ведём, ведённый, ведёт, ведёте, ведёшь, веди, ведите, ведомый, веду, ведут, ведущий, ведши, ведший, ве-  
дя, вела, вели, вело, вести, поведём, поведённый, поведёт, поведёте, поведёшь, поведи, поведите, поведу,  
поведут, поведши, поведший, поведя, повела, повели, повело, повести — “lead” (7b);

вытек, вытекав, вытекавши, вытекавший, вытекаем, вытекает, вытекаете, вытекаешь, вытекай, вы-  
текайт, вытекал, вытекала, вытекали, вытекало, вытекать, вытекаю, вытекают, вытекающий, вы-  
текая, вытеки, вытеките, вытекла, вытекли, вытекло, вытеку, вытекут, вытекши, вытекший, выте-  
чем, вытечет, вытечете, вытечешь, вытеч, теки, теките, текли, текло, теку, текут, текущий,  
течём, течёт, течёте, течёшь, течь, тёк, тёкии, тёкий — “leak” (8a);

береги, берегите, берегла, берегли, берегло, берегу, берегут, берегущий, бережём, бережённый, бережёт, бережёте, бережёшь, беречь, берёг, берёгши, берёгший, сбереги, сберегите, сберегла, сберегли, сберегло, сберегу, сберегут, сбережём, сбережённый, сбережёт, сбережёте, сбережёшь, сберечь, сберёг, сберёгши, сберёгший, убереги, уберегите, уберегла, уберегли, уберегло, уберегу, уберегут, убережём, убережённый, убережёт, убережёте, убережёшь, уберечь, уберёг, уберёгши, уберёгший — “guard” (8b);

вымер, вымерев, вымереть, вымерла, вымерли, вымерло, вымерши, вымерший, вымирав, вымиравши, вымиравший, вымираем, вымирает, вымираете, вымираешь, вымирай, вымирайте, вымирали, вымирала, вымирали, вымирало, вымират, вымираю, вымирают, вымирающий, вымирая, вымрем, вымрет, вымрете, вымрешь, вымы, вымрите, вымру, вымрут — “becomeextinct” (9a);

потерев, потереть, потёр, потёрга, потёрги, потёрло, потёртый, потёриши, потёришай, потрём, потрёт, потрёте, потрёшишь, потри, потрите, потру, потрут, тереть, тёр, тёргла, тёрги, тёрло, тёртый, тёриши, тёришай, трём, трёт, трёте, трёшишь, три, трите, тру, трут, трущий — “rub” (9b);

выпальывав, выпальывавши, выпальывавший, выпальываем, выпальываемый, выпальывает, выпальываете, выпальываешь, выпальывай, выпальывайте, выпальвал, выпальвали, выпальвало, выпальывало, выпальывать, выпальываю, выпальывают, выпальывающий, выпальывая, выполем, выполет, выполните, выполнешь, выполи, выполните, выполнов, выполновши, выполновший, выполнол, выполнола, выполноли, выполноло, выполнотый, выполнить, выполню, выполнют — “weed” (10a);

вкалывав, вкалывавши, вкалывавший, вкалываем, вкалываемый, вкалывает, вкалываете, вкалываешь, вкалывай, вкалывайте, вкалывал, вкалывала, вкалывали, вкалывало, вкалывать, вкалываю, вкалывают, вкалывающий, вкалывая, вколем, вколет, вколете, вколеши, вколи, вколите, вколов, вколовши, вколовший, вколол, вкололи, вкололо, вколотый, вколоть, вколою, вколют — “stick” (10c);

вылейся, вылейтесь, выливавшийся, выливавши, выливается, выливаетесь, выливается, выливаешься, выливайся, выливайтесь, выливалась, выливались, выливалось, выливался, выливаться, выливаюсь, выливаются, выливающийся, выливаясь, вылившийся, вылившись, вылилась, вылились, вылилось, вылился, выльться, выльемся, выльетесь, выльешься, выльюсь, выльются — “spill over” (11a);

добей, добейте, добив, добивав, добивавши, добивавший, добиваем, добиваемый, добивает, добиваете, добиваешь, добивай, добивайте, добивал, добивала, добивали, добивало, добивать, добиваю, добивают, добивающий, добивая, добивши, добивший, добил, добила, добили, добило, добитый, добить, добьём, добьёт, добьёте, добьёши, добью, добьют — “crush” (11b);

вымоем, вымоете, вымоешь, вымой, вымойте, вымою, вымоят, вымыв, вымывши, вымывший, вымыл, вымыла, вымыли, вымыло, вымытый, вымыть, моем, моемый, моет, моете, моешь, мой, мойте, мою, моют, моющий, моя, мыв, мывши, мывший, мыл, мыла, мыли, мыло, мытый, мыть, помоем, помоет, помоете, помоешь, помой, помойте, помою, помоют, помыв, помывши, помывший, помыл, помыла, помыли, помыло, помытый, помыть, умоем, умоет, умоетесь, умоишь, умой, умойте, умою, умоют, умыв, умывши, умывший, умыл, умыла, умыли, умыло, умытый, умыть — “wash” (12a);

гнив, гнивши, гнивший, гниём, гниёт, гниёте, гниёшь, гнил, гнила, гнили, гнило, гнить, гнию, гниют, гниющий, сгнив, сгнивши, сгнивший, сгниём, сгниёт, сгниёте, сгниёшь, сгнил, сгнила, сгнили, сгнило, сгнить, сгнию, сгниют — “putrefy” (12b);

дав, давав, дававши, дававший, даваемый, давай, давайте, давал, давала, давали, давало, давать, давая, давши, давший, дадим, дадите, дадут, даём, даёт, даёте, даёшишь, дай, дайте, дал, дала, дали, дало, дам, данный, даст, дать, даши, даю, дают, дающий — “give” (13);

выжав, выжавши, выжавший, выжал, выжала, выжали, выжало, выжатый, выжать, выжимав, выжимавши, выжимавший, выжимаем, выжимаемый, выжимает, выжимаете, выжимаешь, выжимай, выжимайте, выжимал, выжимала, выжимали, выжимало, выжимать, выжимаю, выжимают, выжимающий, выжимая, выжмем, выжмет, выжмете, выжмешь, выжми, выжмите, выжму, выжмут — “wring” (14a);

жав, жавши, жавший, жал, жала, жали, жало, жатый, жать, жмём, жмёт, жмёте, жмёшь, жми, жмите, жму, жмут, жмущий, жмя, пожав, пожавши, пожавший, пожал, пожала, пожали, пожало, пожатый, пожать, пожмём, пожмёт, пожмёте, пожмёшишь, пожми, пожмите, пожму, пожмут, сжав, сжавши, сжавший, сжал, сжала, сжали, сжало, сжатый, сжать, сожмём, сожмёт, сожмёте, сожмёшь, сожми, сожмите, сожму, сожмут — “squeeze” (14b);

обнимав, обнимавши, обнимавший, обнимаем, обнимаемый, обнимает, обнимаете, обнимаешь, обнимай, обнимайте, обнимал, обнимала, обнимали, обнимало, обнимать, обнимаю, обнимают, обнимающий, обнимая, обнимем, обнимете, обнимешь, обними, обнимите, обниму, обнимут, обняв, обнявши, обнявший, обнял, обняла, обняли, обняло, обнятый, обнять — “hug” (14c);

встав, вставав, встававши, встававший, вставай, вставайте, вставал, вставала, вставали, вставало, вставать, вставая, вставши, вставший, встаём, встаёт, встаёте, встаёшь, встал, встала, встали, встало, встанем, встанет, встанете, встанешь, встану, встанут, встань, встаньте, встать, встаю, встают, встающий — “rise” (15);

выжив, выживав, выживавши, выживавший, выживаем, выживает, выживаете, выживашь, выживай, выживайте, выживал, выживала, выживали, выживало, выживать, выживую, выживают, выживают, выживший, выживая, выживем, выживет, выживете, выживешь, выживи, выживите, выживу, выживут, выживши, выживший, выжил, выжила, выжили, выжило, выжить — “survive” (16a);

жив, живём, живёт, живёте, живёшь, живи, живите, живу, живут, живущий, живши, живший, живя, живл, живла, живли, живло, жить, пожив, поживём, поживёт, поживёте, поживешь, поживи, поживите, поживу, поживут, поживши, поживший, пожил, пожила, пожили, пожило, пожить — “live” (16b).

These verbs are clustered by our algorithms as follows:

{береги, берегите, берегла, берегли, берегло, берегу, берегут, берегущий, бережём, бережённый, бережёт, бережёте, бережёшь, беречь, берёг, берёгши, берёгший, сбереги, сберегите, сберегла, сберегли, сберегло, сберегу, сберегут, сбережём, сбережённый, сбережёт, сбережёте, сбережёшь, сберечь, сбёг, сберёгши, сберёгший, убереги, уберегите, уберегла, уберегли, уберегу, уберегут, убережём, убережённый, убережёт, убережёшь, уберечь, уберёг, уберёгши, уберёгший},

{блевав, блевавши, блевавший, блевал, блевала, блевали, блевало, блеванём, блеванёт, блеванёте, блеванёшь, блевани, блеваните, блевану, блеванув, блеванувши, блеванувший, блеванул, блеванула, блеванули, блевануло, блеванут, блевануть, блевать, блёванный, блюём, блюёт, блюёте, блюёшь, блюй, блюйте, блюю, блюют, блюющий, блюя},

{бренча, бренчав, бренчавши, бренчавший, бренчал, бренчала, бренчали, бренчало, бренчат, бренчать, бренчащий, бренчи, бренчим, бренчит, бренчите, бренчишь, бренчу},

{ведём, ведённый, ведёт, ведёте, ведёшь, веди, ведите, ведомый, ведши, ведший, ведя, вела, вели, вело, вести, поведём, поведённый, поведёт, поведёте, поведёшь, поведи, поведите, поведу, поведут, поведши, поведший, поведя, повела, повели, повело, повести},

{веду, ведут, ведущий},

{веем, веемый, веете, веешь, вей, вейте, вею, веют, веющий, вея, веяв, веяви, веявиши, веял, веяла, веали, веяло, веянный, веять},

{веет},

{взглядывав, взглядыавши, взглядыавший, взглядываем, взглядывает, взглядываете, взглядаешь, взглядывой, взглядывой, взглядывайт, взглядывал, взглядывала, взглядывали, взглядывало, взглядывать, взглядыва, взглядывают, взглядывающ, взглядывая, взглянем, взглянет, взглянете, взглянешь, взгляни, взгляните, взгляну, взглянув, взглянуши, взглянувший, взглянул, взглянула, взглянули, взглянуло, взглянут, взглянуть, глядев, глядевши, глядевший, глядел, глядела, глядели, глядело, глядеть, гляди, глядим, глядит, глядите, глядишь, глядя, глядят, глядящий, гляжу},

{взыдав, взыдавши, взыдавший, взыдаем, взыдает, взыдываете, взыдаешь, взыдай, взыдайте, взыдал, взыдала, взыдвали, взыдало, взыдвать, взыдаю, взыдают, взыдают, взыдающий, взыдая, взыдем, взыдеш, взыдите, взыдешь, взыди, взыдите, взыду, взыдуют},

{взыскав, взыскавши, взыскавший, взысканный, взыскать, взыскивав, взыскивавши, взыскивавший, взыскавем, взыскаваемый, взыскивает, взыскиваете, взыскиваешь, взыскивай, взыскивайт, взыскивал, взыскавала, взыскавали, взыскавало, взыскавать, взыскиваю, взыскивают, взыскивающий, взыскивав},

{взыскал, взыскала, взыскали, взыскало},

{вкалывав, вкалывавши, вкалывавший, вкалываем, вкалываемый, вкалывает, вкалываете, вкалываешь, вкалывай, вкалывайте, вкалывал, вкалывала, вкалывали, вкалывало, вкалывать, вкалываю, вкалывают, вкалывающий, вкалывая, вколем, вколет, вколете, вколеши, вcoli, вколите, вколов, вколовши, вколовший, вколол, вколола, вкололи, вкололо, вколотый, вколоть, вколо, вколют},

{влез, влезав, влезавши, влезавший, влезаем, влезает, влезаете, влезаешь, влезай, влезайте, влезал, влезала, влезали, влезало, влезать, влезаю, влезают, влезающий, влезая, влезем, влезет, влезете, влезешь, влезла, влезли, влезло, влезть, влезу, влезут, влезши, влезший, влезь, влезьте},

{воззвав, воззвавши, воззвавший, воззвал, воззвала, воззвали, воззвало, воззванный, воззвать, воззовём, воззовёт, воззовёте, воззовёшь, воззови, воззовите, воззову, воззовут},

{встав, вставав, встававши, встававший, вставай, вставайте, вставил, вставала, вставали, вставало, вставать, вставая, вставши, вставший, встаём, встаёт, встаёте, встаёшь, встал, встала, встали, встало, встанем, встанет, встанете, встанешь, встану, встанут, встань, встаньте, встать, встаю, встают, встающий},

{выжав, выжавши, выжавший, выжал, выжала, выжали, выжало, выжатый, выжать, выжив, выжива-  
вав, выживавши, выживавший, выживаем, выживает, выживаете, выживашь, выживай, выживайте,  
выживал, выживала, выживали, выживало, выживать, выживая, выживают, выживавший, выживая,  
выживем, выживет, выживете, выживешь, выживи, выживите, выживу, выживут, выживши, выжив-  
ший, выжил, выжила, выжили, выжило, выжимав, выжимавши, выжимавший, выжимаем, выжимае-  
мый, выжимает, выжимаете, выжимашь, выжимай, выжимайте, выжимал, выжимала, выжимали,  
выжимало, выжимать, выжимаю, выжимают, выжимающий, выжимая, выжить, выжмем, выжмет,  
выжмете, выжмешь, выжми, выжмите, выжму, выжмут},

{вылейся, вылейтесь, выливавшийся, выливавшихся, выливаемся, выливаетесь, выливается, выливаешься,  
выливайся, выливайтесь, выливалась, выливались, выливалось, выливался, выливаться, выливаюсь, выли-  
ваются, выливающийся, выливаясь, вылившийся, вылившись, вылилась, вылились, вылилось, вылился, вы-  
льиться, выльемся, выльетесь, выльется, выльешься, выльюсь, выльются},

{вымер, вымерев, вымереть, вымерла, вымерли, вымерло, вымерши, вымерший, вымирав, вымиравши, вы-  
миравший, вымираем, вымирает, вымираете, вымираешь, вымирай, вымирайте, вымирал, вымирала, вы-  
мирали, вымирало, вымират, вымираю, вымирают, вымирающий, вымирая, вымрем, вымрет, вымрете,  
вымрешь, вымири, вымрите, вымру, вымрут},

{вымоем, вымоет, вымоете, вымоешь, вымой, вымойте, вымою, вымоют, вымыв, вымывиши, вымывши,  
вымыл, вымыла, вымыли, вымыло, вымытый, вымыть, моем, моемый, моет, моете, моешь, мой, мойте,  
мою, моют, моющий, моя, мыв, мывши, мывши, мывший, мыл, мыла, мыли, мыло, мытый, мыть, помоем,  
помоет, помоете, помоешь, помой, помойте, помою, помоют, помыв, помывши, помывши, помыл, помыла,  
помыли, помыло, помытый, помыть, умоем, умоет, умоете, умоишь, умой, умойте, умою, умоют, умыв,  
умывши, умывши, умыл, умыла, умыли, умыло, умытый, умыть},

{выпалывав, выпалывавши, выпалывавший, выпалываем, выпалываемый, выпалывает, выпалываете, вы-  
палываешь, выпалывай, выпалывайте, выпалывал, выпалывала, выпалывали, выпалывало, выпалывать,  
выпалываю, выпалывают, выпалывающий, выпалывая, выполем, выполет, выполните, выполнешь, выполи,  
выполните, выполнов, выполновши, выполновши, выполнол, выполнола, выполноли, выполноло, выполнолый, выполн-  
лоть, выполню, выполнют},

{вытек, вытекав, вытекавши, вытекавший, вытекаем, вытекает, вытекаете, вытекаешь, вытекай, вы-  
текайте, вытекал, вытекала, вытекали, вытекало, вытекать, вытекаю, вытекают, вытекающий, вы-  
текая, вытеки, вытеките, вытекла, вытекли, вытекло, вытеку, вытекут, вытекши, вытекший, выте-  
чем, вытечет, вытечете, вытечешь, вытечь, теки, теките, текли, текло, теку, текут, текущий,  
течём, течёт, течёте, течёши, течь, тёк, тёкиши, тёкий},

{гиб, гибла, гибли, гибло, гибнем, гибнет, гибните, гибнешь, гибни, гибните, гибну, гибнув, гибнувши, гиб-  
нувший, гибнул, гибнут, гибнуть, гибнущий, погиб, погибла, погибли, погибло, погибнем, погибнет, погиб-  
ните, погибнешь, погибни, погибните, погибну, погибнут, погибнуть, погибши, погибший},

{гнив, гнивши, гнивши, гниём, гниёт, гниёте, гниёшь, гнил, гнила, гнили, гнило, гниТЬ, гнию, гниют, гни-  
ющий, сгнив, сгнивши, сгнивши, сгниём, сгниёт, сгниёте, сгниёшь, сгнил, сгнила, сгнили, сгнило, сгниТЬ,  
сгнию, сгниют},

{дав, давав, дававши, дававший, даваемый, давай, давайте, давал, давала, давали, давало, давать, давая, давши, давший, дадим, дадите, дадут, даём, даёт, даёте, даёшь, дай, дайте, дал, дала, дали, дало, дам, данный, даст, дать, даешь, даю, дают, дающий},

{делав, делавши, делавший, делаем, делаемый, делает, делаете, делаешь, делай, делайте, делал, делала, делали, делало, деланный, делать, делаю, делают, делающий, делая, сделав, сделавши, сделавший, сделаем, сделает, сделаете, сделаешь, сделай, сделайте, сделал, сделала, сделали, сделало, сделанный, сделать, сделаю, сделают},

{добей, добейте, добив, добивав, добивавши, добивавший, добиваем, добиваемый, добивает, добиваете, добиваешь, добивай, добивайте, добивал, добивала, добивали, добивало, добивать, добиваю, добивают, добивающий, добивая, добивши, добивший, добил, добила, добили, добило, добитый, добить, добьём, добьёт, добьёте, добьёшь, добью, добьют},

{жасв, жасвиши, жасвий, жасл, жасла, жасленный, жасли, жаслив, жаслил, жасила, жасили, жасило, жаслим, жаслимый, жаслит, жаслите, жаслить, жаслиш, жасло, жасль, жасльте, жаслю, жасля, жаслят, жаслящий, жастый, жасть, жасём, жасёт, жасёте, жасёшь, жасми, жасмите, жасму, жасмут, жасущий, жасмя, пожасв, пожасвиши, пожасвий, пожасл, пожасала, пожасали, пожасло, пожастый, пожасть, пожасмёт, пожасмёте, пожасмёшь, пожасми, пожасмите, пожасму, пожасмут, сжасв, сжасвиши, сжасший, сжал, сжасала, сжасали, сжало, сжастый, сжасть, сожасмёт, сожасмёте, сожасмёшь, сожасми, сожасмите, сожасму, сожасмут},

{жасливши, жасливший, ужасливши, ужасливший},

{жив, живём, живёт, живёте, живёшь, живи, живите, живу, живут, живущий, живши, живший, живя, живил, жила, жили, жило, жить, пожив, поживём, поживёт, поживёте, поживёшь, поживи, поживите, поживу, поживут, поживши, поживший, пожил, пожила, пожили, пожило, пожить},

{изгнав, изгнавши, изгнавший, изгнал, изгнала, изгнали, изгнало, изгнанный, изгнать, изгони, изгоним, изгонит, изгоните, изгонишь, изгоню, изгоняв, изгонявши, изгонявший, изгоняем, изгоняется, изгоняет, изгояете, изгоняешь, изгоняй, изгоняйте, изгонял, изгоняла, изгоняли, изгоняло, изгонят, изгонять, изгоняю, изгоняют, изгоняющий, изгоняя},

{люби, любив, любивши, любивший, любил, любила, любили, любило, любим, любимый, любит, любите, любить, любишь, любленный, люблю, любя, любят, любящий, полюби, полюбив, полюбивши, полюбивший, полюбил, полюбила, полюбили, полюбило, полюбим, полюбит, полюбите, полюбить, полюбишь, полюблённый, полюблю, полюбят},

{нарисовав, нарисовавши, нарисовавший, нарисовал, нарисовала, нарисовали, нарисовало, нарисованный, нарисовать, нарисуем, нарисует, нарисуете, нарисуешь, нарисуй, нарисуйте, нарисую, нарисуют, рисовав, рисовавши, рисовавший, рисовал, рисовала, рисовали, рисовало, рисованный, рисовать, рисуем, рисуемый, рисует, рисуете, рисуешь, рисуй, рисуйте, рисую, рисуют, рисующий, рисуя},

{обнимав, обнимавши, обнимавший, обнимаем, обнимает, обнимаете, обнимаешь, обнимай, обнимайте, обнимал, обнимала, обнимали, обнимало, обнимать, обнимаю, обнимают, обнимающий, обнимая, обнимем, обнимете, обнимешь, обними, обнимите, обниму, обнимут, обняв, обнявши, обнявший, обнял, обняла, обняли, обняло, обнятый, обнять},

{потерев, потереть, потёр, потёрла, потёрли, потёрло, потёртый, потёриши, потёришний, потрём, потрёт, потрёте, потрёшь, потри, потрите, потру, потрут, тереть, тёр, тёра, тёрила, тёрло, тёртый, тёриши, тёришний},

{пощади, пощадив, пощадивши, пощадивший, пощадил, пощадила, пощадили, пощадило, пощадим, пощадит, пощадите, пощадить, пощадишь, пощадят, пощажённый, пощажжу, щади, щадив, щадивши, щадивший, щадил, щадила, щадили, щадило, щадим, щадимый, щадит, щадите, щадить, щадишь, щадя, щадят, щадящий, щажённый, щажжу},

{искнём, рискнёт, рискнёте, рискнёшь, рискни, рискните, рискну, рискнув, рискнувш, рискнувший, рискнул, рискнула, рискнули, рискнуло, рискнут, рискнуть, рисковав, рисковавши, рисковавший, рисковал, рисковала, рисковали, рисковало, рисковать, рискуем, рискует, рискуете, рискуешь, рискуй, рискуйте, рискую, рискуют, рискующий, рискуя},

{слыша, слышав, слышавши, слышавший, слышал, слышала, слышали, слышало, слышанный, слышат, слышать, слышащий, слышим, слышимый, слышит, слышите, слышишь, слышу, услышав, услышавши, услышавший, услышал, услышала, услышали, услышало, услышанный, услышат, услышать, услышим, услышит, услышите, услышишь, услышу, услышь, услышьте},

{трём, трёт, трёте, трёшь, три, трите, тру, трут, труций},

{ужаленный, ужалив, ужалил, ужалила, ужалили, ужалило, ужалим, ужалит, ужалите, ужалить, ужалишь, ужаль, ужальте, ужалю, ужалият},

Please beware that the verb clustering example above contains a few notable errors, partly due to interference between verbs with similar-looking stems.

*Example 7.26.3.* In Fig. S11, we further apply the aforementioned word clustering algorithm to topic extraction and machine translation (see Table S1 for text sources).

In Russian, “bat” (a kind of flying mammal) is called *летучая мышь* (literally “flying mouse”), so we consider *bat* an exact match to *летучая* “flying” in Fig. S11b”.

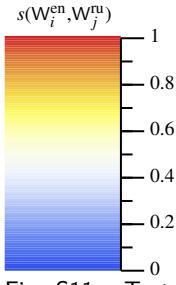
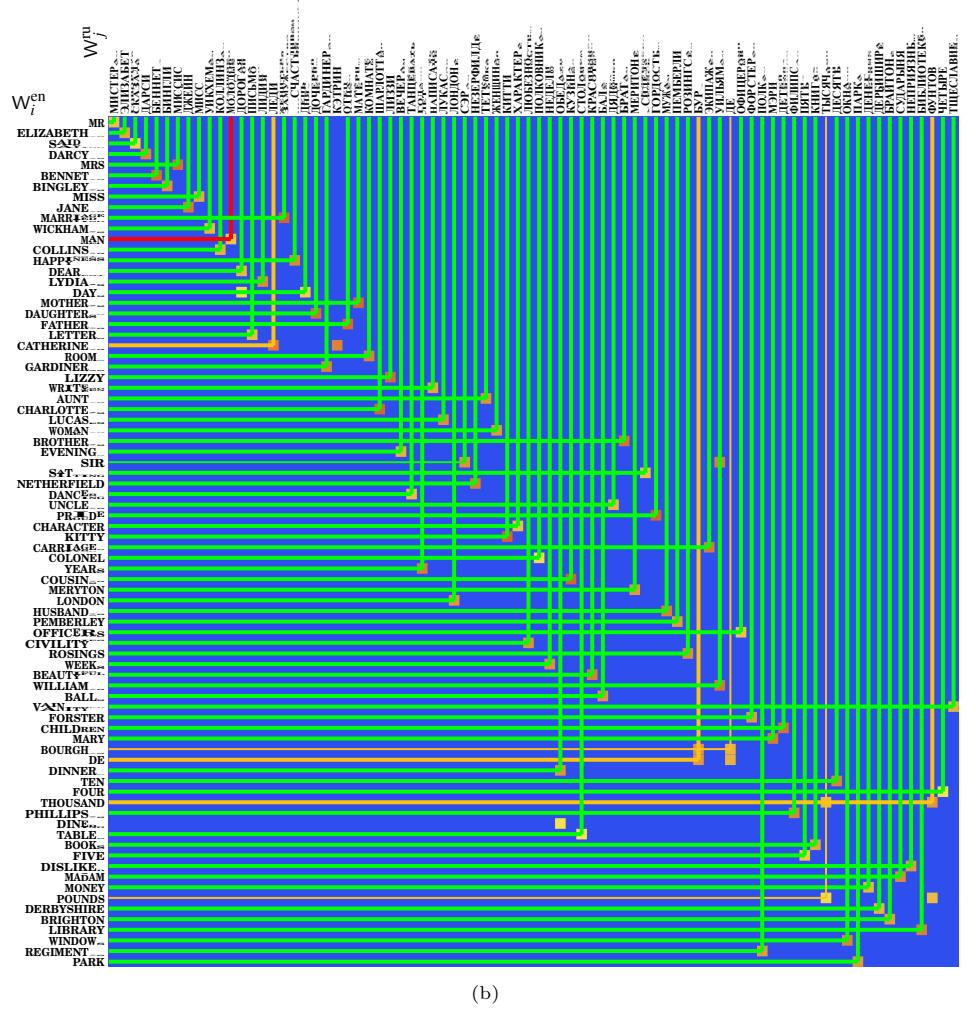


Fig. S11. Text mining in Russian. (a) Statistically identified topics ( $n_{ii} \geq 20$ ) in a Russian version of *Pride and Prejudice*, with the same color encoding scheme as Fig. S3. (b) Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{ru}})$  between selected topics in English and Russian versions of *Pride and Prejudice*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) ranking of the similarity score for the optimal node. Green (resp. red) cross-hair indicates a correct (resp. incorrect) match. Amber cross-hair marks a link between distinct concepts that share the same hypernyms.

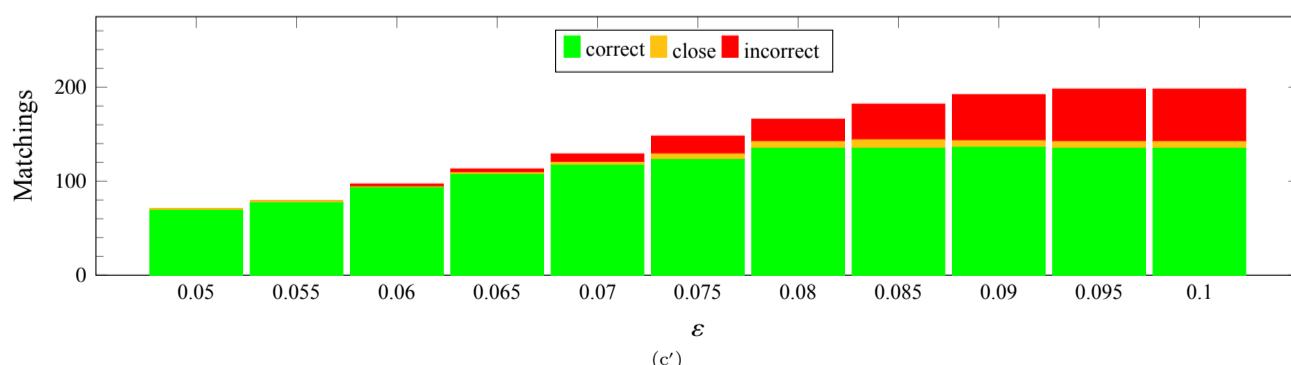
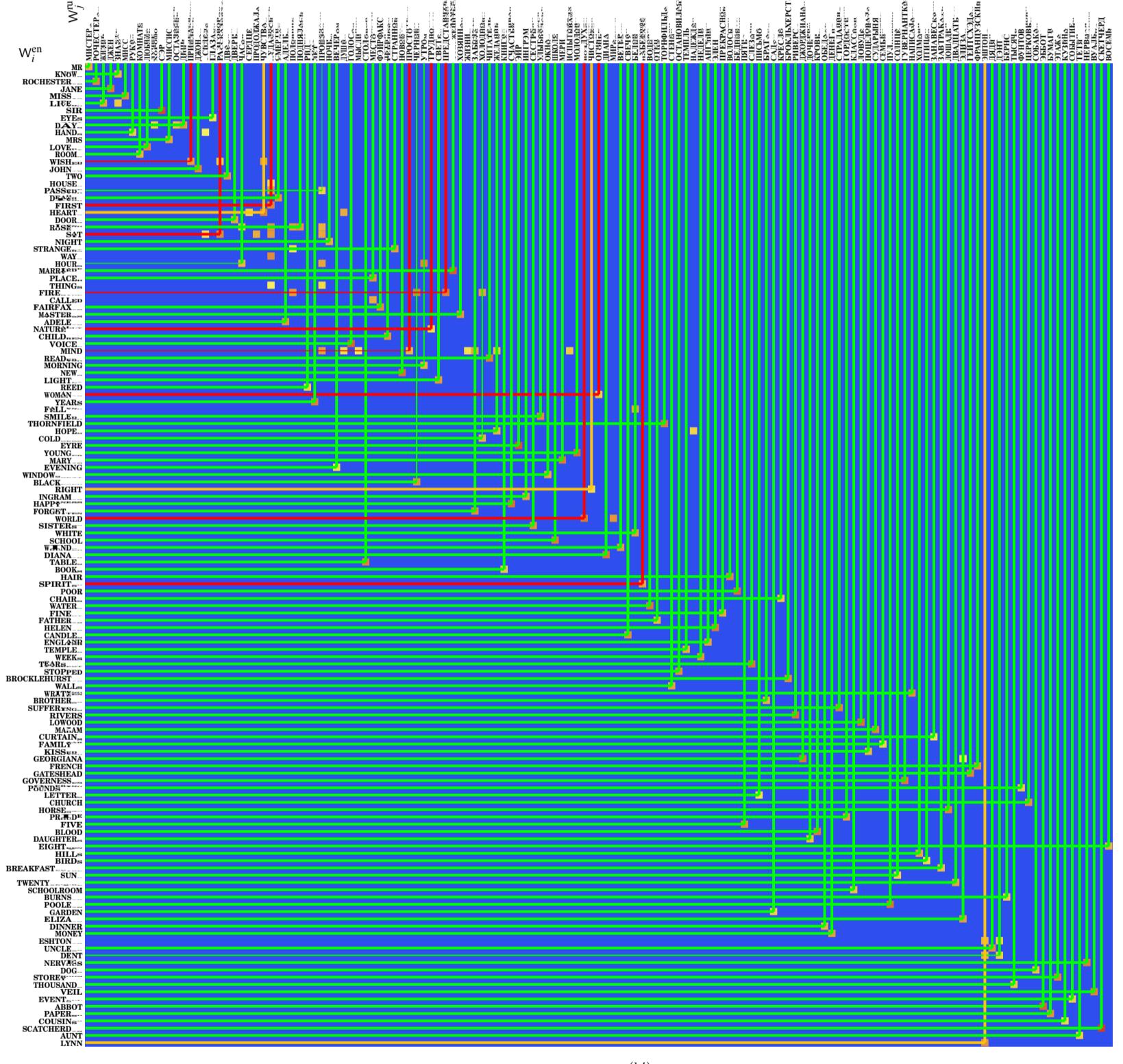


Fig. S11. Text mining in Russian. (Continued)  
(a') Statistically identified topics ( $n_{ii} \geq 20$ ) in a Russian version of *Jane Eyre*, with the same color encoding scheme as Fig. S3. (b') Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{ru}})$  between selected topics in English and Russian versions of *Jane Eyre*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) rank of the similarity score for the optimal node. Green (resp. red) cross-hair indicates a correct (resp. incorrect) match. Amber cross-hair marks a link between distinct concepts that share the same hypernyms. (c') Results from control experiments with different choices of the  $\varepsilon$ -parameter in ball-park screening criteria (1.13).



## 8 Approximate word clustering in selected Uralic languages

In this section, we present the approximate clustering algorithms for two representative Uralic languages: Finnish and Hungarian.

The major challenges in information processing for Finnish and Hungarian are the following:

- Both languages have extremely complex case systems (15 for Finnish and 18 for Hungarian, not counting rarely used cases), and the situation is compounded by the fact that several case endings may attach sequentially to the same word.
- Finnish has consonant gradation (mutation of stem-final consonants) during declension. (Estonian, another Uralic language not treated here, behaves similarly.) Hungarian case endings sometimes also cause changes in the stem-final consonants.
- Their verb morphologies have a similar complexity as the Romance languages.
- Similar to the Germanic counterpart, both Finnish and Hungarian orthographies require that certain compounds be written as single words, with no spaces or hyphens between their constituents.

Finnish and Hungarian are two major agglutinative languages spoken in Europe. Basque, a language isolate spoken in France and Spain, also has agglutinative morphology. The agglutinative nature of Turkish (as well as other Turkic languages spoken in Eurasia) is similar to Finnish and Hungarian. We will cover the word clustering algorithms for Basque and Turkish in §9.

### 8.1 Modified Porter stemming algorithm for Finnish

**Definition 8.1** (Finnish stop words). If a word belongs to the following list<sup>96</sup>:

aika, aina, ainainen, ainoa, ainoaa, ainoain, ainoan, ainoana, ainoata, ainoaten, ainoihin, ainoiksi, ainoil-  
la, ainoille, ainoilta, ainoin, ainoina, ainoine, ainoisiin, ainoissa, ainoista, ainoita, ainoitta, ainoitten, ainuiden,  
ainuihin, ainuiksi, ainuilla, ainuille, ainuulta, ainuin, ainuina, ainuineen, ainuisiin, ainuissa, ainuista, ainuita, ai-  
nuutta, ainuitten, ainut, ainutta, ainuksi, ainuulla, ainuulle, ainuulta, ainuun, ainuuna, ainuuseen, ainuussa, ai-  
nuusta, ainuut, ainuutta, aivan, ala, alakkain, alapuolella, alas, alatusten, alhaalla, alhaalle, alhaalta, ali, alitse,  
alla, alle, allekkain, alta, ansiosta, asti, edelle, edellä, edeltä, edes, edessä, edestä, editse, eduksi, ehkä, ei, elati-  
ve, eli, ellei, emme, en, enemmän, eniten, ennen, ennenkuin, ennestä, ennestään, ensi, ensin, entä, enää, erittäin,  
eräiden, eräihin, eräksi, eräille, eräillä, eräiltä, eräin, eräine, eräinä, eräisiin, eräissä, eräistä, eräitten, eräittä,  
eräitä, eräs, erästä, eräaksi, eräälle, eräällä, eräältä, erään, eräänä, eräaseen, eräässä, eräästä, eräät, eräättä,  
esi, esiin, et, eteemme, eteen, eteensä, etenkin, etenkään, ette, ettei, etten, ettet, että, halki, he, heidän, heidät, hei-  
hin, heille, heillä, heiltä, heissä, heistä, heitä, heti, hiukan, hyvin, hän, häneen, hänelle, hänelä, hänelttä, hänessä,  
hänestä, hänet, häntä, ihan, ikään, ilman, itse, itseemme, itseeni, itseenne, itseensä, itseesi, itseksesi, itsellesi, it-  
selläsi, itsellään, itseltäsi, itseltään, itsen, itsenä, itsenäsi, itsenään, itsesi, itsessäsi, itsessään, itsestäsi, itsestään,  
itseä, itseäsi, itseään, ja, jk, jkhun, jkn, jkta, jo, johdosta, johon, johonkuhun, joiden, joidenkuiden, joihin, joihin-  
kuihin, joiksi, joiksikuksi, joilla, joillain, joillakuilla, joille, joillekuille, joilta, joiltain, joiltakuulta, join, joina,  
joinain, joinakuina, joine, joinekuine, joinkuin, joissa, joissain, joissakuissa, joista, joistain, joistakuista, joita,  
joitain, joitakuita, joitse, joitta, joittakuitta, joitten, joittenkuitten, jokainen, jokaiseen, jokaiseksi, jokaisella, jo-  
kaiselle, jokaiselta, jokaisen, jokaisena, jokaisessa, jokaisesta, jokaiset, jokaisetta, jokaisia, jokaisien, jokaisiin,  
jokaisiksi, jokaisilla, jokaisille, jokaisilta, jokaisin, jokaisina, jokaisine, jokaisissa, jokaisista, jokaisitta, jokaista,  
jokaisten, joksikuksi, joku, jollain, jollakulla, jollei, jollekulle, jollen, jollet, jolloin, joltain, joltakulta, jommaksi-  
kumaksi, jommallakummalla, jommallekummalle, jommatakummalta, jommankumman, jommassakummassa,  
jommastakumma, jommatakummat, jommatakummat, jommsikummi, jommilkummi, jommilekum-  
mille, jommiltakummi, jomminkummin, jommissakummissa, jommistikummi, jommittakummitta, jompaa-  
kumpaa, jompaankumpaan, jompanakumpana, jompiakumpia, jompienkumpi, jompinakumpina,  
jompinekumpine, jona, jonain, jonakuna, jonka, jonkin, jonkun, jos, joskus, jossain, jossakussa, jostain, josta-  
kusta, jota,  
jäljessä, jälkeen, jälleen, kai, kaikesi, kaikella, kaikelle, kaikelta, kaiken, kaikessa, kaikesta, kaiket, kaiketta, kai-  
kaksi, kaikilla, kaikille, kaikilta, kaikissa, kaikista, kaikitta, kaikkea, kaikkein, kaikkena, kaikki, kaikkia,  
kaikkiin, kaikkina, kaikkine, kanssa, kautta, keiden, keihin, keksi, keille, keillä, keiltä, keinä, keissä, keistä, kei-  
tä, keneen, keneksi, kenelle, kenellä, keneltä, kenen, kenenä, kenessä, kenestä, kenties, kerran, keskelle,

<sup>96</sup>Our list of Finnish stop words is based on <http://snowball.tartarus.org/algorithms/finnish/stop.txt>, with extensive additions to roughly match their counterparts in English.

keskellä, keskeltä, kesken, keskenään, keskeä, keski, keskitse, keskuudeksi, keskuudella, keskuudelle, keskuudelta, keskuuden, keskuudessa, keskuudesta, keskuudet, keskuudetta, keskuuksia, keskuuksien, keskuuksiin, keskuuksiksi, keskuuksilla, keskuuksille, keskuuksilta, keskuuksin, keskuuksina, keskuuksineen, keskuuksissa, keskuuksista, keskuuksitta, keskuus, keskuuteen, keskuutena, keskuutta, ketkä, ketä, kohdaksi, kohdalla, kohdalle, kohdalta, kohdan, kohdassa, kohdasta, kohdat, kohdatta, kohdiksi, kohdilla, kohdille, kohdilta, kohdin, kohdissa, kohdista, kohditta, kohta, kohtaa, kohtain, kohtana, kohti, kohtia, kohtiin, kohtina, kohtineen, koko, koska, kovin, kuhunkin, kuidenkin, kuihinkin, kuiksikin, kuillakin, kullekin, kultakin, kummaksi, kummalla, kummalle, kummalta, kumman, kummassa, kummasta, kummat, kummatta, kummiksi, kummilla, kummille, kummilta, kummin, kummissa, kummista, kummitta, kumpaa, kumpaan, kumpain, kumpana, kumpi, kumpia, kumpiin, kumpina, kumpine, kun, kunakin, kunnakin, kunes, kussakin, kustakin, kutakin, kuten, kutkin, kuttakin, kykene, kykenee, kykeneminen, kykenemistä, kykenemää, kykenemäisillään, kykenemän, kykenemättömien, kykenemättömiin, kykenemättömiksi, kykenemättömille, kykenemättömillä, kykenemättömiltä, kykenemättömin, kykenemättömine, kykenemättöminä, kykenemättömissä, kykenemättömistä, kykenemättömittä, kykenemättömiä, kykenemättömäksi, kykenemättömälle, kykenemättömällä, kykenemättömältä, kykenemättömän, kykenemättömänä, kykenemättömässä, kykenemättömästä, kykenemättömät, kykenemättömättä, kykenemättömään, kykenemätön, kykenemätönen, kykenemätöntä, kykenemää, kykenen, kykenet, kykenevien, kykeneviin, kykeneviksi, kykeneville, kykenevillä, kykeneviltä, kykenevin, kykenevine, kykenevinä, kykenevissä, kykenevistä, kykenevittä, kykeneviä, kykenevä, kykenevän, kykenevänä, kykenevää, kykenevänän, kykeni, kykenin, kykenisi, kykenisin, kykenisit, kykenit, kyllin, kyllä, lain, liene, lienee, lienen, lienet, liian, luo, luokse, luona, luota, lähekkäin, lähelle, lähellä, lähetä, lähes, lähetysten, lähi, lähitse, läpi, läsnä, lävitse, me, meidän, meidät, meihin, meille, meillä, meiltä, meissä, meistä, meitä, melko, mieluumin, mihin, miksi, mikä, mikäli, mikään, millainen, millaiseen, millaiseksi, millaisella, millaiselle, millaiselta, millaisen, millaisena, millaisessa, millaisesta, millaiset, millaisetta, millaisia, millaisien, millaisiin, millaisiksi, millaisilla, millaisille, millaisilta, millaisin, millaisina, millaisineen, millaisissa, millaisista, millaisitta, millaista, millaisten, mille, milloin, millä, millään, miltä, miltään, minkä, minkälainen, minkälaiseen, minkälaiseksi, minkälaisella, minkälaiselle, minkälaiselta, minkälaisen, minkälaisena, minkälaisessa, minkälaisesta, minkälaiset, minkälaisetta, minkälaisia, minkälaisien, minkälaisiin, minkälaisiksi, minkälaisilla, minkälaisille, minkälaisilta, minkälaisin, minkälaisina, minkälaisineen, minkälaisissa, minkälaisista, minkälaisitta, minkäläista, minkäläisten, minkään, minne, minua, minulla, minulta, minut, minun, minussa, minusta, minut, minuun, minä, minään, missä, missään, mistä, mistään, miten, mitkä, mitkään, mitä, mitään, moinen, moiseen, moiseksi, moisella, moiselle, moiselta, moisen, moisena, moisessa, moisesta, moiset, moisetta, moisia, moisi, moisiin, moiski, moisilla, moisille, moisilta, moisin, moisina, moisine, moissa, moisista, moisitta, moista, moisten, molemmaksi, molemmaalla, molemmalle, molemmalta, molemman, molemmaassa, molemmasta, molemmat, molemmatta, molemmiksi, molemmilla, molemmille, molemmilta, molemin, molemmissa, molemmista, molemitta, molempaa, molempaan, molempana, molempi, molempia, molempiin, molempina, molempine, moneen, moneksi, monella, monelle, monelta, monen, monena, monessa, monesta, monesti, monet, monetta, moni, monia, moniaalla, moniaalle, moniaalta, moniin, monin, monina, monine, monna, monta, muiden, muihin, muiksi, muilla, muille, multa, muin, muina, muine, muissa, muista, muita, muitta, mutitten, mukaan, mukana, mutta, mutte, muttei, mutten, muttet, muu, muualla, muualle, muualta, muuhun, muulloin, muun, muuna, muut, muuta, muutamaa, muutamain, muutaman, muutamana, muutamat, muutamia, muutamien, muutamiin, muutamiksi, muutamilla, muutamille, muutamilta, muutamin, muutamina, muutamineen, muutamissa, muutamista, muutamitta, myös, myötä, ne, niiden, niihin, niiksi, niille, niillä, niiltä, niin, niine, niinkuin, niinä, niissä, niistä, niitten, niittä, niitä, no, noiden, noihin, noksi, noilla, noille, noilta, noin, noina, noissa, noista, noita, nuo, nyt, näiden, näihin, näiksi, näille, näillä, näiltä, näin, näine, näinä, näissä, näistä, näittä, näitä, nämä, ohi, ohitse, oikein, oitis, ole, olema, olemaa, olemain, olemaisillaan, oleman, olemana, olemat, olematon, olematonta, olemattonen, ole mattomaan, ole mattomaksi, ole mattomalla, ole mattomalle, ole mattomalta, ole mattoman, ole mattomana, ole mattomassa, ole mattomasta, ole mattomat, ole mattomatta, ole mattomia, ole mattomien, ole mattomiin, ole mattomiksi, ole mattomilla, ole mattomille, ole mattomilta, ole mattomin, ole mattomina, ole mattomine, ole mattomissa, ole mattomista, ole mattomitta, ole mia, ole mien, ole miin, ole miksi, ole milla, ole mille, ole milta, ole min, ole mina, ole mineen, ole missa, ole mista, ole mittä, olen, olet, oleva, olevaa, olevain, olevan, olevana, olevasti, oli, olin, olisi, olisin, olisit, olit, olkaa, olko, olkoon, olkoot, olla, olleksi, ollella, olleelle, ollelta, olleen, olleena, olleeseen, olleessa, olleesta, olleet, olleetta, olleiden, olleihin, olleiksi, olleilla, olleille, ollelta, ollein, olleina, ollein, olleistin, olleissa, olleista, olleita, olleitta, olleitten, ollen, ollessa, olluksi, olluilla, olluille, olluulta, olluin, olluissa, olluista, olluitta, olluksi, ollulla, ollulle, ollulta, ollun, ollussa, ollusta, ollut, ollutta, oltaessa, oltaisi, oltaisiin, oltako, oltakoon, oltaman, oltane, oltava, oltiin, oltu, oltua, oltuihin, oltuina, oltuine, oltuja, oltujen, oltuna, oltuun, oma, omaa, omain, oman, omana, omat, omia, omien, omiin, omiksi, omilla, omille, omilta, omin, omina, omine, omissa, omista, omitta, on, ovat, paitsi, paljo, paljoa, paljoihin, paljoiksi, paljoilla, paljoilta, paljoin,

paljoina, paljoineen, paljoissa, paljoista, paljoitta, paljoja, paljojen, paljon, paljona, paljoon, paljot, pelkaksi, pelkille, pelkillä, pelkiltä, pelkin, pelkissä, pelkistä, pelkittä, pelkkien, pelkkine, pelkkinä, pelkkiä, pelkkä, pelkkänä, pelkkää, pelkkään, pelkäksi, pelkälle, pelkällä, pelkältä, pelkän, pelkässä, pelkästä, pelkät, pelkättä, perässä, pian, pitkin, poikki, pois, poissa, pysty, pystyessä, pystyi, pystyin, pystysi, pystyisin, pystyisit, pystyit, pystykää, pystyköön, pystyköt, pystyminen, pystymistä, pystymä, pystymäisillään, pystymän, pystymätön, pystymään, pystyn, pystyne, pystyneet, pystyneenä, pystyneeseen, pystyneet, pystyneiden, pystyneihin, pystyneiksi, pystyneille, pystyneillä, pystyneiltä, pystynein, pystyneine, pystyneinä, pystyneisiin, pystyneissä, pystyneistä, pystyneitten, pystyneittä, pystyneit, pystynet, pystynyttä, pystyt, pystyttiin, pystytty, pystytäessä, pystytäisi, pystytäisiin, pystytäkö, pystytäköön, pystytämän, pystytäne, pystytävä, pystytä, pystytään, pystyvä, pystyy, pystyä, pän, rinnalla, saa, saada, saaden, saadessa, saakaa, saakka, saakoon, saakoot, saama, saamaisillaan, saaman, saamaton, saamisen, saamista, saan, saane, saaneet, saanet, saanut, saat, saataessa, saataisi, saataisiin, saatako, saatakoon, saataman, saatane, saatava, saatin, saatu, saava, sai, sain, sasis, sasis, saisit, sait, sama, samaa, saman, samana, samat, samoithin, samoiksi, samoilla, samoilta, samoin, samoina, samoine, samoissa, samoista, samoitta, samoja, samojen, sangen, se, sellainen, sellaiseen, sellaiseksi, sellaisella, sellaiselle, sellaiselta, sellaisen, sellaisessa, sellaisesta, sellaiset, sellaisetta, sellaisia, sellaisien, sellaisiin, sellaisiksi, sellaisilla, sellaisille, sellaisilta, sellaisin, sellaisina, sellaisineen, sellaisissa, sellaisista, sellaisitta, sellaista, sellaisten, semmoinen, semmoiseen, semmoiseksi, semmoisella, semmoiselle, semmoiselta, semmoisen, semmoisena, semmoisessa, semmoiset, semmoisetta, semmoisia, semmoisien, semmoisiin, semmoisiksi, semmoisilla, semmoisille, semmoisilta, semmoisin, semmoisina, semmoisine, semmoissä, semmoisista, semmoisitta, semmoista, semmoisten, sen, sentään, senvuoksi, siellä, sieltä, siihen, stinä, siis, siitä, sijaan, siksi, sille, silloin, sillä, siltä, sinne, sinua, sinulla, sinulta, sinulta, sinun, sinussa, sinusta, sinut, sinuun, sinä, sisä, sisäkkäin, sisätysten, sisään, siten, sitten, sitä, taa, taakse, taanoin, taas, tahansa, tahi, tai, taitse, taka, takaa, takaisin, takana, takia, te, tee, teet, tehdent, tehdessä, tehdä, tehdään, tehkää, tehkö, tehköön, tehköt, tehne, tehnee, tehneet, tehnent, tehnyt, tehtiin, tehty, tehtäessä, tehtäisi, tehtäisiin, tehtäkö, tehtäköön, tehtämän, tehtäne, tehtävä, teidän, teidät, teihin, teille, teillä, teiltä, teimme, tein, teissä, teistä, teit, teitte, teitä, tekee, tekeminen, tekemistä, tekemä, tekemäillään, tekemän, tekemätön, tekemään, tekevä, tekevät, teki, tekisi, tekisin, tekisit, todella, toinen, toisaalla, toisaalle, toisaalta, toiseen, toiseksi, toisella, toiselle, toiselta, toisen, toisena, toisensa, toisessa, toisesta, toiset, toisetta, toisia, toisien, toisiensa, toisiin, toisiinsa, toisikseen, toisiksen, toisiksi, toisilla, toisilta, toisin, toisine, toisinsa, toisista, toisitta, toista, toisten, toistensa, tuo, tuohon, tuollainen, tuollaiseen, tuollaiseksi, tuollaisella, tuollaiselle, tuollaiselta, tuollaisen, tuollaisena, tuollaisessa, tuollaisesta, tuollaiset, tuollaisetta, tuollaisia, tuollaisien, tuollaisiin, tuollaisiksi, tuollaisilta, tuollaisille, tuollaisilta, tuollaisin, tuollaisina, tuollaisineen, tuollaisissa, tuollaisista, tuollaisitta, tuollaista, tuollaisten, tuon, tuona, tuota, tuskin, tähden, tähän, täksi, tällainen, tällaiseen, tällaiseksi, tällaisella, tällaiselle, tällaiselta, tällaisen, tällaisessa, tällaisesta, tällaiset, tällaisetta, tällaisia, tällaisien, tällaisiin, tällaisksi, tällaisilla, tällaisille, tällaisilta, tällaisin, tällaisina, tällaisineen, tällaisissa, tällaisista, tällaisitta, tällaista, tällaisten, tälle, tällä, tällön, tältä, tämä, tämän, täinne, täänä, tässä, tästä, täten, täitä, tädy, täytyi, täytyisi, täytykö, täytyköön, täytyminen, täytymistä, täytyne, täytyne, täytynyt, täytyy, täytyä, täällä, täältä, ulkoa, ulkona, ulos, usea, useaa, useain, useain, useammaksi, useammalla, useammalle, useammalta, useamman, useammassa, useammas-ta, useammat, useammatta, useammiksi, useammilla, useammille, useammlta, useammin, useammissa, useam-mista, useammitta, useampaan, useampain, useampaan, useampi, useampia, useampiin, useampina, useampine, usean, useana, useat, useata, useiden, useihin, useiksi, useilla, useilta, useimmaksi, useim-malla, useimmalta, useimman, useimmassa, useimmasta, useimmat, useimmatta, useimmiksi, useimmilla, useim-mille, useimnilta, useimmin, useimmissa, useimmista, useimmitta, useimpaan, useimpain, useimpana, useimpia, useimpien, useimpiin, useimpina, useimpine, usein, useina, useine, useinta, useinten, useisiin, useissa, useista, useita, useitta, useitten, uudelleen, udestaan, vaan, vai, vaikka, vaikkei, vaikkemme, vaikken, vaikket, vaikkette, vain, vallan, varsin, vasta, vastedes, vasten, vastikään, vielä, viereen, vierekäin, vierellä, viereltä, vieren, vieressä, vierestä, vieret, vieretysten, vieri, vieriin, vieritse, vieriä, viertä, voi, voiden, voidessa, voikaa, voikoon, voikoot, voima, voimaisillaan, voiman, voimaton, voiminen, voimista, vain, voine, voinee, voineet, voinet, voinut, voisi, voisit, voit, voitaessa, voitaisi, voitaisiin, voitako, voitakoon, voitaman, voitane, voitava, voitiin, voitu, voiva, vuoksi, väliin, välillä, välissä, välistä, yhdeksi, yhdelle, yhdellä, yhdeltä, yhden, yhdessä, yhdestä, yhdestä, yhdet, yhdettä, yhteen, yhtenä, yhtä, yhtälle, yhtälä, yhtälä, yhtää, yhtää, yhää, yksi, yksin, yksin, yksine, yksinä, yksittäin, yksiiä, ylen, yli, ylitse, yläpuolella, ylös, ympäri, ympäriiinsä, ympärítse, älkää, älköön, älkööt, ällös, älä,

then we consider it a Finnish stop root (notation: **FinnishStopRoot**). A Finnish stop word matches the following string pattern:

**FinnishStopRoot**( $\emptyset|ksilla|lle|llä|ltä|ssa|ssä|sta|stää|tta|ttä|$ )( $\emptyset|kse|$ )  
 $(\emptyset|an|en|mme|ni|nne|nsa|nsä|si|tte|vat|vät|)$ ( $\emptyset|han|hän|ka|kaan|kin|ko|kä|kää|kō|pa|pä|$ )<sub>m</sub>,

or (*jok|vast*) **X**.<sup>97</sup> All the Finnish stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list.  $\square$

Up to some modifications, our method is modeled after the Porter stemming algorithm for Finnish (<http://snowball.tartarus.org/algorithms/finnish/stemmer.html>), devised by Martin F. Porter. Our modifications of the Porter stemming algorithm aim at two goals: to fully take care of verb conjugations and to fully accommodate to consonant/vowel alternations in declensions.

### 8.1.1 Effective spelling and essential root

It is assumed that all Finnish words are converted to lowercase before going through any of the procedures below.

Unlike the situation in German, diacritical marks in Finnish can affect the meaning of a word immensely. For example, we have *rohkaista* “to encourage” vs. *röhkäistä* “to grunt”. Therefore, we are going to retain diacritical marks in the construction of Finnish effective spellings and essential roots, contrary to the practice for German.

Before we start, we recapitulate Martin F. Porter’s definition of Finnish vowels, restricted vowels and long vowels.

**Definition 8.2** (Finnish Vowels, Restricted Vowels and Long Vowels). Hereafter in §8.1, the symbol **V** stands for any member from the list of Finnish vowels {*a*, *e*, *i*, *o*, *u*, *y*, *ä*, *ö*}. In line with the multiplicity notations introduced in Definition 3.3, the symbol **V<sub>m</sub>** stands for a text string formed by consecutive appearance of one or more (not necessarily identical or distinct) members from the set of Finnish vowels.

Dual to the notations above, the symbol **C** stands for any character that does not belong to the list {*a*, *e*, *i*, *o*, *u*, *y*, *ä*, *ö*}, and **C<sub>m0</sub>** stands for a text string formed by consecutive appearance of zero or more (not necessarily identical or distinct) characters that do not belong to the same list.

The symbol **V\*** stands for any member from the list of Finnish restricted vowels {*a*, *e*, *i*, *o*, *u*, *ä*, *ö*}, while **V\*\*** represents any one in the list of Finnish long vowels {*aa*, *ee*, *ii*, *oo*, *uu*, *ää*, *öö*}.  $\square$

In the following, we define the first and second protected ranges of a Finnish word, following Martin F. Porter’s definition of regions *R1* and *R2*.

**Definition 8.3** (Finnish First protected range). Let  $\hat{\sigma}$  be a Finnish text string, then its first protected range  $\text{ProtRg}_1(\hat{\sigma})$  is specified through the following procedures:

- Look for the string pattern  $\mathbf{C}_{m_0} \mathbf{V}_m \mathbf{C} \sim$  in the string  $\hat{\sigma}$ ;
- If the string pattern above is found, the last position occupied by such a string defines  $\text{ProtRg}_1(\hat{\sigma})$ ; otherwise, set  $\text{ProtRg}_1(\hat{\sigma}) = \ell(\hat{\sigma})$ .

Once  $\text{ProtRg}_1(\hat{\sigma})$  is determined, the text string  $\hat{\sigma}$  needs to be explicitly marked with a separator signifying the position of the first protected range. Concretely speaking, we define  $\text{MarkRg}_1(\hat{\sigma})$  by inserting a dagger symbol (“†”) at position  $\text{ProtRg}_1(\hat{\sigma}) + 1$  to the string  $\hat{\sigma}$ .  $\square$

**Definition 8.4** (Finnish Second protected range). Let  $\hat{\sigma}'$  be a Finnish text string marked explicitly with its first protected range, then its second protected range  $\text{ProtRg}_2(\hat{\sigma}')$  is determined by the following procedures:

- Look for the string pattern  $\dagger \mathbf{C}_{m_0} \mathbf{V}_m \mathbf{C}$  in the string  $\hat{\sigma}'$ ;
- If the string pattern above is found, the last position occupied by such a string defines  $\text{ProtRg}_2(\hat{\sigma}')$ ; otherwise, set  $\text{ProtRg}_2(\hat{\sigma}') = \ell(\hat{\sigma}')$ .

Once  $\text{ProtRg}_2(\hat{\sigma}')$  is determined, the text string  $\hat{\sigma}'$  needs to be further marked with a separator signifying the position of the second protected range. Concretely speaking, we define  $\text{MarkRg}_2(\hat{\sigma}')$  by inserting a double dagger symbol (“‡”) at position  $\text{ProtRg}_2(\hat{\sigma}') + 1$  to the string  $\hat{\sigma}'$ .  $\square$

*Example 8.4.1.* We use some English words to illustrate the workings of  $\text{ProtRg}_1$ ,  $\text{MarkRg}_1$ ,  $\text{ProtRg}_2$  and  $\text{MarkRg}_2$ :

$\hat{\sigma}$	<i>beau</i>	<i>beauty</i>	<i>beautiful</i>	<i>beautifully</i>
$\text{ProtRg}_1(\hat{\sigma})$	4	5	5	5
$\text{MarkRg}_1(\hat{\sigma})$	<i>beau†</i>	<i>beaut†y</i>	<i>beaut†iful</i>	<i>beaut†ifully</i>
$\text{ProtRg}_2(\text{MarkRg}_1(\hat{\sigma}))$	5	7	8	8
$\text{MarkRg}_2(\text{MarkRg}_1(\hat{\sigma}))$	<i>beau†‡</i>	<i>beaut†y‡</i>	<i>beaut†if‡ul</i>	<i>beaut†if‡ully</i>

<sup>97</sup>Here, the apostrophe refers to either APOSTROPHE **U+0027** or RIGHT SINGLE QUOTATION MARK **U+2019**, as is appropriate for the particular text mining task.

**Algorithm 8.5** (Finnish effective spelling). For a Finnish word  $\hat{\sigma}$  (converted to lower case form), its effective spelling  $\text{EffSpell}(\hat{\sigma})$  is constructed in sequential steps:

- (1) Convert to lowercase and do  $\sim \text{ksest}(a|\ddot{a}) \rightarrow s$ ,  $\sharp \mathbf{X} \rightarrow \emptyset$ .
- (2) Do  $\text{iudelleen} \rightsquigarrow \rho e$ .
- (3) Do  $\text{nt}(o|\ddot{o})(i|j)(\emptyset|h)i(\emptyset|n)(\emptyset|a|en|\ddot{a}) \rightarrow nta$ ,  $\sim \hat{\chi}_1 \text{eeseen} \mathbf{X} \rightarrow \hat{\chi}_1 e$ ,  $\sim \hat{\chi}_2 m(a|\ddot{a}) \text{sill} \mathbf{X} \rightarrow \hat{\chi}_2$ ,  $\sim \hat{\chi}_3 \text{sest}(a|\ddot{a}) \mathbf{X} \rightarrow \hat{\chi}_3 \text{nen}$ ,  $\mathbf{X} \hat{\chi}_m(a|\ddot{a})(t|tt)(o|\ddot{o})(m|n)(\emptyset|a|aa|i|t|\ddot{a}|\ddot{a})(\emptyset|en|iin|ksi|ll|l|n|ss|st|t|tt)(\emptyset|a|e|\ddot{a}) \rightarrow \tau v\tau - \mathbf{X} \hat{\chi}$ .
- (4) Replace

$(sisar velj)enpo(i j)\mathbf{X} \sim$		$(sisar velj)entyt\mathbf{X} \sim$		$amerika$		$avioliit(\emptyset t) \sim$		$etei(n s) \sim$	
$\tilde{v}efw$		$\tilde{v}iece$		$aμerika$		$μarp$		$aiλs$	
help~	henkilö	$iham\mathbf{X} \sim$	ihmi	isänne	john~	$kai(\emptyset t)o$	$kai(\emptyset k)u\mathbf{X} \sim$	$kaika(a isi)\mathbf{X} \sim$	kalle~
ηελρ	περσօ	ωο̄δρε	inhimi	isä	γιω̄ν	λοσσο	εχօ	εχօ	expnse
kalli	kasvo	kats $\mathbf{X}$	kirjasto	koir	kuol~	kylj~	kä(d t)e~	käyt(∅ t)(e i ō)(∅ j)	luopu
expnsi	φaso	kato $\sigma$	λιβρօ	dog	διελ	kyly	käsi	υσā	luovu
läpikotaisin	maine	musta(ank sukk)~	naimis~	ou(d t)o~	pilot~	pool~	puh(e u) $\mathbf{X}$	puista	puolu
θury	fame	jaλωσ	μarp	weirdo	πιλοτ	πωολ	λογγ	ξακα	defu
puolust	päiv	selv~	seuralai(n s)	silm	sisaru	sito~	sulo~	summa	syntyperā
defens	πäiv	celar	σοclais	εγε	σιβλυ	biñdo	μελτα	ξармо	σιμта syntyenperā
taho~	tall	tarkoi~	tarkoi~	tullut $\mathbf{X} \sim$	tult	tun(n t)o	tunnust	tuntei(s t) $\mathbf{X}$	tuntisi~
auθo	taλll	μεαν්i	μεαν්i	tuli	tul	tuntea	confεβ	tuntea	tunteasi
usko~	uu(de si)m(m p) $\mathbf{X} \sim$		vaikut(∅ t) $\mathbf{X}$		val(l t) $\mathbf{X}^{\epsilon}(a o)$		valli~	viha	vilp(∅ p)~
φaθo	uusi		infλ		γawlt $\mathbf{X}$		δομνi	ανγρα	fardp
vuo(si teen tena tt)		yh(d t)~	miss		nai(∅ si va)(∅ n t)(∅ e ut)(∅ e et mme tte vat)				
jahpoi		yks	neiti			μarp			
uusitte	vuonna	~aikoi	~käyt(iin te)		~käyttää(d t) $\mathbf{X}$		~vuode(ksi lla lle lta n ssa sta t)		
uusi	jahp	aio	käy	baχs			jahp		

- (5) Replace

$(mamm muts äip äipp äisk) \sim$		$\mathbf{X}^{\epsilon}(V s)tapa$		$\mathbf{X}^{\epsilon}(\emptyset lähi)omai(n s) \sim$	
$\tilde{a}it$		Xtava		Xpelτβa	
$\mathbf{X}^{\epsilon}(\emptyset oiko)lu(e i) \sim$	ahdist	aiko~	alu	asem~	asu
Xluki	ayiσt	aio	alku	σοτωεμ	δερσυ
har(r t)~	harm~	harv~	heng~	huone	ihm~
δowt	hram	par	henk	ζιμε	ηημ
juo(ks sk ss)~	kahv	kakist~	katker~	kaun(a oi oj)~	kes(i ā)~
juost	kofφ	kof	βitter	χατα	εετεα
kitty~	ko(d t)(e i)	kohtalo	kuul~	kuun(n t)el~	käyttä
κιττу	koti	φατο	kuul	kuul	υσā
loist~	lom~	luist~	luk(∅ k)o	luki(nn t ts tt)	luku
σiŋiň	ηολδ	σiλδ	λοκο	λοκ	ταζπι
mail(e ei ej a i)~		maini(n nn t ts tt)~		mary~	matk
μαιλi		μεντζion		μαργ	τεριπ
nei(d k m t)~	nei(d t)o	oiv(a i)~	onnist	ost~	parv
μarp	μαιδο	φινα	σukc	βυη	gupp
pisi(mm mp n nt)~	pitk $\mathbf{X}^{\epsilon}(i ā) \sim$	pitu~	puh(d t)	puisto	päivälli~
λοῦγ	λοῦγX	λοῦγу	κελανd	πарко	διεπlli

<i>sal(e eɪ eja i)</i>	<i>seinsem~</i>	<i>selk~</i>	<i>seurustel</i>	<i>sie(d t)~</i>	<i>siev(i ä)~</i>	<i>suvai(n nn nt ts it)</i>
<i>ζaaλi</i>	<i>7em</i>	<i>sel</i>	<i>σocl</i>	<i>σotd</i>	<i>περττi</i>	<i>τολρατ</i>
<i>säily~</i>	<i>sääst~</i>	<i>tahd</i>	<i>talou</i>	<i>tapa~</i>	<i>tapp~</i>	<i>tee~</i>
<i>räμv̄y</i>	<i>σavat</i>	<i>taht</i>	<i>εekou</i>	<i>tava</i>	<i>κiλ</i>	<i>ξai</i>
<i>tuha(n nn nt t tt)</i>	<i>tulkχ̄(a o)~</i>	<i>tun(n t)(ei en ev u)(Ø u)</i>	<i>tunti(n s)</i>	<i>turh~</i>	<i>tyt(Ø t)är~</i>	
<i>θat</i>	<i>τuλkkχ̄</i>	<i>tuntea</i>		<i>τuntis</i>	<i>ÿiλh</i>	<i>δoχτr</i>
<i>tyt(Ø t)ö~</i>	<i>tädi~</i>	<i>tänään~</i>	<i>vapi~</i>	<i>varh</i>	<i>varm</i>	<i>veli~</i>
<i>γipλö</i>	<i>täti</i>	<i>τοδηjan</i>	<i>vavi</i>	<i>εarλh</i>	<i>κuσρ</i>	<i>deβta</i>
<i>aie(Ø tta)</i>	<i>eilen</i>	<i>enoX</i>	<i>huomen(issa na)</i>	<i>leskirouvX</i>	<i>nainen</i>	
<i>aike</i>	<i>eilinen</i>	<i>ÿoklo</i>	<i>τoμρoω</i>	<i>leski</i>	<i>naisen</i>	
<i>~aik(a aa oi ojen)(Ø in n na neen)</i>	<i>~kaunis</i>	<i>~käytä(Ø mme tte)</i>	<i>~käytös(Ø ten)</i>	<i>~varu(ks s st)X</i>	<i>~veli</i>	
<i>aja</i>	<i>kaunii</i>	<i>uσä</i>	<i>baχs</i>	<i>kiβmt</i>	<i>veljen</i>	

## (6) Replace

<i>(hyv parah paras parh)~</i>	<i>V'</i>	<i>X<sup>ε</sup>(yö öiC)~</i>	<i>é</i>	<i>hdeks</i>	<i>herist~</i>	<i>hert~</i>	<i>herä~</i>	<i>huom~</i>	<i>ilo~</i>	<i>inh~</i>
<i>parempi</i>	<i>VQ</i>	<i>qX</i>	<i>e</i>	<i>hdex</i>	<i>ξak</i>	<i>σewt</i>	<i>wakä</i>	<i>σim</i>	<i>ilzo</i>	<i>inzh</i>
<i>is(i ä)~</i>	<i>jano~</i>	<i>juo(k s)C~</i>	<i>kirje~</i>	<i>kuum~</i>	<i>lm</i>	<i>miel</i>	<i>sall~</i>	<i>sellatX<sup>ε</sup>(n s)~</i>	<i>siro~</i>	<i>siskoX~</i>
<i>faija</i>	<i>γiano</i>	<i>qojo</i>	<i>girje</i>	<i>guum</i>	<i>λμ</i>	<i>μiel</i>	<i>sallq</i>	<i>zenlaiX</i>	<i>σiro</i>	<i>sisar</i>
<i>toiv~</i>	<i>tun(n t)(ej i)χ̄(e s)</i>	<i>tunne(i r)</i>	<i>tunnis(s t)a</i>	<i>tunnei(hin na neen)</i>		<i>tuom(Ø ar)i~</i>				
<i>χoπ</i>	<i>τuntiχ̄</i>	<i>τunti</i>	<i>τunti</i>	<i>τunti</i>		<i>τunti</i>			<i>δoμi</i>	
<i>vaj~</i>	<i>ystäv(Ø yks yst)(Ø i)~</i>	<i>hra</i>	<i>nti</i>	<i>rva</i>	<i>~(Ø si)tte</i>		<i>~X<sup>ε</sup>(χ̄<sub>m</sub>)zχ̄<sub>m0</sub></i>	<i>~C'X</i>		
<i>waj</i>	<i>ystew</i>	<i>herra</i>	<i>neiti</i>	<i>rouva</i>	<i>Ø</i>		<i>χz</i>	<i>C</i>		
<i>~X<sup>ε</sup>(l s)k(o ö)</i>	<i>X</i>	<i>~X<sup>ε</sup>(u y)(ksin s tta ttä)</i>	<i>~ej(a ä)</i>	<i>~et</i>	<i>~isk(aa ää)</i>	<i>~lleen</i>	<i>~mies</i>	<i>~tunti</i>		
	<i>Xd</i>		<i>e</i>	<i>etkin</i>	<i>isk</i>	<i>ll</i>	<i>mieh</i>	<i>tunti</i>		
<i>(hyv parah paras parh)~</i>	<i>V'</i>	<i>X<sup>ε</sup>(yö öiC)~</i>	<i>é</i>	<i>hdeks</i>	<i>herist~</i>	<i>hert~</i>	<i>herä~</i>	<i>huom~</i>	<i>ilo~</i>	<i>inh~</i>
<i>parempi</i>	<i>VQ</i>	<i>qX</i>	<i>e</i>	<i>hdex</i>	<i>ξak</i>	<i>σewt</i>	<i>wakä</i>	<i>σim</i>	<i>ilzo</i>	<i>inzh</i>
<i>is(i ä)~</i>	<i>jano~</i>	<i>juo(k s)C~</i>	<i>kirje~</i>	<i>kuum~</i>	<i>lm</i>	<i>miel</i>	<i>sall~</i>	<i>sellatX<sup>ε</sup>(n s)~</i>	<i>siro~</i>	<i>siskoX~</i>
<i>faija</i>	<i>γiano</i>	<i>qojo</i>	<i>girje</i>	<i>guum</i>	<i>λμ</i>	<i>μiel</i>	<i>sallq</i>	<i>zenlaiX</i>	<i>σiro</i>	<i>sisar</i>
<i>toiv~</i>	<i>tun(n t)(ej i)χ̄(e s)</i>	<i>tunne(i r)</i>	<i>tunnis(s t)a</i>	<i>tunnei(hin na neen)</i>		<i>tuom(Ø ar)i~</i>				
<i>χoπ</i>	<i>τuntiχ̄</i>	<i>τunti</i>	<i>τunti</i>	<i>τunti</i>		<i>τunti</i>			<i>δoμi</i>	
<i>vaj~</i>	<i>ystäv(Ø yks yst)(Ø i)~</i>	<i>hra</i>	<i>nti</i>	<i>rva</i>	<i>~(Ø si)tte</i>		<i>~X<sup>ε</sup>(χ̄<sub>m</sub>)zχ̄<sub>m0</sub></i>	<i>~C'X</i>		
<i>waj</i>	<i>ystew</i>	<i>herra</i>	<i>neiti</i>	<i>rouva</i>	<i>Ø</i>		<i>Xz</i>	<i>C</i>		
<i>~X<sup>ε</sup>(l s)k(o ö)</i>	<i>X</i>	<i>~X<sup>ε</sup>(u y)(ksin s tta ttä)</i>	<i>~ej(a ä)</i>	<i>~et</i>	<i>~isk(aa ää)</i>	<i>~lleen</i>	<i>~mies</i>	<i>~tunti</i>		
	<i>Xd</i>		<i>e</i>	<i>etkin</i>	<i>isk</i>	<i>ll</i>	<i>mieh</i>	<i>tunti</i>		

(7) If the result so far start with **C<sub>m</sub>VV~** (where the two occurrences of **V** may or may not represent the same vowel), perform the following replacements on the word initial pattern **C<sub>m</sub>VV~**:<sup>98</sup>

<i>ξ̄ε(V)i</i>	<i>V**</i>	<i>ie</i>	<i>uo</i>	<i>yö</i>
<i>ξ̄ijξi</i>	<i>(V**)⁽¹⁾iji</i>	<i>eije</i>	<i>oijo</i>	<i>öijö</i>

If the result after step (6) does not start with **C<sub>m</sub>VV~**, leave it as is.

## (8) Mark the first protected range with †.

## (9) Replace

<i>~X<sup>ε</sup>((a e i n o t u y ä ö)(Ø †))(han hän kaan kääñ kin ko kö pa pää)m</i>	<i>X</i>
--	----------

<sup>98</sup>Here, the two instances of *ξ̄* in *ξ̄ijξi* represent the same letter; *(V\*\*)⁽¹⁾iji* stands the first letter in the long vowel **V\*\***, followed immediately by *iji*.

$$\sim(\emptyset|s)(iess(a|i\ddot{a})k(aa|\ddot{a}\ddot{a})(\emptyset|mme)|k(oo|\ddot{o}\ddot{o})(n|t)|minen|nee(\emptyset|t)|n(e|u|y)t|si(mme|n|t)tt(aessa|\ddot{a}ess\ddot{a}|\ddot{a}isi|aman|\ddot{a}m\ddot{a}n|ane|\ddot{a}ne|aneen|\ddot{a}neen|iin|u|y)v(a|\ddot{a}|at|\ddot{a}t))$$

$$\begin{array}{ccc} \sim\dagger\mathbf{X}(l|v)(a|\ddot{a})i(n|s)\hat{\chi}_{m_0} & \sim\dagger\mathbf{X}hk\hat{\chi}_{m_0} & \sim\dagger\mathbf{X}lli(n|s)\hat{\chi}_{m_0} \\ \dagger\mathbf{X} & \dagger\mathbf{X} & \dagger\mathbf{X} \\ \sim\dagger\mathbf{X}luont(a|e)\hat{\chi}_{m_0} & \sim\dagger\mathbf{X}suh(d|t)\hat{\chi}_{m_0} & \sim\dagger\mathbf{X}t(a|\ddot{a})r\hat{\chi}_{m_0} \\ \dagger\mathbf{X} & \dagger\mathbf{X} & \dagger\mathbf{X} \end{array}$$

(10) *Mark the second protected range with  $\ddot{\cdot}$ .*

(11) *Do  $\sim\dagger\mathbf{X}sti \rightarrow \ddot{\mathbf{X}}$ ,  $\sim\dagger\mathbf{X}'(ks|ll|lt|ss|st)\mathbf{X}'' \rightarrow \ddot{\mathbf{X}}'$ ,  $(k\ddot{\cdot}s|l\ddot{\cdot}l|l\ddot{\cdot}t|s\ddot{\cdot}s|s\ddot{\cdot}t)\mathbf{X}''' \rightarrow \emptyset$ .*

(12) *Do  $\ddot{\cdot} \rightarrow \emptyset$ .*

(13) *Replace*

$$\begin{array}{ccccc} \sim\dagger it & \sim\dagger\mathbf{XC}it & \sim\dagger\mathbf{X}(mme|nne|nsa|ns\ddot{a}) & \sim\dagger\mathbf{X}_{\hat{\chi}\notin}(k)si & \sim\dagger\mathbf{X}^{X_1\in}(\emptyset|kse)ni \\ i & \dagger\mathbf{XC}i & \dagger\mathbf{X} & \dagger\mathbf{X}\hat{\chi} & \sim\dagger\mathbf{XX}'_1 \\ \sim\mathbf{X}_2^{\in}((\dagger\mathbf{X}(in|ll))|i\dagger n|l\dagger t)een & & \sim\mathbf{X}_3^{\in}((\dagger\mathbf{X}(ll|lt|n|ss|st|t))|l\dagger l|l\dagger t|n\dagger|s\dagger s|s\dagger t|t\dagger) & \mathbf{X}_3\mathbf{X}_4^{\{1\}} & \sim\mathbf{X}_4^{\in}(aa|\ddot{a}\ddot{a})n \\ \mathbf{X}_2e & & & & \end{array}$$

*where one obtains  $\mathbf{X}'_1$  after doing  $kse \rightarrow ksi$  on  $\mathbf{X}_1$ .*

(14) *Do  $\dagger \rightarrow \emptyset$ .*

(15) *Do  $mai\sim \rightarrow eap\theta$ ,  $mm\hat{\chi}^{\in}(a|i|\ddot{a})(\emptyset|n) \rightarrow M\hat{\chi}$ ,  $sisar\sim \rightarrow \sigma i\sigma a\rho$ ,  $\tau v\tau-j\mathbf{V}_m \rightarrow \tau v\tau-Z\mathbf{V}_m$ ,  $\tau v\tau-käy \rightarrow \tau v\tau-kkäy$ ,  $\sim nt(a|\ddot{a}) \rightarrow nt$ ,  $\sim x(a|\ddot{a})s \rightarrow x$ ,  $kahdexi \rightarrow kaksi$ .*

**Algorithm 8.6** (Finnish essential root). *Let  $\hat{\sigma}$  be the effective spelling of a Finnish word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:*

(1) *Mark the first protected range with  $\dagger$ .*

(2) *Do  $\sim\dagger\mathbf{X}\hat{\chi}h\hat{\xi}n \rightarrow \dagger\mathbf{X}\hat{\xi}^*$  if  $\hat{\xi} \in \{a, e, i, o, u, \ddot{a}, \ddot{o}\}$ . Make the following replacements in the same sweep:*

$$\begin{array}{cccc} \sim\dagger\mathbf{X}\mathbf{V}^*i(den|siin|tten) & \sim\mathbf{X}_1^{\in}(\dagger\mathbf{X}(\mathbf{V}^{**}|(a|\ddot{a})k))seen & \sim\mathbf{X}_2^{\in}(\dagger\mathbf{X}\mathbf{C}\mathbf{V}|\mathbf{C}\dagger\mathbf{V})(\emptyset|ej)(a|\ddot{a}) & \sim\dagger\mathbf{X}\mathbf{V}tt(a|\ddot{a}) \\ \dagger\mathbf{X}\mathbf{V}^*i^* & \mathbf{X}_1^* & \mathbf{X}_2^* & \dagger\mathbf{X}\mathbf{V}^* \\ \sim\dagger\mathbf{X}(ine|ksi|ll(a|\ddot{a}|e)|lt(a|\ddot{a})|n(a|\ddot{a})|ss(a|\ddot{a})|st(a|\ddot{a})) & & & \sim\mathbf{X}_3^{\in}(\dagger\mathbf{X})n \\ \dagger\mathbf{X}^* & & & \dagger\mathbf{X}_3^* \end{array}$$

*where one constructs  $\mathbf{X}'_1$  by doing  $\sim(a|\ddot{a})k \rightarrow \emptyset$  on  $\mathbf{X}_1$ , and  $\mathbf{X}'_3$  by doing  $\sim\mathbf{X}^{\in}(aa|ee|ie|ii|oo|uu|\ddot{a}\ddot{a}|\ddot{o}\ddot{o}) \rightarrow \mathbf{X}^{\{1\}}$  on  $\mathbf{X}_3$ .*

(3) *Do  $\sim\dagger\mathbf{X}(m(a|\ddot{a})(\emptyset|isi)|m(ato|\ddot{a}t\ddot{o})|mi|ne|t(a|\ddot{a})|tt(a|\ddot{a})i \rightarrow \dagger\mathbf{X}^*$ .*

(4) *Mark the second protected range with  $\ddot{\cdot}$ .*

(5) *If the result so far does not end with  $*$ , append a symbol # at the end of the string.*

(6) *Replace*

$$\begin{array}{cc} \sim\dagger\mathbf{Xi}(M|mp)(a|i|\ddot{a}) & \sim\dagger\mathbf{X}(M|mp)(a|i|\ddot{a}) \\ \dagger\mathbf{X} & \dagger\mathbf{X} \text{ if } \mathbf{X} \text{ does not end with } \sim po \end{array}$$

(7) *Replace*

$$\begin{array}{cc} \sim\dagger(\hat{\chi}|\ddot{\cdot})_{m_0}(i|j)* & \sim\mathbf{X}^{\in}(\dagger(\hat{\chi}|\ddot{\cdot})_{m_0}\mathbf{V}|\mathbf{V}(\dagger(\hat{\chi}|\ddot{\cdot})_m)t(\emptyset|\ddot{\cdot})\# \\ \dagger(\hat{\chi}|\ddot{\cdot})_{m_0}* & \mathbf{X}\tau \end{array}$$

(8) *If the pattern  $\sim\dagger\mathbf{X}Ma\tau$  is found, replace it with  $\ddot{\mathbf{X}}'$ . Here,  $\mathbf{X}'$  is obtained by removing  $\sim(\emptyset|i)Ma$  from  $\mathbf{X}Ma$  if  $\mathbf{X}$  does not end with  $\sim po$ ; otherwise, set  $\mathbf{X}'$  as  $\mathbf{X}Ma$ .*

(9) *Do  $(\ddot{\cdot}|*|\#\tau) \rightarrow \emptyset$ .*

- (10)  $Do \sim \dagger \mathbf{X} \mathbf{V}^{**} \rightarrow \dagger \mathbf{X} (\mathbf{V}^{**})^{\{1\}}.$
- (11)  $Do \sim \dagger \mathbf{X} \mathbf{C}(a|e|i|\ddot{a}) \rightarrow \dagger \mathbf{X} \mathbf{C}.$
- (12)  $Do \sim \dagger \mathbf{X} \hat{\chi}^\epsilon(o|u)j \rightarrow \dagger \mathbf{X} \hat{\chi}.$
- (13)  $Do \sim \dagger \mathbf{X}_1 i(l)_m \mathbf{X}_2 \rightarrow \dagger \mathbf{X}_1, \sim \dagger \mathbf{X} C it \rightarrow \dagger \mathbf{X} Ci, \sim \dagger \mathbf{X} jo \rightarrow \dagger \mathbf{X} j, \sim \dagger \mathbf{X} ks \rightarrow \dagger \mathbf{X} d, \sim \dagger \mathbf{X} (st|ts) \rightarrow \dagger \mathbf{X}.$
- (14)  $Do \dagger \rightarrow \emptyset.$
- (15) Replace

$ng$	$nk$	$kah(\emptyset d dest t)$	$kä(y)\mathbf{X}$	$nä(e emm en et h hkö ht hty ijä ijäi ke ki$	$tun(\emptyset n s)$	$vase(\emptyset nt)$	$\sim oi$	$\sim ti$
$\gamma$	$\kappa$	$kaks$	$käy$	$nähd$	$tunt$	$vaseM$	$\emptyset$	$t$

- (16)  $Do \hat{\chi}_{\times 2} \rightarrow \hat{\chi}^+, \text{ that is, replace all double letters by the capital form of the corresponding single letter.}$

- (17)  $Do U(d|t) \sim \rightarrow Usi, \sim \hat{\chi}(U|Y)(\emptyset|d|s|t) \rightarrow \hat{\chi}, \underline{Us} \rightarrow Usi.$

For clustering purposes, we will also define an operation to undouble final consonants in the effective spelling of a Finnish word.

**Algorithm 8.7** (Reduction of Final Double Consonants). *Let  $\hat{\sigma}$  be the effective spelling of a Finnish word, then UndCons( $\hat{\sigma}$ ) is defined by the following procedure:*

- Look for the pattern  $\sim \hat{\chi}_{\times 2} \mathbf{V}_{m_0}$  in  $\hat{\sigma}$ .
- If the pattern is not found or if  $\hat{\chi}$  is one of the Finnish vowels {a, e, i, o, u, y, ä, ö}, do nothing; otherwise, replace the pattern by  $\hat{\chi} \mathbf{V}_{m_0}$ .

### 8.1.2 Admissible mutation and approximate clustering

**Algorithm 8.8** (Simple heredity test). *Construct  $\hat{\alpha}'$  by doing  $\sim(d|e|i|n|t|v) \rightarrow \emptyset$  on  $\hat{\alpha}$ . Define  $\hat{\beta}'$  similarly. The Boolean-valued function SimpHrdTest( $\hat{\alpha}, \hat{\beta}$ ) returns **TRUE** if the lowercase form of  $\hat{\alpha}$  contains at least one instance of {a, e, i, o, u, y, ä, ö} **AND** at least one of the following three conditions holds.<sup>99</sup>*

- (i)  $\hat{\alpha} = \hat{\beta};$
- (ii)  $\hat{\alpha}' = \hat{\beta}';$
- (iii)  $\ell(\hat{\beta}) > \ell(\hat{\alpha}) \geq \frac{\ell(\hat{\beta})}{2}$  **AND**  $\hat{\beta} = \hat{\alpha} \mathbf{V} \mathbf{X}$ .

In what follows, we define SuffixNW( $\hat{\alpha}, \hat{\beta}$ ), RootNW( $\hat{\alpha}, \hat{\beta}$ ), NW\*( $\hat{\alpha}, \hat{\beta}$ ) and SuffixSW( $\hat{\alpha}, \hat{\beta}$ ), SuffixSW( $\hat{\alpha}, \hat{\beta}$ ), SW\*( $\hat{\alpha}, \hat{\beta}$ ) as what is done in the Danish case (Algorithm 5.7), which was also applied to two other Germanic languages treated in §5.

Consonant gradations refer to consonant mutation patterns found in Finnish and its close relative, Estonian. (While Hungarian is also a Uralic language, it does not exhibit systematic consonant mutations as in Finnish.) In order to accommodate to this particular feature of Finnish, we need to define admissible consonant alternation patterns in Finnish effective spellings and essential roots.

**Definition 8.9** (Finnish consonant gradations). Define  $\mathbb{G} = [\emptyset, \emptyset][[k, \emptyset]][[k, j]][[k, Q]][[k, v]][[K, k]][[L, l]][[L, lt]][[M, mp]][[mp, nt]][[N, n]][[N, nt]][[p, v]][[P, p]][[R, r]][[s, d]][[s, t]][[S, s]][[S, st]][[t, d]][[t, s]][[T, s]][[T, t]][[\kappa, \gamma]]$  as admissible Finnish consonant gradations.

**Algorithm 8.10** (Admissible suffix mismatch and vowel alternation). *For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function*

$$\text{AdmMut}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta}))$$

*returns **TRUE** if at least one of the following three conditions holds:*

- (i)  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = [\emptyset|h|\mathbf{V}_m, \emptyset|h|\mathbf{V}_m] \text{ AND } (\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \mathbb{G} \text{ OR } \text{NW}^*(\hat{\beta}, \hat{\alpha}) = \mathbb{G} \text{ OR } (\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \emptyset \text{ AND the lowercase form of } \text{RootNW}(\hat{\alpha}, \hat{\beta}) \text{ contains at least one instance of } \mathbf{V});$
- (ii)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \emptyset \text{ AND } \text{SuffixNW}'(\hat{\alpha}, \hat{\beta}) = \mathbb{G}[[a, o]][[m, nt]][[n, t]], \text{ where one obtains } \text{SuffixNW}'(\hat{\alpha}, \hat{\beta}) \text{ by doing } \sim(h|s|\mathbf{V}_m) \rightarrow \emptyset \text{ on both components of } \text{SuffixNW}(\hat{\alpha}, \hat{\beta});$

<sup>99</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical **OR**.

- (iii)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \emptyset$  **AND**  $\text{SuffixNW}''(\hat{\beta}, \hat{\alpha}) = \mathbb{G}$ , where  $\text{SuffixNW}''(\hat{\beta}, \hat{\alpha})$  results from doing  $\sim(h|\mathbf{V}_m) \rightarrow \emptyset$  on both components of  $\text{SuffixNW}(\hat{\beta}, \hat{\alpha})$ .

Similarly, one can evaluate another Boolean-valued function

$$\text{AdmMut}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta}))$$

by trading all the occurrences of NW in the statements above with SW.

The following heredity test function is very similar to the German version (Algorithm ), except that we additionally require that the two strings in question must start with the same letter — constant gradations do not apply to the initial consonants.

**Algorithm 8.11** (Heredity test function). *In what follows,  $\ell([\hat{\sigma}, \hat{\tau}]) = \min\{\ell(\hat{\sigma}), \ell(\hat{\tau})\}$ . For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if at least one of the following three conditions holds:*

- (i)  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta}) = \text{TRUE}$ ;
- (ii)  $\ell(\text{RootNW}(\hat{\alpha}, \hat{\beta})) + \ell(\text{SuffixNW}(\hat{\alpha}, \hat{\beta})) + \ell(\text{NW}^*(\hat{\alpha}, \hat{\beta})) \geq \frac{\max\{\ell(\hat{\alpha}), \ell(\hat{\beta})\}}{2}$  **AND**  $\text{AdmMut}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta})) = \text{TRUE}$  **AND** both  $\hat{\alpha}$  and  $\hat{\beta}$  start with the same letter;
- (iii)  $\ell(\text{RootSW}(\hat{\alpha}, \hat{\beta})) + \ell(\text{SuffixSW}(\hat{\alpha}, \hat{\beta})) + \ell(\text{SW}^*(\hat{\alpha}, \hat{\beta})) \geq \frac{\max\{\ell(\hat{\alpha}), \ell(\hat{\beta})\}}{2}$  **AND**  $\text{AdmMut}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta})) = \text{TRUE}$  **AND** both  $\hat{\alpha}$  and  $\hat{\beta}$  start with the same letter.

**Algorithm 8.12** (Approximate clustering of Finnish words). *The approximate clustering of a list of Finnish words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  is completed in two stages:*

- (1) *We sort the list  $\{(\hat{\alpha}_1, \text{EffSpell}(\hat{\alpha}_1), \text{UndCons}(\text{EffSpell}(\hat{\alpha}_1))), \dots, (\hat{\alpha}_N, \text{EffSpell}(\hat{\alpha}_N), \text{UndCons}(\text{EffSpell}(\hat{\alpha}_N)))\}$  alphabetically according to the third component (with higher priority) and the second component (with lower priority). If two consecutive neighbors in the alphabetized list  $\{(\hat{\alpha}_{(1)}, \text{EffSpell}(\hat{\alpha}_{(1)}), \text{UndCons}(\text{EffSpell}(\hat{\alpha}_{(1)}))), \dots, (\hat{\alpha}_{(N)}, \text{EffSpell}(\hat{\alpha}_{(N)}), \text{UndCons}(\text{EffSpell}(\hat{\alpha}_{(N)})))\}$  satisfies both  $\text{HrdTest}(\text{EffSpell}(\hat{\alpha}_{(n)}), \text{EffSpell}(\hat{\alpha}_{(n+1)})) = \text{FALSE}$  and  $\text{HrdTest}(\text{EffSpell}(\hat{\alpha}_{(n+1)}), \text{EffSpell}(\hat{\alpha}_{(n)})) = \text{FALSE}$ , where  $n \in \mathbb{Z} \cap [1, N]$ , then add a demarcation line between these two entries. In this way, the alphabetized list is divided into separate groups of words with tags:  $\{g_1 = \{(\hat{\alpha}_{(1,1)}, \text{EffSpell}(\hat{\alpha}_{(1,1)}), \dots), \dots, (\hat{\alpha}_{(1,n_1)}, \text{EffSpell}(\hat{\alpha}_{(1,n_1)}), \dots)\}, \dots, g_M = \{(\hat{\alpha}_{(M,1)}, \text{EffSpell}(\hat{\alpha}_{(M,1)}), \dots), \dots, (\hat{\alpha}_{(M,n_M)}, \text{EffSpell}(\hat{\alpha}_{(M,n_M)}), \dots)\}\}$ , where each sublist also preserves alphabetic order. In particular, we have  $\hat{\alpha}_{(1,1)} = \hat{\alpha}_{(1)}$  and  $\hat{\alpha}_{(M,n_M)} = \hat{\alpha}_{(N)}$ .*
- (2) *For each group of words  $g_m = \{(\hat{\alpha}_{(m,1)}, \text{EffSpell}(\hat{\alpha}_{(m,1)}), \dots), \dots, (\hat{\alpha}_{(m,n_m)}, \text{EffSpell}(\hat{\alpha}_{(m,n_m)}))\}$  where  $m \in \mathbb{Z} \cap [1, M]$ , we augment it into a tagged entry  $G_m = (g_m, \text{EffSpell}(\hat{\alpha}_{(m,1)}), \text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}), T_m))$ , where one constructs the tag  $T_m$  from  $\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))$  in two steps:*
  - Replace all capital letters (including those in the standard Latin alphabet, as well as  $\ddot{A}$  and  $\ddot{O}$ ) by double letters in lowercase form, and do  $\sim o \rightarrow a$ .
  - Apply  $\text{UndCons}$ .

*The list  $\{G_1, \dots, G_M\}$  is sorted alphabetically, with respect to  $T_m$  (with highest priority),  $\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))$  (with medium priority), and  $\text{EffSpell}(\hat{\alpha}_{(m,1)})$  (with lowest priority). If two consecutive neighbors in the alphabetized list  $\{G_{(1)} = (g_{(1)}, \hat{\gamma}'_{(1)}, \hat{\gamma}''_{(1)}, \hat{\gamma}'''_{(1)}), \dots, G_{(M)} = (g_{(M)}, \hat{\gamma}'_{(M)}, \hat{\gamma}''_{(M)}, \hat{\gamma}'''_{(M)})\}$  satisfy*

$$\text{HrdTest}(\hat{\gamma}''_{(m)}, \hat{\gamma}''_{(m+1)}) = \text{FALSE}$$

**AND**

$$\text{HrdTest}(\hat{\gamma}'_{(m)}, \hat{\gamma}'_{(m+1)}) = \text{FALSE}$$

*where  $m \in \mathbb{Z} \cap [1, M]$ , then add a demarcation line between these two entries. In this manner, the alphabetized list  $\{G_{(1)}, \dots, G_{(M)}\}$  is divided into separate groups of tagged clusters  $\{G_{(1)} = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, G_{(K)} = \{G_{(K,1)}, \dots, G_{(1,m_K)}\}\}$ . Finally, generate a list of word clusters  $\{\check{G}_1 = (G_{(1,1)}, \dots), \dots, \check{G}_K = (G_{(K,1)}, \dots)\}$  by discarding all the tags.*

*Example 8.12.1.* There are 78 inflection classes in the Finnish language, according to the classification scheme of KOTUS ([https://en.wiktionary.org/wiki/Appendix:Finnish\\_nominal\\_inflection](https://en.wiktionary.org/wiki/Appendix:Finnish_nominal_inflection) and [https://en.wiktionary.org/wiki/Appendix:Finnish\\_conjugation](https://en.wiktionary.org/wiki/Appendix:Finnish_conjugation)). Among these classes, the first 51 classes describe nominal declensions (nouns, adjectives and numerals) while the rest are classes for verb conjugations. We pick one representative word from each inflection class, with or without gradation (consonant alternation in the stem).

The KOTUS class number will be shown parenthetically after the English translation for each word family of Finnish nominals below.

*valo, valoa, valoihin, valoaksi, valoilla, valoille, valolita, valoin, valoina, valoineen, valoissa, valoista, valoittha, valoja, valojen, valoksi, valolla, valolle, valolta, valon, valona, valoon, valossa, valosta, valot, valotta* — “light” (1);

*palvelu, palvelua, palveluiden, palveluihin, palveluksi, palveluilla, palveluille, palveluitta, palveluin, palveluina, palveluineen, palveluissa, palveluista, palveluita, palveluitten, palveluja, palvelujen, palveluksi, palvelulla, palvelulle, palvelulta, palvelun, palveluna, palvelussa, palvelusta, palvelut, palvelutta, palveluun* — “service” (2);

*valtio, valtioiden, valtioihin, valtioiksi, valtioilla, valtioille, valtioilta, valtioin, valtioina, valtioineen, valtioissa, valtioista, valtioita, valtioitta, valtioitten, valtioksi, valtiolla, valtiolle, valtiolta, valtion, valtiona, valtioon, valtiossa, valtiosta, valtiot, valtiota, valtiotta* — “government” (3);

*laatikko, laatikkoa, laatikkoihin, laatikkoina, laatikkoineen, laatikkoja, laatikkojen, laatikkona, laatikkoon, laatikoiden, laatikoihin, laatikoiksi, laatikoilla, laatikoille, laatikointa, laatikoiissa, laatikosta, laatikota, laatikotta, laatikottien, laatikoksi, laatikolla, laatikolle, laatikolta, laatikon, laatikossa, laatikosta, laatikot, laatikotta* — “box” (4);

*risteihiin, risteiksi, risteille, risteillä, risteiltä, ristein, risteineen, risteinä, risteissä, risteistä, risteittä, ristejä, risti, ristien, ristiin, ristiksi, ristille, ristillä, ristiltä, ristin, ristinä, ristissä, rististä, ristit, ristittä, ristiä* — “cross” (5);

*papereiden, papereihin, papereiksi, papereilla, papereille, papereilta, paperein, papereina, papereineen, papeissa, papereista, papereita, papereitta, papereitten, papereja, paperi, paperia, paperien, paperiin, paperiksi, paperilla, paperille, paperilta, paperin, paperina, paperissa, paperista, paperit, paperitta* — “paper” (6);

*ovea, oveen, oveksi, ovella, ovella, ovelta, oven, ovensa, ovessa, ovesta, ovet, ovetta, ovi, ovia, ovien, oviin, oviksi, ovilla, oville, ovilta, ovin, ovina, ovineen, ovissa, ovista, ovitta* — “door” (7);

*nalle, nallea, nalleen, nalleihin, nalleksi, nalleilla, nalleille, nalleilta, nallein, nalleina, nalleineen, nalleissa, nalleista, nalleitta, nalleja, nallejen, nalleksi, nallella, nallelle, nalleltä, nallen, nallena, nallessa, nallesta, nallet, nalletta* — “teddy bear” (8);

*kala, kalaan, kalaan, kalain, kalaksi, kalalla, kalalle, kalalta, kalan, kalana, kalassa, kalasta, kalat, kalatta, kaloihin, kaloiksi, kaloilla, kaloille, kaloilta, kaloin, kaloina, kaloineen, kaloissa, kaloista, kaloitta, kaloja, kalojen* — “fish” (9);

*kaira, kiraan, koiraan, koirain, koiraksi, koiralla, koiralte, koiranta, koirana, koirassa, koirasta, koirat, koiratta, koria, koirien, koiriin, koiriksi, koirilla, koirille, koirulta, koirin, koirina, koirineen, koirissa, koirista, koiritta* — “dog” (10);

*omena, omenaa, omenaan, omenain, omenaksi, omenalla, omenalle, omenalta, omenan, omenana, omenassa, omenasta, omenat, omenatta, omenia, omenien, omeniin, omeniksi, omenilla, omenille, omenilta, omenin, omenina, omenineen, omenissa, omenista, omenitta, omenoiden, omenoihin, omenoiksi, omenoilla, omenoille, omenoilita, omenoin, omenoina, omenoineen, omenoissa, omenoista, omenoita, omenoitta, omenoitten, omenoja, omenojen* — “apple” (11);

*kulkija, kulkijaa, kulkijaan, kulkijain, kulkijaksi, kulkijalla, kulkijalle, kulkijalta, kulkijan, kulkijana, kulkijassa, kulkijasta, kulkijat, kulkijatta, kulkijoiden, kulkijoihin, kulkijoiksi, kulkijoilla, kulkijoille, kulkijoilta, kulkijoin, kulkijoina, kulkijoineen, kulkijoissa, kulkijoista, kulkijoita, kulkijoitta, kulkijoitten* — “traveller” (12);

*katiska, katiskaa, katiskaan, katiskain, katiskaksi, katiskalla, katiskalle, katiskalta, katiskan, katiskana, katiskassa, katiskasta, katiskat, katiskatta, katiskoiden, katiskoihin, katiskoiksi, katiskoilla, katiskoille, katiskoilta, katiskoin, katiskoina, katiskoineen, katiskoissa, katiskoista, katiskoita, katiskoitta, katiskoitten, katiskoja, katiskojen* — “fishtrap” (13);

*solakaksi, solakalla, solakalle, solakalta, solakan, solakassa, solakasta, solakat, solakatta, solakka, solakkaa, solakkaan, solakkain, solakkana, solakkoihin, solakkoina, solakkoineen, solakkoja, solakkojen, solakoiden, solakoihin, solakoiksi, solakoilla, solakoille, solakoilta, solakoin, solakoissa, solakoista, solakoita, solakoitta, solakoitten* — “slender” (14);

*korkea, korkeaa, korkeaan, korkeain, korkeaksi, korkealla, korkealle, korkealta, korkean, korkeana, korkeassa, korkeasta, korkeat, korkeata, korkeatta, korkeiden, korkeihin, korkeiksi, korkeilla, korkeille, korkeilta, korkein, korkeina, korkeine, korkeisiin, korkeissa, korkeista, korkeita, korkeitta, korkeitten* — “tall” (15);

*vanhemaksi, vanhemmallla, vanhemmalle, vanhemmalta, vanhemman, vanhemmassa, vanhemmasta, vanhemmat, vanhemmatta, vanhemmiksi, vanhemmillla, vanhemmille, vanhemmilta, vanhemmin, vanhemmissa, vanhemmista, vanhemmitta, vanhempaa, vanhempaan, vanhempain, vanhempana, vanhempa, vanhempia, vanhempien, vanhempii, vanhempina, vanhempine* — “senior” (16);

*vapaa, vapaaksi, vapaalla, vapaalle, vapaalta, vapaan, vapaana, vapaaseen, vapaassa, vapaasta, vapaat, vapaata, vapaatta, vapaiden, vapainhin, vapaaksi, vapailla, vapaille, vapailta, vapain, vapaina, vapaine, vapaisiin, vapaissa, vapaista, vapaita, vapaitta, vapaitten* — “free” (17);

*maa, maahan, maaksi, maalla, maalle, maaalta, maan, maana, maassa, maasta, maat, maata, maatta, maiden, maihin, maaksi, mailla, mailta, main, maina, maineen, maissa, maista, maita, maitta, maitten* — “earth” (18);

*soiden, soihin, soiksi, soilla, soille, soilta, soin, soina, soineen, soissa, soista, soita, soitten, suo, suohon, suoksi, suolla, suolle, suulta, suon, suona, suossa, suosta, suot, suota, suotta* — “swamp” (19);

*filee, fileehen, fileeksi, fileelle, fileellä, fileeltä, fileen, fileenä, fileeseen, fileessä, fileestä, fileet, fileettä, fileetä, fileiden, fileihin, fileiksi, fileille, fileillä, fileiltä, filein, fileineen, fileinä, fileisiin, fileissä, fileistä, fileitten, fileittä, fileitä* — “fillet” (20);

*rosé, roséhen, roséiden, roséihin, roséiksi, roséilla, roséille, roséulta, roséin, roséina, roséineen, roséissa, roséista, roséita, roséitta, roséitten, roséksi, rosélla, rosélla, roséltä, rosén, roséna, roséssa, roséstä, roséta, roséttä* — “pinkwine” (21);

*parfait, parfait'hen, parfait'iden, parfait'ihin, parfait'iksi, parfait'illa, parfait'ilta, parfait'in, parfait'ina, parfait'ineen, parfait'issa, parfait'ista, parfait'ita, parfait'itta, parfait'itten, parfait'ksi, parfait'lla, parfait'lle, parfait'ita, parfait'n, parfait'na, parfait'ssa, parfait'sta, parfait't, parfait'ta, parfait'tta* — “parfait” (22);

*tiileen, tiileksi, tiilelle, tiilellä, tiileltä, tiilen, tiilenä, tiilessä, tiilestä, tiilet, tiilettä, tiili, tiilien, tiiliin, tiiliksi, tiilille, tiilillä, tiililtä, tiilin, tiilineen, tiilinä, tiilissä, tiilistä, tiilitä, tiiliä, tiiltä* — “brick” (23);

*uneen, uneksi, unella, unelle, unelta, unen, unena, unessa, unesta, unet, unetta, uni, unia, unien, uniin, uniksi, unilla, unille, unilta, unin, unina, unineen, unissa, unista, unitta, unta, unten* — “dream” (24);

*toimea, toimeen, toimeksi, toimella, toimelle, toimelta, toimen, toimena, toimessa, toimesta, toimet, toimetta, toimi, toimia, toimien, toimiin, toimiksi, toimilla, toimille, toimilta, toimin, toimina, toimineen, toimissa, toimista, toimitta, tointa, tointen* — “chore” (25);

*pieneen, pieneksi, pienelle, pienellä, pieneltä, pienen, pienenä, pienessä, pienestä, pienet, pienettä, pieni, pienien, pieniin, pieniksi, pienille, pienillä, pieniltä, pienin, pienine, pieninä, pienissä, pienistä, pienittä, pieniä, pienent, pienä — “small” (26);*

*kädeksi, kädelle, kädellä, kädeltä, käden, kädessä, kädestä, kädet, kädettä, käsi, käsien, käsiin, käsiksi, käsille, käsillä, käsiltä, käsin, käsineen, käsinä, käsisä, käsistä, käsittä, käsiä, käteen, kätenä, kätten, kättä* — “arm” (27);

*kynneksi, kynnelle, kynnellä, kynnetä, kynnen, kynnessä, kynnestä, kynnet, kynnettä, kynsi, kynsien, kynsiin, kynsiksi, kynsille, kynsillä, kynsiltä, kynsin, kynsineen, kynsinä, kynsissä, kynsistä, kynsittä, kynsiä, kynteen, kyntenä, kyntten, kynttä* — “fingernail” (28);

*lapseen, lapseksi, lapsella, lapselle, lapselta, lapsen, lapsena, lapsessa, lapsesta, lapset, lapsetta, lapsi, lapsia, lapsien, lapsiin, lapsiksi, lapsilla, lapsille, lapsilta, lapsin, lapsina, lapsineen, lapsissa, lapsista, lapsitta* — “child” (29);

*veisten, veistä, veitseen, veitseksi, veitselle, veitsellä, veitseltä, veitsen, veitsenä, veitsessä, veitsestä, veitset, veitsettä, veitsi, veitsien, veitsiin, veitsiksi, veitsille, veitsillä, veitsiltä, veitsin, veitsineen, veitsinä, veitsissä, veitsistä, veitsittä, veitsiä* — “knife” (30);

*kahdeksi, kahdella, kahdelle, kahdelta, kahden, kahdessa, kahdesta, kahdesti, kahdet, kahdetta, kahta, kahtaalla, kahtaalle, kahtalta, kahteen, kahtena, kahtia, kaksi, kaksia, kaksien, kaksiin, kaksiksi, kaksilla, kaksille, kaksilta, kaksin, kaksina, kaksine, kaksissa, kaksista, kaksitta, kaksittain* — “two” (31);

*sisar, sisareen, sisareksi, sisarella, sisarelle, sisarela, sisaren, sisarena, sisaressa, sisaresta, sisaret, sisaretta, sisaria, sisarien, sisariin, sisariksi, sisarilla, sisarille, sisarilta, sisarin, sisarina, sisarineen, sisarissa, sisarista, sisaritta, sisarta, sisarten* — “sister” (32);

*kytkimeen, kytkimeksi, kytkimelle, kytkimellä, kytkimeltä, kytkimen, kytkimenä, kytkimessä, kytkimestä, kytkimet, kytkimettä, kytkimien, kytkimiin, kytkimksi, kytkimille, kytkimellä, kytkimiltä, kytkimin, kytkimineen, kytkiminä, kytkimissä, kytkimistä, kytkimittä, kytkimiä, kytkin, kytkinten, kytkintä* — “switch” (33);

*onneton, onnetonta, onnetonten, onnettomaan, onnettomaksi, onnettomalla, onnettomalle, onnettomalta, onnettoman, onnettomana, onnettomassa, onnettomasta, onnettomat, onnettomatta, onnettomia, onnettomien, onnettomii, onnettomiksi, onnettomilla, onnettomille, onnettomilta, onnettomin, onnettomina, onnettomine, onnettomisa, onnettomista, onnettomitta* — “unhappy” (34);

*lämmin, lämmintä, lämpimien, lämpimiin, lämpimiksi, lämpimille, lämpimillä, lämpimiltä, lämpimin, lämpimine, lämpiminä, lämpimissä, lämpimistä, lämpimittä, lämpimiä, lämpimän, lämpimäksi, lämpimälle, lämpimälä, lämpimältä, lämpimän, lämpimänä, lämpimässä, lämpimästä, lämpimät, lämpimättä, lämpimäään* — “warm” (35);

*sisimmiksi, sisimille, sisimmillä, sisimiltä, sisimmin, sisimissä, sisimmistä, sisimittä, sisimmäksi, sisimmälle, sisimmällä, sisimmältä, sisimmän, sisimmässä, sisimmästä, sisimmät, sisimmättä, sisimpien, sisimpiin, sisimpine, sisimpinä, sisimpiä, sisimpän, sisimpänä, sisipäään, sisin, sisinten, sisintä* — “innermost” (36);

*vasemmaksi, vasemmalla, vasemmale, vasemmalta, vasemman, vasemmassa, vasemasta, vasemmat, vasemmatta, vasemmiksi, vasemmillä, vasemmile, vasemmita, vasemmin, vasemmissa, vasemmista, vasemmitta, vasempaa, vasempaan, vasempain, vasempana, vasempia, vasempien, vasempiin, vasempina, vasempine, vasen, vasenta, vasenten* — “left” (37);

*nainen, naiseen, naiseksi, naisella, naiselle, naiselta, naisen, naisena, naisessa, naisesta, naiset, naisetta, naisia, naisien, naisiin, naisiksi, naisilla, naisille, naisulta, naisin, naisina, naisineen, naisissa, naisista, naisitta, naista, naisten* — “woman” (38);

*vastaukseen, vastaukseksi, vastauksella, vastaukselle, vastaukselta, vastaukseni, vastauksena, vastauksessa, vastauksesta, vastaukset, vastauksetta, vastauksia, vastauksien, vastauksiin, vastauksiksi, vastauksilla, vastauksille, vastauksilta, vastauksin, vastauksina, vastauksineen, vastauksissa, vastauksista, vastauksitta, vastaus, vastauta, vastausten* — “answer” (39);

*kalleudeksi, kalleudella, kalleudelle, kalleudelta, kalleuden, kalleudessa, kalleudesta, kalleudet, kalleudetta, kalleuksia, kalleuksien, kalleuksiin, kalleuksiksi, kalleuksilla, kalleuksille, kalleuksilta, kalleuksin, kalleuksina, kalleuksineen, kalleuksissa, kalleuksista, kalleuksitta, kalleus, kalleuteen, kalleutena, kalleutta* — “expensiveness” (40);

*vieraaksi, vieraalla, vieraalle, vieraalta, vieraan, vieraana, vieraaseen, vieraassa, vieraasta, vieraat, vieraatta, vieraiden, vieraihin, vieraaksi, vierailla, vieraille, vierailta, vierain, vieraina, vieraine, vieraisiin, vieraissa, vieraista, vieraita, vieraitten, vieraitten, vieras, vierasta* — “unfamiliar” (41);

*mieheen, mieheksi, miehelle, miehellä, mieheltä, miehen, miehenä, miehessä, miehestä, miehet, miehettä, miehien, miehiin, miehiksi, miehille, miehillä, miehiltä, miehin, miehineen, miehinä, miehissä, miehistä, miehittä, miehiä, mies, miesten, miestä* — “man” (42);

*ohueen, ohueksi, ohuella, ohuelle, ohuelta, ohuen, ohuena, ohuessa, ohuesta, ohuet, ohuetta, ohuiden, ohuihin, ohuiksi, ohuilla, ohuille, ohuilta, ohuin, ohuina, ohuine, ohuisiin, ohuissa, ohuista, ohuita, ohuitta, ohuitten, ohut, ohutta* — “thin” (43);

*keväiden, keväihin, keväaksi, keväille, keväillä, keväiltä, keväin, keväineen, keväinä, keväisiin, keväissä, keväistä, keväitten, keväitä, keväitä, kevät, kevättä, kevääksi, keväälle, keväällä, keväiltä, kevään, keväänä, kevääseen, keväässä, keväästä, keväät, keväättä* — “spring” (44);

*kahdeksanneksi, kahdeksannella, kahdeksannelle, kahdeksannelta, kahdeksannen, kahdeksannessa, kahdeksanesta, kahdeksannet, kahdeksannetta, kahdeksansia, kahdeksansien, kahdeksansiin, kahdeksansiksi, kahdeksansilla, kahdeksansille, kahdeksansilta, kahdeksansin, kahdeksansina, kahdeksansineen, kahdeksansissa, kahdeksansista, kahdeksansitta, kahdeksanteen, kahdeksantena, kahdeksas, kahdeksatta* — “eighth” (45);

*tuhanneksi, tuhannella, tuhannelle, tuhannelta, tuhannen, tuhannessa, tuhannesta, tuhannet, tuhannetta, tuhansia, tuhansien, tuhansiin, tuhansiksi, tuhansilla, tuhansille, tuhansilta, tuhansin, tuhansina, tuhansine, tuhansissa, tuhansista, tuhansitta, tuhanteen, tuhanten, tuhantena, tuhat, tuhatta* — “thousand” (46);

*kuolleeksi, kuolleella, kuolleelle, kuolleelta, kuolleen, kuolleena, kuolleeseen, kuolleessa, kuolleesta, kuolleet, kuolleetta, kuolleiden, kuolleihin, kuolleiksi, kuolleilla, kuolleille, kuolleilta, kuollein, kuolleina, kuolleine, kuoleisiin, kuolleissa, kuolleista, kuolleita, kuolleitta, kuolleitten, kuollut, kuollutta* — “dead” (47);

*hame, hameeksi, hameella, hameelle, hameelta, hameen, hameena, hameeseen, hameessa, hameesta, hameet, hameetta, hameiden, hameihin, hameiksi, hameilla, hameille, hameita, hamein, hameina, hameineen, hameisiin, hameissa, hameista, hameita, hameitta, hameitten, hametta* — “skirt” (48);

*askel, askele, askeleeksi, askeleella, askeleelle, askeleelta, askeleen, askeleena, askeleeseen, askeleessa, askeleesta, askeleet, askeleetta, askeleiden, askeleiksi, askeleilla, askeleille, askelelta, askelein, askeleina, askeleineen, askeleisiin, askeleissa, askeleista, askeleita, askeleitta, askeleitten, askeleksi, askelella, askelelle, askeleita, askelen, askelena, askelessa, askelest, askelet, askeleta, askelia, askelien, askeliin, askeliksi, askelilla, askelille, askelilta, askelin, askelina, askelineen, askelissa, askelista, askelitta, askelta, askelten* — “step” (49);

*isoäideiksi, isoäideille, isoäideillä, isoäideiltä, isoäidein, isoäideissä, isoäideistä, isoäideittä, isoäidiksi, isoäidil-  
le, isoäidillä, isoäidiltä, isoäidin, isoäidissä, isoäidistä, isoäidit, isoäidittä, isoäiteihin, isoäiteineen, isoäiteinä,  
isoäitejä, isoäiti, isoäitiin, isoäitinä, isoäitiä* — “grandmother” (50);

*nuorenpariin, nuoreksipariksi, nuorellaparilla, nuorelleparille, nuoreltaparilta, nuorenparina, nuorenparin,  
nuoressaparissa, nuorestaparista, nuoretparit, nuoretparitta, nuoriapareja, nuorienparien, nuoriinpareihin, nuo-  
riksipareiksi, nuorillapareilla, nuorillepareille, nuoriltapareilta, nuorinapareina, nuorinpareineen, nuorinpa-  
rein, nuoripari, nuorissapareissa, nuoristapareista, nuorittapareitta, nuortaparia, nuortenparien* — “newly-wed  
couple” (51).

Throwing all the words above to our algorithm, we receive the following output:

{*valo, valoa, valoihin, valoiksi, valoilla, valoille, valoulta, valoin, valoina, valoineen, valoissa, valoista, valoitta,  
valoja, valojen, valoksi, valolla, valolle, valolta, valon, valona, valoon, valossa, valosta, valot, valotta*} — (1) 26  
words;

{*palvelu, palvelua, palveluiden, palveluihin, palveluiksi, palveluilla, palveluille, palveluilta, palveluin, palvelui-  
na, palveluineen, palveluissa, palveluista, palveluita, palveluitta, palveluitten, palveluja, palvelujen, palveluksi,  
palvelulla, palvelulle, palvelulta, palvelun, palveluna, palvelussa, palvelusta, palvelut, palvelutta, palveluun*} —  
(2) 29 words;

{*valtio, valtioiden, valtioihin, valtioiksi, valtioilla, valtioille, valtioilta, valtioin, valtioina, valtioineen, valtiois-  
sa, valtioista, valtioita, valtioitta, valtioitten, valtioksi, valtiolla, valtiolle, valtiolta, valtion, vältion, vältioon,  
vältiassa, vältiosta, vältiot, vältiota, vältiotta*} — (3) 27 words;

{*laatikko, laatikkoo, laatikkoihin, laatikkoina, laatikkoineen, laatikkoja, laatikkojen, laatikkona, laatikkoon, laa-  
tikoiden, laatikoihin, laatikoiksi, laatikoilla, laatikolla, laatikolle, laatikointa, laatikoiissa, laatikoista, laatikoita,  
laatikoitta, laatikoitien, laatikoksi, laatikolla, laatikolle, laatikonta, laatikona, laatikossa, laatikosta, laatikot, laa-  
tikotta*} — (4) 30 words;

{*risteihin, risteiksi, risteille, risteillä, risteiltä, ristein, risteineen, risteinä, risteissä, risteistä, risteittä, ristejä,  
risti, ristien, ristiin, ristiksi, ristille, ristillä, ristiltä, ristin, ristinä, ristissä, rististä, ristit, ristittä, ristiä*} — (5) 26  
words;

{*papereiden, papereihin, papereiksi, papereilla, papereille, papereilta, paperein, papereina, papereineen, pape-  
reissa, papereista, papereita, papereitta, papereitten, papereja, paperi, paperia, paperien, paperiin, paperiksi,  
paperilla, paperille, paperilta, paperin, paperina, paperissa, paperista, paperit, paperitta*} — (6) 29 words;

{*ovea, oveen, oveksi, ovella, ovelle, ovelta, oven, ovena, ovessa, ovesta, ovet, ovetta, ovi, ovia, ovien, oviin, oviksi,  
ovilla, oville, ovilta, ovin, ovina, ovineen, ovissa, ovista, ovitta*} — (7) 26 words;

{*nalle, nallea, nalleen, nalleihin, nalleiksi, nalleilla, nalleille, nalleilta, nallein, nalleina, nalleineen, nalleissa,  
nalleista, nalleitta, nalleja, nallejen, nalleksi, nallella, nallelle, nalleltta, nallen, nallena, nallessa, nallesta, nallet,  
nalletta*} — (8) 26 words;

{*kala, kalaa, kalaan, kalain, kalaksi, kalalla, kalalle, kalalta, kalan, kalana, kalassa, kalasta, kalat, kalatita, kaloihin, kaloiksi, kaloilla, kaloille, kaloilta, kaloin, kaloina, kaloineen, kaloissa, kaloista, kaloitta, kaloja, kalojen*} — (9) 27 words;

{*kaira, koiraa, koiraan, koirain, koiraksi, koiralla, koiralle, koiralta, koiran, koirana, koirassa, koirasta, koirat, koiratta, koria, koirien, koirin, koiraksi, koirilla, koirille, koirilta, koirin, koirina, koirineen, koirissa, koirista, koiritta*} — (10) 27 words;

{*omena, omenaa, omenaan, omenain, omenaksi, omenalla, omenalte, omenalta, omenan, omenana, omenassa, omenasta, omenat, omenatta, omenia, omenien, omeniin, omeniksi, omenilla, omenille, omenulta, omenin, omenina, omenineen, omenissa, omenista, omenitta, omenoiden, omenoihin, omenoiksi, omenoilla, omenoille, omenoiltta, omenoin, omenoina, omenoineen, omenoissa, omenoista, omenoita, omenoitta, omenoitten, omenoja, omenojen*} — (11) 43 words;

{*kulkija, kulkijaa, kulkijaan, kulkijain, kulkijaksi, kulkijalla, kulkijalle, kulkijalta, kulkijan, kulkijana, kulkijassa, kulkijasta, kulkijat, kulkijatta, kulkijoiden, kulkijoihin, kulkijoiksi, kulkijoilla, kulkijoille, kulkijoilta, kulkijoin, kulkijoina, kulkijoineen, kulkijoissa, kulkijoista, kulkijoitta, kulkijoitten*} — (12) 28 words;

{*katiska, katiskaan, katiskaan, katiskain, katiskaksi, katiskalla, katiskalle, katiskalta, katiskan, katiskana, katiskassa, katiskasta, katiskat, katiskatta, katiskoiden, katiskoihin, katiskoiksi, katiskoilla, katiskoille, katiskoilta, katiskoin, katiskoina, katiskoineen, katiskoissa, katiskoista, katiskoita, katiskoitta, katiskoitten, katiskoja, katiskojen*} — (13) 30 words;

{*solakaksi, solakalla, solakalle, solakalta, solakan, solakassa, solakasta, solakat, solakatta, solakka, solakkaa, solakkaan, solakkain, solakkana, solakkoihin, solakkoina, solakkoineen, solakkoja, solakkojen, solakoiden, solakoihin, solakoiksi, solakoilla, solakoille, solakoilta, solakoin, solakoissa, solakoista, solakoita, solakoitta, solakoitten*} — (14) 31 words;

{*korkea, korkeaa, korkeaan, korkeain, korkeaksi, korkealla, korkealle, korkealta, korkean, korkeana, korkeassa, korkeasta, korkeat, korkeata, korkeatta, korkeiden, korkeihin, korkeiksi, korkeilla, korkeilta, korkein, korkeina, korkeine, korkeisiin, korkeissa, korkeista, korkeita, korkeitta, korkeitten*} — (15) 30 words;

{*vanhemaksi, vanhemmall, vanhemmalle, vanhemmalta, vanhemman, vanhemmassa, vanhemmasta, vanhemmat, vanhemmatta, vanhemmiksi, vanhemmill, vanhemmilta, vanhemmin, vanhemmissa, vanhemmista, vanhemmitta, vanhempaa, vanhempaan, vanhempain, vanhempana, vanhemi, vanhempia, vanhempien, vanhempipi, vanhempina, vanhempine*} — (16) 27 words;

{*vapaa, vapaaksi, vapaalla, vapaalle, vapaalta, vapaan, vapaana, vapaaseen, vapaassa, vapaasta, vapaat, vapaata, vapaatta, vapaiden, vapaihin, vapaaksi, vapailla, vapaille, vapailta, vapain, vapaina, vapaine, vapaistiin, vapaissa, vapaista, vapaita, vapaitta, vapaitten*} — (17) 28 words;

{*maineen*} — (18) 1 word;

{*maa, maahan, maaksi, maalla, maalle, maalta, maan, maana, maassa, maasta, maat, maata, maatta, maiden, maihin, maiksi, mailla, mailta, main, maina, maissa, maista, maita, maitta, maitten*} — (18) 26 words;

{*soiden, soihin, soiksi, soilla, soille, soilta, soin, soina, soineen, soissa, soista, soita, soitta, soitten, suo, suohon, suoksi, suolla, suolle, suolta, suon, suona, suossa, suosta, suot, suota, suotta*} — (19) 27 words;

{*filee, fileehen, fileeksi, fileelle, fileellä, fileeltä, fileen, fileenä, fileeseen, fileessä, fileestä, fileet, fileettä, fileetä, fileiden, fileihin, fileiksi, fileille, fileillä, fileiltä, filein, fileineen, fileinä, fileisiin, fileissä, fileistä, fileitten, fileittä, fileitä*} — (20) 29 words;

{*rosé, roséhen, roséiden, roséihin, roséiksi, roséilla, roséille, roséilda, roséin, roséina, roséineen, roséissa, roséista, roséita, roséitta, roséitten, roséksi, rosélla, rosélle, rosélda, rosén, roséna, roséssa, rosést, rosét, roséttä, roséttä*} — (21) 27 words;

{*parfait, parfait'hen, parfait'iden, parfait'ihin, parfait'iksi, parfait'illa, parfait'ilta, parfait'itta, parfait'in, parfait'ina, parfait'ineen, parfait'issa, parfait'ista, parfait'ita, parfait'itten, parfait'ksi, parfait'lla, parfait'lle, parfait'ita, parfait'n, parfait'na, parfait'ssa, parfait'sta, parfait't, parfait'ta, parfait'tta*} — (22) 27 words;

{*tiileen, tiileksi, tiilelle, tiilellä, tiileltä, tiilen, tiilenä, tiilessä, tiilestää, tiilet, tiilettä, tiili, tiilien, tiiliin, tiiliksi, tiilille, tiilllä, tiililtä, tiilin, tiilineen, tiilinä, tiilissä, tiilistää, tiilitä, tiiliä, tiiltä*} — (23) 26 words;

{*uneen, uneksi, unella, unelle, unelta, unen, unena, unessa, unesta, unet, unetta, uni, unia, unien, uniin, uniksi, unilla, unille, unilta, unin, unina, unineen, unissa, unista, unitta, unta, unten*} — (24) 27 words;

{*toimea, toimeen, toimeksi, toimella, toimelle, toimelta, toimen, toimena, toimessa, toimesta, toimet, toimetta, toimi, toimia, toimien, toimiin, toimiksi, toimilla, toimille, toimlta, toimin, toimina, toimineen, toimissa, toimista, toimitta, tointa, tointen*} — (25) 28 words;

{*pieneen, pieneksi, pienelle, pienellä, pieneltä, pienen, pienenä, pienessä, pienestää, pienet, pienettä, pieni, pienien, pieniin, pieniksi, pienille, pienillä, pieniltä, pienin, pienine, pieninä, pienissä, pienistä, pienittä, pieniä, pienien, pientä*} — (26) 27 words;

{*kädeksi, kädelle, kädellä, kädeltä, käden, kädessä, kädestä, kädet, kädettä, käsi, käsien, käsiin, käsiksi, käsille, käsillä, käsiltä, käsin, käsineen, käsinä, käsisä, käsistä, käsittä, käsiä, käteen, kätenä, kätten, kättä*} — (27) 27 words;

{*kynneksi<sub>28</sub>, kynnelle<sub>28</sub>, kynnellä<sub>28</sub>, kynneltä<sub>28</sub>, kynnen<sub>28</sub>, kynnessä<sub>28</sub>, kynnestä<sub>28</sub>, kynnet<sub>28</sub>, kynnettä<sub>28</sub>, kynsi<sub>28</sub>, kynsien<sub>28</sub>, kynsiin<sub>28</sub>, kynski<sub>28</sub>, kynsille<sub>28</sub>, kynsillä<sub>28</sub>, kynsiltä<sub>28</sub>, kynsin<sub>28</sub>, kynsineen<sub>28</sub>, kynsinä<sub>28</sub>, kynssä<sub>28</sub>, kynsistä<sub>28</sub>, kynsittä<sub>28</sub>, kynsiä<sub>28</sub>, kynteen<sub>28</sub>, kyntenä<sub>28</sub>, kyntten<sub>28</sub>, kynttä<sub>28</sub>, kytkin<sub>33</sub>*} — (28) 27 words; (33) 1 word;

{*lapseen, lapseksi, lapsella, lapselle, lapselta, lapsen, lapsena, lapsessa, lapsesta, lapset, lapsetta, lapsi, lapsia, lapsien, lapsiin, lapsiksi, lapsilla, lapsille, lapsilta, lapsin, lapsina, lapsineen, lapsissa, lapsista, lapsitta*} — (29) 25 words;

{*veisten, veistä, veitseen, veitseksi, veitselle, veitsellä, veitseltä, veitsen, veitsenä, veitsessä, veitsestä, veitset, veitsettä, veitsi, veitsien, veitsiin, veitsiksi, veitsille, veitsillä, veitsiltä, veitsin, veitsineen, veitsinä, veitsissä, veitsistä, veitsittä, veitsiä*} — (30) 27 words;

{*kahdeksi, kahdella, kahdelle, kahdelta, kahden, kahdessa, kahdesta, kahdesti, kahdet, kahdetta, kahta, kahtaalla, kahtaalle, kahtaalta, kahteen, kahtena, kahtia, kaksi, kaksia, kaksien, kaksiin, kaksiksi, kaksilla, kaksille, kaksilta, kaksin, kaksina, kaksine, kaksissa, kaksista, kaksitta, kaksittain*} — (31) 32 words;

{*sisar, sisareen, sisareksi, sisarella, sisarelle, sisarelda, sisaren, sisarena, sisaressa, sisaresta, sisaret, sisaretta, sisaria, sisarien, sisariin, sisariksi, sisarilla, sisarille, sisarilta, sisarin, sisarina, sisarineen, sisarissa, sisarista, sisaritta, sisarta, sisarten*} — (32) 27 words;

{*kytkimeen, kytkimeksi, kytkimelle, kytkimellä, kytkimeltä, kytkimen, kytkimenä, kytkimessä, kytkimestä, kytkimet, kytkimettä, kytkimien, kytkimiin, kytkimiksi, kytkimille, kytkimillä, kytkimittä, kytkimin, kytkimeen, kytkiminä, kytkimissä, kytkimistä, kytkimittä, kytkimiä, kytkinten, kytkintä*} — (33) 26 words;

{*onneton, onnetonta, onnetonten, onnettomaan, onnettomaaksi, onnettomaalla, onnettomalte, onnettomaalta, onnettoman, onnettoman, onnettomaassa, onnettomaasta, onnettomat, onnettomaatta, onnettomi, onnettomi, onnettomien, onnettomii, onnettomi, onnettomi, onnettomailla, onnettommille, onnettomialta, onnettomin, onnettomin, onnettomin, onnettomin, onnettomaissa, onnettomaista, onnettommitta*} — (34) 27 words;

{*lämin, lämmintä, lämpimien, lämpimiin, lämpimiksi, lämpimille, lämpimillä, lämpimiltä, lämpimin, lämpimine, lämpiminä, lämpimissä, lämpimistä, lämpimittä, lämpimiä, lämpimän, lämpimäksi, lämpimälle, lämpimällä, lämpimältä, lämpimän, lämpimänä, lämpimässä, lämpimästä, lämpimät, lämpimättä, lämpimään*} — (35) 27 words;

{*sisimmiksi, sisimille, sisimillä, sisimmiltä, sisimmin, sisimissä, sisimmistä, sisimittä, sisimmäksi, sisimälle, sisimmällä, sisimmältä, sisimmän, sisimmässä, sisimmästä, sisimmät, sisimmättä, sisimpien, sisimpiin, sisimpine, sisimpiinä, sisimpiä, sisimpän, sisimpänä, sisimpää, sisin, sisinten, sisintä*} — (36) 28 words;

{*vasemmaksi, vasemalla, vasemmalle, vasemmalta, vasemman, vasemmassa, vasemasta, vasemmat, vasematta, vasemaksi, vasemilla, vasemelle, vasemmlta, vasemmin, vasemmissa, vasemista, vasemmitta, vasempaa, vasempaan, vasempain, vasempana, vasempia, vasempien, vasepiin, vasempina, vasempine, vasen, vasenta, vasenten*} — (37) 29 words;

{*naisin*} — (38) 1 word;

{*nainen, naiseen, naiseksi, naisella, naiselle, naiselta, naisen, naisena, naisessa, naisesta, naiset, naisetta, naisia, naisien, naisiin, naisiksi, naisilla, naisille, naisilta, naisina, naisineen, naisissa, naisista, naisitta, naista, naisten*} — (38) 26 words;

{*vastaukseen, vastaukseksi, vastauksella, vastaukselle, vastaukselta, vastauksen, vastauksena, vastauksessa, vastauksesta, vastaukset, vastauksetta, vastauksia, vastauksien, vastauksiin, vastauksiksi, vastauksilla, vastauksille, vastauksilta, vastauksin, vastauksina, vastauksineen, vastauksissa, vastauksista, vastauksitta, vastaus, vastausta, vastausten*} — (39) 27 words;

{*kalleudeksi, kalleudella, kalleudelle, kalleudelta, kalleuden, kalleudessa, kalleudesta, kalleudet, kalleudetta, kalleuksia, kalleuksien, kalleuksiin, kalleuksiksi, kalleuksilla, kalleuksille, kalleuksilta, kalleuksin, kalleuksina, kalleuksineen, kalleuksissa, kalleuksista, kalleuksitta, kalleus, kalleuteen, kalleutena, kalleutta*} — (40) 26 words;

{*vieraaksi, vieraalla, vieraalle, vieraalta, vieraan, vieraana, vieraaseen, vieraassa, vieraasta, vieraat, vieraatta, vieraiden, vieraihin, vieraiksi, vierailla, vieraille, vierailta, vierain, vieraina, vieraine, vieraasiin, vieraissa, vieraista, vieraita, vieraitta, vieraitten, vieras, vierasta*} — (41) 28 words;

{*mieheen, mieheksi, miehelle, miehellä, mieheltä, miehen, miehenä, miehessä, miehestä, miehet, miehettä, miehien, miehiin, miehiksi, miehille, miehillä, miehiltä, miehin, miehineen, miehinä, miehissä, miehistä, miehittä, miehiä, mies, mesten, mestä*} — (42) 27 words;

{*ohueen, ohueksi, ohuelle, ohuelle, ohuelta, ohuen, ohuena, ohuessa, ohuesta, ohuet, ohuetta, ohuiden, ohuihin, ohuiksi, ohuilla, ohuille, ohuulta, ohuin, ohuina, ohuine, ohuisiin, ohuissa, ohuista, ohuita, ohuitta, ohuitten, ohut, ohutta*} — (43) 28 words;

{*keväiden, keväihin, keväiksi, keväille, keväillä, keväiltä, keväin, keväineen, keväinä, keväisiin, keväissä, keväistä, keväitten, keväittä, keväitä, kevät, kevättä, kevääksi, keväälle, keväällä, keväiltä, kevään, keväänä, kevääseen, keväässä, keväästä, keväät, keväättä*} — (44) 28 words;

{*kahdeksanneksi, kahdeksannella, kahdeksannelle, kahdeksannelta, kahdeksannen, kahdeksannessa, kahdeksanesta, kahdeksannet, kahdeksannetta, kahdeksansia, kahdeksansien, kahdeksansiin, kahdeksansiksi, kahdeksansilla, kahdeksansille, kahdeksansilta, kahdeksansin, kahdeksansina, kahdeksansineen, kahdeksansissa, kahdeksansista, kahdeksansitta, kahdeksanteen, kahdeksantena, kahdeksas, kahdeksatta*} — (45) 26 words;

{*tuhanneksi, tuhannella, tuhannelle, tuhannelta, tuhannen, tuhannessa, tuhannesta, tuhannet, tuhannetta, tuhansia, tuhansien, tuhansiin, tuhansiksi, tuhansilla, tuhansille, tuhansilta, tuhansin, tuhansina, tuhansine, tuhansissa, tuhansista, tuhansitta, tuhanteen, tuhanten, tuhantena, tuhat, tuhatta*} — (46) 27 words;

{*kuolleeksi, kuolleella, kuolleelle, kuolleelta, kuolleen, kuolleena, kuolleeseen, kuolleessa, kuolleesta, kuolleet, kuolleetta, kuolleiden, kuolleihin, kuolleiksi, kuolleilla, kuolleille, kuolleilta, kuollein, kuolleina, kuolleine, kuollesiin, kuolleissa, kuolleista, kuolleita, kuolleitta, kuolleitten, kuollut, kuollutta*} — (47) 28 words;

{*hame, hameeksi, hameella, hameelle, hameelta, hameen, hameena, hameeseen, hameessa, hameesta, hameet, hameetta, hameiden, hameihin, hameiksi, hameilla, hameille, hameita, hameina, hameineen, hameisiin, hameissa, hameista, hameita, hameitta, hameitten, hametta*} — (48) 28 words;

{*askel, askele, askeleeksi, askeleella, askeleelle, askeleelta, askeleen, askeleena, askeleeseen, askeleessa, askeleesta, askeleet, askeleetta, askeleiden, askeleksi, askelella, askeleille, askeleilta, askelein, askeleina, askelein, askeleinen, askeleisiin, askeleissa, askeleista, askeleita, askeleitta, askeleitten, askeleksi, askelella, askelelle, askeleta, askelen, askelena, askelessa, askeleta, askelet, askeletta, askelia, askelien, askeliin, askeliksi, askelilla, askelille, askelilta, askelin, askelina, askelineen, askelissa, askelista, askelitta, askelta, askelten*} — (49) 52 words;

{*isoäideiksi, isoäideille, isoäideillä, isoäideiltä, isoäidein, isoäideissä, isoäideistä, isoäideittä, isoäidiksi, isoäidille, isoäidillä, isoäidiltä, isoäidin, isoäidissä, isoäidistä, isoäidit, isoäidittä, isoäiteihin, isoäiteineen, isoäiteinä, isoäitejä, isoäiti, isoäitiin, isoäitinä, isoäitiä*} — (50) 26 words;

{*nuoreenpariin*} — (51) 1 word;

{*nuoreksipariksi*} — (51) 1 word;

{*nuorellaparilla*} — (51) 1 word;  
 {*nuorelleparille*} — (51) 1 word;  
 {*nuoreltaparilta*} — (51) 1 word;  
 {*nuorenaparina*} — (51) 1 word;  
 {*nuorenparin*} — (51) 1 word;  
 {*nuoretparit*} — (51) 1 word;  
 {*nuorettaparitta*} — (51) 1 word;  
 {*nuoriapareja*} — (51) 1 word;  
 {*nuorienparien*} — (51) 1 word;  
 {*nuoriinpareihin*} — (51) 1 word;  
 {*nuoriksipareiksi*} — (51) 1 word;  
 {*nuorinapareina*} — (51) 1 word;  
 {*nuorinepareineen*} — (51) 1 word;  
 {*nuorinparein*} — (51) 1 word;  
 {*nuoripari*} — (51) 1 word;  
 {*nuorittapareitta*} — (51) 1 word;  
 {*nuortaparia*} — (51) 1 word;  
 {*nuortenparien*} — (51) 1 word;  
 {*nuoressaparissa, nuorestaparista*} — (51) 2 words;  
 {*nuorissapareissa, nuoristapareista*} — (51) 2 words;  
 {*nuorillapareilla, nuorillepareille, nuoriltapareilta*} — (51) 3 words;

Here, as we present the results above, each cluster is followed by a brief description about the KOTUS affiliations of the words therein. When a cluster contains words from more than one KOTUS class, the words within the cluster are labeled with class numbers in subscripts.

The word *maineen* is both the comitative plural form of *maa* “earth” and the genitive singular form of *maine* “fame”. Since the comitative case is rarely used, it is understandable that our algorithm separated *maineen* from the other inflected forms of *maa*.

For a similar reason, the word *naisin* was not treated as the instructive plural (a rarely used inflection) form of *nainen* “woman”, but was considered as the first personal singular conditional present of *naida* “marry”.

For more information regarding the compound word in KOTUS-51, see Example 8.14.1.

As we test our algorithm against the following sample verbs:

*sano, sanoa, sanoakseen, sanoen, sanoessa, sanoi, sanoimme, sanoin, sanoisi, sanoisimme, sanoisin, sanoisit, sanoisitte, sanoisivat, sanoit, sanoitte, sanoivat, sanokaa, sanokaamme, sanoko, sanokoon, sanokoot, sanoma, sanomaan, sanomaisillaan, sanomalla, sanoman, sanomassa, sanomasta, sanomatton, sanomatta, sanominen, sanomista, sanomme, sanon, sanone, sanonee, sanoneet, sanonemme, sanonen, sanonet, sanonette, sanonevat, sanonut, sanoo, sanot, sanota, sanotaan, sanottaessa, sanottaisi, sanottaisiin, sanottako, sanottakoon, sanottaman, sanottane, sanottaneen, sanottava, sanotte, sanottiin, sanottu, sanova, sanovat* — “tell” (52);

*muista, muistaa, muistaakseen, muistaen, muistaessa, muistaisi, muistaisimme, muistaisin, muistaisit, muistaisitte, muistaisivat, muistakaa, muistakaamme, muistako, muistakoon, muistakoot, muistama, muistamaan, muistamaisillaan, muistamalla, muistaman, muistamassa, muistamasta, muistamatton, muistamatta, muistaminen, muistamista, muistamme, muistan, muistane, muistanee, muistaneet, muistanemme, muistanan, muistananet, muistannerette, muistanevat, muistanut, muistat, muistatte, muistava, muistavat, muisteta, muistetaan, muistettaessa, muistettai-si, muistettaiiin, muistettako, muistettakoon, muistettaman, muistettane, muistettaneen, muistettava, muistettiin, muistettu, muisti, muistimme, muistin, muistit, muistitte, muistivat* — “remember” (53);

*huuda, huudamme, huudan, huudat, huudatte, huudeta, huudetaan, huudettaessa, huudettaisi, huudettaisiin, huudettako, huudettakoon, huudettaman, huudettane, huudettaneen, huudettava, huudettiin, huudetu, huusi, huusime, huusin, huusit, huusitte, huusivat, huutaa, huutaakseen, huutaen, huutaessa, huutaisi, huutaisimme, huutaisin, huutaisit, huutaisitte, huutaisivat, huutakaa, huutakaamme, huutako, huutakoon, huutakoot, huutama, huutamaan, huutamaisillaan, huutamalla, huutaman, huutamassa, huutamasta, huutamaton, huutamatta, huutaminen, huutamista, huutane, huutaneet, huutanemme, huutanen, huutanet, huutanette, huutanevat, huutanut, huutava, huutavat* — “shout” (54);

*souda, soudamme, soudan, soudat, soudatte, soudeta, soudetaan, soudettaessa, soudettaisi, soudettaisiin, soudettako, soudettakoon, soudettaman, soudettane, soudettaneen, soudettava, soudettiin, soudettu, soudimme, soudin, soudit, souditte, sousi, sousimme, sousin, sousit, sousitte, sousivat, soutaa, soutaakseen, soutaen, soutaessa, soutaisi, soutaisimme, soutaisin, soutaisit, soutaisitte, soutaisivat, soutakaa, soutakaamme, soutako, soutakoon, soutakoot, soutama, soutamaan, soutamaisillaan, soutamalla, soutaman, soutamassa, soutamasta, soutamaton, soutamatta, soutaminen, soutamista, soutane, soutanee, soutaneet, soutanemme, soutanen, soutanet, soutanette, soutanevat, soutanut, soutava, soutavat, souti, soutivat* — “row” (55);

*kaiva, kaivaa, kaivaisi, kaivaisimme, kaivaisin, kaivaisit, kaivaisitte, kaivaisivat, kaivakaa, kaivakaamme, kaivako, kaivakoon, kaivakoot, kaivamme, kaivan, kaivaneet, kaivanut, kaivat, kaivatte, kaivavat, kaiveta, kaivetaan, kaivettaisi, kaivettaisiin, kaivettako, kaivettakoon, kaivettiin, kaivettu, kaivoi, kaivoimme, kaivoin, kaivoit, kai-voitte, kaivoivat* — “dig” (56);

*saarra, saarramme, saarran, saarrat, saarratte, saarreta, saarrettaan, saarrettaessa, saarrettaisi, saarrettaisiin, saarrettako, saarrettakoon, saarrettaman, saarrettane, saarrettaneen, saarrettava, saarrettiin, saarrettu, saar-roiimme, saarroin, saarroit, saarroitte, saarsi, saarsimme, saarsin, saarsit, saarsitte, saarsivat, saartaa, saar-taakseen, saartaen, saartaessa, saartaisi, saartaisimme, saartaisin, saartaisit, saartaisitte, saartaisivat, saarta-kaa, saartakaamme, saartako, saartakoon, saartakoot, saartama, saartamaan, saartamaisillaan, saartamalla, saartaman, saartamassa, saartamasta, saartamaton, saartamatta, saartaminen, saartamista, saartane, saarta-nee, saartaneet, saartanemme, saartanen, saartanet, saartanette, saartanevat, saartanut, saartava, saartavat, saarto, saartoivat* — “surround” (57);

*laske, laskea, laskeakseen, laskee, laskekaa, laskekaamme, laskeko, laskekoon, laskekoot, laskema, laskemaan, laskemaisillaan, laskemalla, laskeman, laskemassa, laskemasta, laskematon, laskematta, laskeminen, laskemista, laskemme, lasken, laskene, laskenee, laskeneet, laskenemme, laskenen, laskenet, laskenette, laskenevat, laskenut, lasket, lasketa, lasketaan, laskettaessa, laskettaisi, laskettaisiin, laskettako, laskettakoon, laskettaman, laskettane, laskettaneen, laskettava, laskette, laskettiin, laskettu, laskeva, laskevat, laski, laskien, laskiessa, laskimme, laskin, laskisi, laskisimme, laskisin, laskisit, laskisitte, laskisivat, laskit, laskitte, laskivat* — “drop” (58);

*tunne, tunnemme, tunnen, tunnet, tunnetta, tunnetaan, tunnettaessa, tunnettaisi, tunnettaisiin, tunnettako, tunnet-takoon, tunnettaman, tunnettane, tunnettaneen, tunnettava, tunnette, tunnettiin, tunnettu, tunsi, tunsimme, tunsin, tunsit, tunsitte, tunsivat, tunea, tuneakseen, tunee, tunteka, tuntekaamme, tunteko, tuntekoon, tuntekoot, tun-tema, tuntemaan, tuntemaisillaan, tuntemalla, tunteman, tuntemassa, tuntemasta, tuntematon, tuntematta, tunteminen, tuntemista, tuntene, tuntenee, tunteneet, tuntenemme, tuntenen, tuntenet, tuntenette, tuntenevat, tuntenut, tunteva, tuntevat, tuntien, tuntiessa, tuntisi, tuntisimme, tuntisin, tuntisit, tuntisitte, tuntisivat* — “feel” (59);

*lähde, lähdemme, lähdens, lähdet, lähdette, lähdettiin, lähdetty, lähdettäässä, lähdettäisi, lähdettäisiin, lähdettäkö, lähdettäköön, lähdettämän, lähdettäne, lähdettäneen, lähdettävä, lähdetä, lähdetään, lähdimmme, lähdin, lähdit, lähditte, lähtee, lähtekää, lähtekäämme, lähtekö, lähteköön, lähteköt, lähteminen, lähtemistä, lähtemä, lähte-mäisillään, lähtemällä, lähtemän, lähtemässä, lähtemästä, lähtemättä, lähtemätön, lähtemään, lähtene, lähtenee, lähteneet, lähtenemme, lähtenen, lähtenet, lähtenette, lähtenevät, lähtenyt, lähtevä, lähtevät, lähteä, lähteäkseen, lähti, lähtien, lähtiessä, lähtisi, lähtisimme, lähtisin, lähtisit, lähtisitte, lähtisivät, lähtivät* — “leave” (60);

*salli, sallia, salliakseen, sallien, sallissa, sallii, sallikaamme, salliko, sallikoon, sallima, salli-maan, sallimaisillaan, sallimalla, salliman, sallimassa, sallimasta, sallimaton, sallimatta, sallimin, sallimista, sallimme, sallin, salline, sallinee, sallineet, sallinemme, sallinen, sallinet, sallinette, sallinevat, sallinut, sallisi, sallisimme, sallisim, sallisit, sallisitte, sallisivat, sallit, sallita, sallitaan, sallittaessa, sallittaisi, sallittaisiin, sal-littako, sallittakoon, sallittaman, sallittane, sallittaneen, sallittava, sallitte, sallittiin, sallitu, salliva, sallivat* — “allow” (61);

*voi, voida, voidaan, voidakseen, voiden, voidessa, voikaa, voikaamme, voiko, voikoon, voikoot, voima, voimaan, voimaisillaan, voimalla, voiman, voimassa, voimasta, voimatton, voimatta, voiminen, voimista, voimme, voin, voi-ne, voinee, voineet, voinemme, voinen, voinet, voinette, voinevat, voinut, voisi, voisimme, voisim, voisit, voisit,*

*voisivat, voit, voitaessa, voitaisi, voitaistin, voitako, voitakoon, voitaman, voitane, voitaneen, voitava, voitiin, voitte, voitu, voiva, voivat* — “can” (62);

*saa, saada, saadaan, saadakseen, saaden, saadessa, saakaa, saakaamme, saako, saakoon, saakoot, saama, saamaan, saamaisillaan, saamalla, saaman, saamassa, saamasta, saamaton, saamatta, saaminen, saamista, saame, saan, saane, saanee, saaneet, saanemme, saanen, saanet, saanette, saanevat, saanut, saat, saataessa, saataisi, saataisiin, saatako, saatakoon, saataman, saatane, saataneen, saatava, saatiin, saatte, saatu, saava, saavat* — “receive” (63);

*juo, juoda, juodaan, juodakseen, juoden, juodessa, juokaa, juokaamme, juoko, juokoon, juokoot, juoma, juomaan, juomaisillaan, juomalla, juoman, juomassa, juomasta, juomaton, juomatta, juominen, juomista, juomme, juon, juone, juonee, juoneet, juonemme, juonen, juonet, juonette, juonevat, juonut, juot, juotaessa, juotaisi, juotaisiin, juotako, juotakoon, juotaman, juotane, juotaneen, juotava, juotiin, juotte, juotu, juova, juovat* — “drink” (64);

*kävi, kävime, kävin, kävisi, kävisimme, kävisin, kävisitte, kävisivät, kävit, kävitte, kävivät, käy, käyden, käydessä, käydä, käydäkseen, käydään, käykää, käykäämme, käykö, käyköön, käyköt, käyminen, käymistä, käyme, käymä, käymäisillään, käymällä, käymän, käymässä, käymästä, käymättä, käymätön, käymäään, käyn, käyne, käyneet, käynemme, käynen, käynet, käynette, käynevat, käynyt, käyt, käytiin, käytte, käyty, käytäässä, käytäisi, käytäisiin, käytäkö, käytämän, käytäne, käytäneen, käytävä, käyvä, kävät* — “go” (65);

*rohkaise, rohkaisee, rohkaisema, rohkaisemaan, rohkaisemillaan, rohkaisemalla, rohkaiseman, rohkaisemasa, rohkaisemasta, rohkaisematon, rohkaisematta, rohkaiseminen, rohkaisemista, rohkaisemme, rohkaisen, rohkaiset, rohkaisette, rohkaiseva, rohkaisevat, rohkaisi, rohkaisimme, rohkaisin, rohkaisisi, rohkaisisimme, rohkaisisin, rohkaisisit, rohkaisisitte, rohkaisisivat, rohkaisit, rohkaisitte, rohkaisivat, rohkaiska, rohkaiskaamme, rohkaisko, rohkaiskoon, rohkaiskoot, rohkaisse, rohkaissee, rohkaisseet, rohkaissemme, rohkaissen, rohkaisset, rohkaissette, rohkaissevat, rohkaissut, rohkaista, rohkaistaan, rohkaistaessa, rohkaistaasi, rohkaistaasiin, rohkaista-ko, rohkaistakoon, rohkaistakseen, rohkaistaman, rohkaistane, rohkaistaneen, rohkaistava, rohkaisten, rohkaistessa, rohkaistiin, rohkaistu* — “encourage” (66);

*tule, tulee, tulema, tulemaan, tulemaisillaan, tulemallla, tuleman, tulemassa, tulemasta, tulematon, tulematta, tuleminen, tulemista, tulemme, tulen, tulet, tulette, tuleva, tulevat, tuli, tulimme, tulin, tulisi, tulismme, tulisin, tulisit, tulisitte, tulisivat, tulit, tulitte, tulivat, tulkaa, tulkaamme, tulko, tulkoon, tulkoot, tulla, tullaan, tullakseen, tulle, tullee, tulleet, tullemme, tullen, tullessa, tullet, tullette, tullevat, tullut, tultaessa, tultaisi, tultaisiin, tultako, tultakoon, tultaman, tultane, tultaneen, tultava, tultiin, tultu* — “come” (67);

*tupakoi, tupakoida, tupakoidaan, tupakoidakseen, tupakoiden, tupakoidessa, tupakoikaa, tupakoikaamme, tupakoiko, tupakoikoon, tupakoikoot, tupakoima, tupakoimaan, tupakoimaisillaan, tupakoimalla, tupakoiman, tupakoimassa, tupakoimasta, tupakoimat, tupakoimatta, tupakoiminen, tupakoimista, tupakoimme, tupakoin, tupakoine, tupakoinee, tupakoineet, tupakoinemme, tupakoinen, tupakoinet, tupakoinette, tupakoinevat, tupakoinut, tupakoisi, tupakoisimme, tupakoisin, tupakoisit, tupakoisitte, tupakoisivat, tupakoit, tupakoitaessa, tupakoitaisi, tupakoitaisiin, tupakoitako, tupakoitakoon, tupakoitaman, tupakoitane, tupakoitanen, tupakoitava, tupakoitiin, tupakoitse, tupakoitsee, tupakoitsema, tupakoitsemaan, tupakoitsemillaan, tupakoitsemalla, tupakoitseman, tupakoitsemassa, tupakoitsemasta, tupakoitsematon, tupakoitsematta, tupakoitseminen, tupakoitsemista, tupakoitsemme, tupakoitsen, tupakoitset, tupakoitsette, tupakoitseva, tupakoitsevat, tupakoitsi, tupakoitsimme, tupakoit-sin, tupakoitsisi, tupakoitsisimme, tupakoitsisin, tupakoitsisit, tupakoitsisitte, tupakoitsisivat, tupakoitsit, tupakoit-sitte, tupakoitsivat, tupakoitte, tupakoitu, tupakoiva, tupakoivat* — “smoke” (68);

*valinne, valinnee, valinneet, valinnemme, valinnen, valinnet, valinnette, valinnevat, valinnut, valita, valitaan, valitakseen, valiten, valitessa, valitkaa, valitkaamme, valitko, valitkoon, valitkoot, valitsee, valitsema, valit-seaan, valitsemaisillaan, valitsemalla, valitseman, valitsemassa, valitsemasta, valitsematon, valitsematta, valit-seminen, valitsemista, valitsemme, valitsen, valitset, valitsette, valitseva, valitsevat, valitsi, valitsimme, valitsin, valitsisi, valitsisimme, valitsisin, valitsisit, valitsisitte, valitsisivat, valitsit, valitsitte, valitsivat, valittaessa, valit-taisi, valittaisiin, valittako, valittakoon, valittaman, valittane, valittaneen, valittava, valittiin, valittu* — “choose” (69);

*juoksee, juoksee, juoksema, juoksemaan, juoksemaisillaan, juoksemalla, juokseman, juoksemassa, juoksemasta, juoksematon, juoksematta, juokseminen, juoksemista, juoksemme, juoksen, juokset, juoksette, juokseva, juokse-vat, juoksi, juoksimme, juoksin, juoksisi, juoksisimme, juoksisin, juoksisit, juoksisitte, juoksisivat, juoksit, juok-sitte, juoksivat, juoskaa, juoskaamme, juosko, juoskoon, juoskoot, juosse, juossee, juosseet, juossemme, juossen, juosset, juossette, juossevat, juossut, juosta, juostaan, juostaessa, juostaisi, juostaisiin, juostako, juostakoon, juostakseen, juostaman, juostane, juostaneen, juostava, juosten, juostessa, juostiin, juostu* — “run” (70);

*näe, näemme, näen, näet, näette, nähden, nähdessä, nähdä, nähdäkseen, nähdään, nähkää, nähkäämme, nähkö, nähköön, nähkööt, nähne, nähnee, nähheet, nähnemme, nähnen, nähnet, nähnette, nähnevät, nähnyt, nähtiin, nähty, nähtäässä, nähtäisi, nähtäisiin, nähtäkö, nähtäköön, nähtämän, nähtäne, nähtäneen, nähtävä, näimme, nän, nät, näitte, näkee, näkeminen, näkemistä, näkemä, näkemäisillään, näkemällä, näkemän, näkemässä, näkemästä, näkemättä, näkemätön, näkemää, näkevät, näkevät, näki, näkisi, näkisimme, näkisin, näkisit, näkisitte, näkisivät, näkivät — “see” (71);*

*vanhene, vanhenee, vanhenema, vanhenemaa, vanhenemaisillaan, vanhenemalla, vanheneman, vanhenemassa, vanhenemasta, vanhenematon, vanhenematta, vanheneminen, vanhenemista, vanhenemme, vanhenen, vanhenet, vanhenette, vanheneva, vanhenevat, vanheni, vanhenimme, vanhenin, vanhenisi, vanhenisimme, vanhenisin, vanhenisit, vanhenisitte, vanhenisivat, vanhenit, vanhenivat, vanhenne, vanheneet, vanhenemme, vanhennen, vanhennet, vanhennette, vanhenevat, vanhennut, vanheta, vanhetaan, vanhetakseen, vanheten, vanhetessa, vanhetkaa, vanhetkaamme, vanhetko, vanhetkoon, vanhetkoot, vanhettaessa, vanhettaisi, vanhettaisiin, vanhettako, vanhettakoon, vanhettaman, vanhettane, vanhettaneen, vanhettava, vanhettiin, vanhettu — “age” (72);*

*salaal, salaama, salaamaan, salaamaisillaan, salaamalla, salaaman, salaamassa, salaamasta, salaamaton, salaamatta, salaaminen, salaamista, salaamme, salaan, salaatt, salaava, salaavat, salaisi, salaisimme, salaisin, salaisit, salaisitte, salaisivat, salanne, salannee, salanneet, salannemme, salannen, salannet, salannette, salannevat, salannut, salasi, salasimme, salasin, salasit, salasitte, salasivat, salata, salataan, salatakseen, salaten, salatessa, salatkaa, salatkaamme, salatko, salatkoon, salatkoot, salattaessa, salattaisi, salattaisiin, salattako, salattakoon, salattaman, salattane, salattaneen, salattava, salattiin, salattu — “conceal” (73);*

*katkea, katkeaa, katkeaisi, katkeaisimme, katkeaisin, katkeaisitte, katkeaisivat, katkeama, katkeamaan, katkeamaisillaan, katkeamalla, katkeaman, katkeamassa, katkeamasta, katkeamaton, katkeamatta, katkeaminen, katkeamista, katkeamme, katkean, katkeat, katkeatte, katkeava, katkeavat, katkeisi, katkeisimme, katkeisin, katkeisit, katkeisitte, katkeisivat, katkenne, katkennee, katkenneet, katkennemme, katkennen, katkennet, katkennette, katkennevät, katkennut, katkesi, katkesimme, katkesin, katkesit, katkesivat, katketa, katketaan, katketseen, katketen, katketessa, katketkaa, katketkaamme, katketko, katketkoon, katketkoot, katkettaessa, katkettaisi, katkettaisiin, katkettako, katkettakoon, katkettaman, katkettane, katkettaneen, katkettava, katkettiin, katkettu — “cut” (74);*

*selvinne, selvinnee, selvinneet, selvinnemme, selvinnen, selvinnet, selvinnette, selvinnevät, selvinnyt, selvisi, selvisimme, selvisin, selvisit, selvisitte, selvisivät, selviten, selvitessä, selvitkää, selvitkäämme, selvitkö, selvitköön, selvitkööt, selvittiin, selvitty, selvittäässä, selvittäisi, selvittäisiin, selvittäkö, selvittäköön, selvittämän, selvittäne, selvittäneen, selvittävä, selvitä, selvitökseen, selvitän, selviä, selviäisi, selviäsimme, selviäisin, selviäisit, selviäisitte, selviäisivät, selviäminen, selviämistä, selviämme, selviämä, selviämäisillään, selviämällä, selviämän, selviämässä, selviämästä, selviämättä, selviämätön, selviämää, selviän, selviät, selviätte, selviävä, selviävät, selviää — “escape” (75);*

*taida, taidamme, taidan, taidat, taidatte, taideta, taidetaan, taidettaa, taidettaessa, taidettaisi, taidettaisiin, taidettako, taidettakoon, taidettaman, taidettane, taidettaneen, taidettava, taidettiin, taidetu, tainne, tainnee, tainneet, tainnemme, tainnen, tainnet, tainnette, tainnevät, tainnut, taisi, taisimme, taisin, taisit, taisitte, taisivat, taitaa, taitakseen, taitaan, taitaessa, taitaisi, taitaisimme, taitaisin, taitaisit, taitaisitte, taitaisivat, taitaka, taitakaamme, taitako, taitakoon, taitakoot, taitama, taitamaan, taitamaisillaan, taitamalla, taitaman, taitamassa, taitamasta, taitamaton, taitamatta, taitaminen, taitamista, taitane, taitanee, taitaneet, taitanemme, taitanen, taitanet, taitanette, taitanevat, taitanut, taitava, taitavat — “master” (76);*

*kumajaa, kumajaisi, kumaji — “boom” (77);*

*kaikaa, kaikaavat, kaikaisi, kaikaisivat — “echo” (78),*

we obtain

{*sanomaton*} — (52) 1 word;

{*sano, sanoa, sanoakseen, sanoen, sanoessa, sanoi, sanoimme, sanoin, sanoisi, sanoisimme, sanoisin, sanoisit, sanoisitte, sanoisivat, sanoit, sanoitte, sanoivat, sanokaa, sanokaamme, sanoko, sanokoon, sanokoot, sanoma, sanomaan, sanomaisillaan, sanomalla, sanoman, sanomassa, sanomasta, sanomatta, sanominen, sanomista, sanomme, sanon, sanone, sanonee, sanoneet, sanonemme, sanonen, sanonet, sanonette, sanonevat, sanonut, sanoo, sanot, sanota, sanotaan, sanottaessa, sanottaisi, sanottaisiin, sanottako, sanottakoon, sanottaman, sanottane, sanottaneen, sanottava, sanotte, sanottiin, sanottu, sanova, sanovat*} — (52) 61 words;

{*muistamaton*} — (53) 1 word;

{*muista, muistaa, muistaakseen, muistaen, muistaessa, muistaisi, muistaisimme, muistaisin, muistaisit, muistaisitte, muistaisivat, muistakaa, muistakaamme, muistako, muistakoon, muistakoot, muistama, muistamaan, muistamillaan, muistamalla, muistaman, muistamassa, muistamasta, muistamatta, muistaminen, muistamista, muistamme, muistan, muistane, muistanee, muistaneet, muistanemme, muistanan, muistonet, muistanne, muistanevat, muistananut, muistat, muistatte, muistava, muistavat, muisteta, muistetaan, muistettaessa, muistettai, muistettaiin, muistettako, muistettakoon, muistettaman, muistettane, muistettaneen, muistettava, muistettiin, muistettu, muisti, muistimme, muistin, muistit, muistite, muistivat*} — (53) 60 words;

{*huutamaton*} — (54) 1 word;

{*huuda, huudamme, huudan, huudat, huudatte, huudeta, huudetaan, huudettaessa, huudettai, huudettaiin, huudettako, huudettakoon, huudettaman, huudettane, huudettaneen, huudettava, huudettiin, huudetu, huusi, huusime, huusin, huusit, huusitte, huusivat, huutaa, huutaakseen, huutaen, huutaessa, huutaisi, huutaisimme, huutaisin, huutaisit, huutaisitte, huutaisivat, huutakaa, huutakaamme, huutako, huutakoon, huutakoot, huutama, huutamaan, huutamillaan, huutamalla, huutaman, huutamassa, huutamasta, huutamatta, huutaminen, huutamista, huutane, huutanee, huutaneet, huutanemme, huutanen, huutanet, huutanette, huutanevat, huutanut, huutava, huutavat*} — (54) 60 words;

{*sousitte*} — (55) 1 word;

{*soutamaton*} — (55) 1 word;

{*souda, soudamme, soudan, soudat, soudatte, soudeta, soudetaan, soudettaessa, soudettai, soudettaiin, soudettako, soudettakoon, soudettaman, soudettane, soudettaneen, soudettava, soudettiin, soudettu, soudimme, soudin, soudit, souditte, sousi, sousimme, sousin, sousit, sousivat, soutaa, soutaakseen, soutaen, soutaessa, soutaisi, soutaisimme, soutaisin, soutaisit, soutaisitte, soutaisivat, soutakaa, soutakaamme, soutako, soutakoon, soutakoot, soutama, soutamaan, soutamillaan, soutamalla, soutaman, soutamassa, soutamasta, soutamatta, soutaminen, soutamista, soutane, soutanee, soutaneet, soutanemme, soutanen, soutanet, soutanette, soutanevat, soutanut, soutava, soutavat, souti, soutivat*} — (55) 65 words;

{*kaiva, kaivaa, kaivaisi, kaivaisimme, kaivaisin, kaivaisit, kaivaisitte, kaivaisivat, kaivakaa, kaivakaamme, kai-vako, kaivakoon, kaivakoot, kaivamme, kaivan, kaivaneet, kaivanut, kaivat, kaivatte, kaivavat, kaiveta, kaivetaan, kaivettai, kaivettaiin, kaivettako, kaivettakoon, kaivettiin, kaivettu, kaivoi, kaivoimme, kaivoin, kaivoit, kai-vitte, kaivoivat*} — (56) 34 words;

{*saartamaton*} — (57) 1 word;

{*saarra, saarramme, saarran, saarrat, saarratte, saarreta, saarretaan, saarrettaessa, saarrettaisi, saarrettaisiin, saarrettako, saarrettakoon, saarrettaman, saarrettane, saarrettaneen, saarrettava, saarrettiin, saarrettu, saar-roimme, saarroin, saarroit, saarroitte, saarsi, saarsimme, saarsin, saarsit, saarsivat, saartaa, saartaakseen, saartaen, saartaessa, saartaisi, saartaisimme, saartaisin, saartaisit, saartaisitte, saartaisivat, saarta-kaa, saartakaamme, saartako, saartakoon, saartakoot, saartama, saartamaan, saartamillaan, saartamalla, saartaman, saartamassa, saartamasta, saartamatta, saartaminen, saartamista, saartane, saartanee, saartaneet, saartanenne, saartanen, saartanet, saartanette, saartanevat, saartanut, saartava, saartavat, saarto, saartoivat*} — (57) 66 words;

{*laskematon*} — (58) 1 word;

{*laske, laskea, laskeakseen, laskee, laskekaa, laskekamme, laskeko, laskekoon, laskekoot, laskema, laskemaan, laskemaisillaan, laskemalla, laskeman, laskemassa, laskemasta, laskematta, laskeminen, laskemista, laskemme, lasken, laskene, laskenee, laskeneet, laskenemme, laskenen, laskenet, laskenette, laskenevat, laskenut, lasket, laska-ta, lasketaan, laskettaessa, laskettai, laskettaiin, laskettako, laskettakoon, laskettaman, laskettane, laskettainen, laskettava, laskette, laskettiin, laskettu, laskeva, laskevat, laski, laskien, laskiessa, laskimme, laskin, laskisi, laskisimme, laskisin, laskisit, laskisitte, laskisivat, laskit, laskitte, laskivat*} — (58) 61 words;

{*tuntematon*} — (59) 1 word;

{*tunne, tunnemme, tunnen, tunnet, tunneta, tunnetaan, tunnettaessa, tunnettaisi, tunnettaisiin, tunnettako, tunnettakoon, tunnettaman, tunnettane, tunnettaneen, tunnettava, tunnette, tunnettiin, tunnettu, tarsi, tunsimme, tunsin, tunsit, tunsitte, tunsivat, tuntea, tunteakseen, tuntee, tuntekaa, tuntekaamme, tunteko, tuntekoon, tuntekoot, tunte-ma, tunteamaan, tuntemaisillaan, tuntemalla, tunteman, tuntemassa, tuntemasta, tuntematta, tunteminen, tuntemis-ta, tuntene, tuntenee, tunteneet, tunteenemme, tundenen, tuntenet, tuntenette, tuntenevat, tuntenut, tunteva, tuntevat, tuntien, tuntiessa, tuntisi, tuntisimme, tuntisin, tuntisit, tuntisivat*} — (59) 61 words;

{*lähtemätön*} — (60) 1 word;

{*lähde, lähdemme, lähden, lähdet, lähdette, lähdettiin, lähdetty, lähdettäässä, lähdettäisi, lähdettäisiin, lähdet-täkö, lähdettäköön, lähdettämän, lähdettäne, lähdettäneen, lähdettävä, lähdetä, lähdetään, lähdimmme, lähdin, lähdit, lähditte, lähee, lähtekää, lähtekäämme, lähtekö, lähteköön, lähtekööt, lähteminen, lähtemistä, lähtemä, lähtemäisillään, lähtemällä, lähtemän, lähtemässä, lähtemästä, lähtemättä, lähtemään, lähtene, lähtenee, lähte-neet, lähtenemme, lähtenen, lähtenet, lähtenette, lähtenevat, lähtenyt, lähtevä, lähtevät, lähteä, lähtekseen, lähti, lähtien, lähtiessä, lähtisi, lähtisimme, lähtisin, lähtisit, lähtisitte, lähtisivät, lähtivät*} — (60) 61 words;

{*sallimaton*} — (61) 1 word;

{*salli, sallia, salliaseen, sallien, salliessa, sallii, sallikaa, sallikaamme, salliko, sallikoon, sallikoot, sallima, sal-limaan, sallimaisillaan, sallimalla, salliman, sallimassa, sallimasta, sallimatta, salliminen, sallimista, sallimme, sallin, salline, sallinee, sallineet, sallinemme, sallinen, sallinet, sallinette, sallinevat, sallinut, sallisi, sallisimme, sallis-in, sallisit, sallisitte, sallisivat, sallit, sallita, sallitaan, sallitaessa, sallitaisi, sallitaisiin, sallittako, sal-littakoon, sallittaman, sallittane, sallittaneen, sallitava, sallitte, sallittiin, sallittu, salliva, sallivat*} — (61) 55 words;

{*voimaton*} — (62) 1 word;

{*voi, voidaan, voidakseen, voiden, voidessa, voikaa, voikaamme, voiko, voikoon, voikoot, voima, voimaan, voimaisillaan, voimalla, voiman, voimassa, voimasta, voimatta, voimin, voimista, voimme, voin, voine, voinee, voineet, voinemme, voinen, voinet, voinette, voinevat, voinut, voisi, voisimme, voisin, voisit, voisitte, voisivat, voit, voitaessa, voitaisi, voitaisiin, voitako, voitakoon, voitaman, voitane, voitaneen, voitava, voitiin, voitte, voitu, voiva, voivat*} — (62) 53 words;

{*saamaton*} — (63) 1 word;

{*saataisiin, saataman, saatane, saataneen, saatiin, saatu*} — (63) 6 words;

{*saa, saada, saadaan, saadakseen, saaden, saadessa, saakaa, saakaamme, saako, saakoon, saakoot, saama, saa-maan, saamaisillaan, saamalla, saaman, saamassa, saamasta, saamatta, saaminen, saamista, saamme, saan, saane, saanee, saaneet, saanemme, saanen, saanet, saanette, saanevat, saanut, saat, saataessa, saataisi, saatako, saatakoon, saatava, saatte, saava, saavat*} — (63) 41 words;

{*juomaton*} — (64) 1 word;

{*juo, juoda, juodaan, juodakseen, juoden, juodessa, juokaa, juokaamme, juoko, juokoon, juokoot, juoma, juo-maan, juomaisillaan, juomalla, juoman, juomassa, juomasta, juomatta, juominen, juomista, juomme, juon, juone, juonee, juoneet, juonemme, juonen, juonet, juonette, juonevat, juonut, juot, juotaessa, juotaisi, juotaisiin, juotako, juotakoon, juotaman, juotane, juotaneen, juotava, juotiin, juotte, juotu, juova, juovat*} — (64) 47 words;

{*käymätön*} — (65) 1 word;

{*kävi, kävimme, kävin, kävisi, kävisimme, kävisin, kävisit, kävisitte, kävisivät, kävit, kävitte, kävivät, käy, käy-den, käydessä, käydä, käydäseen, käydään, käykää, käykäämme, käykö, käyköön, käykööt, käyminen, käymistä, käymme, käymä, käymaisillaan, käymällä, käymän, käymässä, käymästä, käymäään, käyn, käyne, käy-nee, käyneet, käynemme, käynen, käynet, käynette, käynevat, käynt, käyt, käytiin, käytte, käyty, käytäessä, käy-täisi, käytäisiin, käytäkö, käytäköön, käytämän, käytäne, käytäneen, käytävä, käyvä, käyvät*} — (65) 59 words;

{*rohkaisematon*} — (66) 1 word;

{rohkaise, rohkaisee, rohkaisema, rohkaisemaan, rohkaisemillaan, rohkaisemalla, rohkaiseman, rohkaisemasa, rohkaisemasta, rohkaisematta, rohkaiseminen, rohkaisemista, rohkaisemme, rohkaisen, rohkaiset, rohkaisette, rohkaiseva, rohkaisevat, rohkaisi, rohkaisimme, rohkaisin, rohkaisisi, rohkaisisimme, rohkaisisin, rohkaisosit, rohkaisisitte, rohkaisisivat, rohkaisit, rohkaisitte, rohkaisivat, rohkaiska, rohkaiskaamme, rohkaisko, rohkaiskoon, rohkaiskoot, rohkaisse, rohkaissee, rohkaisseet, rohkaisemme, rohkaissen, rohkaisset, rohkaissette, rohkaissevat, rohkaissut, rohkaista, rohkaistaan, rohkaistaessa, rohkaistaisi, rohkaistaisiin, rohkaistako, rohkaistakoon, rohkaistakseen, rohkaistaman, rohkaistane, rohkaistaneen, rohkaistava, rohkaisten, rohkaistessa, rohkaistiin, rohkaistu} — (66) 60 words;

{tulematon} — (67) 1 word;

{tule, tulee, tulema, tulemaan, tulemaillaan, tulemalla, tuleman, tulemassa, tulemasta, tulematta, tuleminen, tulemista, tulemme, tulen, tulet, tulette, tuleva, tulevat, tuli, tulimme, tulin, tulisi, tulismme, tulisin, tulisit, tulisitte, tulisivat, tulit, tulitte, tulivat, tulkaa, tulkaamme, tulko, tulkoon, tulkoot, tulla, tullaan, tullakseen, tulle, tulle, tulleet, tulemme, tullen, tullessa, tullet, tullette, tullevat, tullut, tultaessa, tultaisi, tultaisiin, tultako, tultakoon, tultaman, tultane, tultaneen, tultava, tultiin, tultu} — (67) 59 words;

{tupakoimaton, tupakoitsematon} — (68) 2 words;

{tupakoi, tupakoida, tupakoidaan, tupakoidakseen, tupakoiden, tupakoidessa, tupakoikaa, tupakoikaamme, tupakoiko, tupakoikoon, tupakoikoot, tupakoima, tupakoimaan, tupakoimaisillaan, tupakoimalla, tupakoiman, tupakoimassa, tupakoimasta, tupakoimatta, tupakoimin, tupakoimista, tupakoimme, tupakoin, tupakoine, tupakoinee, tupakoineet, tupakoinemme, tupakoinen, tupakoinet, tupakoinette, tupakoinevat, tupakoinut, tupakoisi, tupakoisimme, tupakoisin, tupakoisit, tupakoisitte, tupakoisivat, tupakoit, tupakoitaessa, tupakoitaisi, tupakoitaisiin, tupakoitako, tupakoitakoon, tupakoitaman, tupakoitane, tupakoitaneen, tupakoitava, tupakoitiin, tupakoitse, tupakoitsee, tupakoitsema, tupakoitsemaan, tupakoitsemaillaan, tupakoitsemalla, tupakoitseman, tupakoitsemasa, tupakoitsemasta, tupakoitsematta, tupakoitseminen, tupakoitsemista, tupakoitsemme, tupakoitsen, tupakoitset, tupakoitsette, tupakoitseva, tupakoitsevat, tupakoitsi, tupakoitsimme, tupakoitsin, tupakoitsisi, tupakoitsisimme, tupakoitsisin, tupakoitsisit, tupakoitsisitte, tupakoitsisivat, tupakoitsit, tupakoitsitte, tupakoitsivat, tupakoitte, tupakoitu, tupakoiva, tupakoivat} — (68) 83 words;

{valitsematon} — (69) 1 word;

{valinne, valinnee, valinneet, valinnemme, valinnen, valinnet, valinnette, valinnevat, valinnut, valita, valitaan, valitakseen, valiten, valitessa, valitkaa, valitkaamme, valitko, valitkoon, valitkoot, valitse, valitsee, valitsema, valitsemaan, valitsemaillaan, valitsemalla, valitseman, valitsemassa, valitsemasta, valitsematta, valitseminen, valitsemista, valitsemme, valitsen, valitset, valitsette, valitseva, valitsevat, valitsi, valitsimme, valitsin, valitsisi, valitsisimme, valitsisin, valitsisit, valitsisitte, valitsisivat, valitsit, valitsitte, valitsivat, valittaessa, valittaisi, valitaisiin, valittako, valittakoon, valittaman, valittane, valittaneen, valittava, valittiin, valittu} — (69) 60 words;

{juoksematon} — (70) 1 word;

{juokse, juoksee, juoksema, juoksemaan, juoksemaillaan, juoksemalla, juokseman, juoksemassa, juoksesta, juoksematta, juokseminen, juoksemista, juoksemme, juoksen, juokset, juoksette, juokseva, juoksevat, juoksi, juoksimme, juoksin, juoksisi, juoksisimme, juoksisin, juoksisit, juoksisitte, juoksisivat, juoksit, juoksitte, juoksivat, juoskaa, juoskaamme, juosko, juoskoon, juoskoot, juosse, juossee, juosseet, juossemme, juossen, juosset, juossette, juossevat, juossut, juosta, juostaan, juostaessa, juostaisi, juostaisiin, juostako, juostakoon, juostakseen, juostaman, juostane, juostaneen, juostava, juosten, juostessa, juostiin, juostu} — (70) 60 words;

{näkemätön} — (71) 1 word;

{näe, näemme, näen, näet, näette, näden, nähdessä, nähdä, nähdäkseen, nähdään, nähkää, nähkääämme, nähkö, nähköön, nähköt, nähne, nähnee, nähneet, nähnemme, nähnen, nähnet, nähnette, nähnevät, nähnyt, nähtiin, nähty, nähtäässä, nähtäisi, nähtäisiin, nähtäkö, nähtäköön, nähtämän, nähtäne, nähtäneen, nähtävä, näimme, näin, näit, näitte, näkee, näkeminen, näkemistä, näkemä, näkemäisillään, näkemällä, näkemän, näkemässä, näkemästä, näkemättä, näkemään, näkevä, näkevät, näki, näkisi, näkisimme, näkisin, näkisit, näkisitte, näkisivät, näkivät} — (71) 60 words;

{vanhenematon} — (72) 1 word;

{vanhene, vanhenee, vanhenema, vanhenemaan, vanhenemaisillaan, vanhenemalla, vanheneman, vanhenemas-  
sa, vanhenemasta, vanhenematta, vanheneminen, vanhenemista, vanhenemme, vanhenen, vanhenet, vanhenet-  
te, vanheneva, vanhenevat, vanheni, vanhenimme, vanhenin, vanhenisi, vanhenisimme, vanhenisin, vanhenisit,  
vanhenisitte, vanhenisivat, vanhenit, vanhenitte, vanhenivat, vanhenne, vanhennee, vanhenneet, vanhennemme,  
vanhennen, vanhennet, vanhennette, vanhenevat, vanhennut, vanheta, vanhetaan, vanhetakseen, vanheteen, van-  
hetessa, vanhetkaa, vanhetkaamme, vanhetko, vanhetkoon, vanhetkoot, vanhettaessa, vanhettaisi, vanhettaisiin,  
vanhettako, vanhettakoon, vanhettaman, vanhettane, vanhettaneen, vanhettava, vanhettiin, vanhettu} — (72) 60  
words;

{salaamaton} — (73) 1 word;

{sala, salaama, salaamaan, salaamaisillaan, salaamalla, salaaman, salaamassa, salaamasta, salaamatta, sa-  
laaminen, salaamista, salaamme, salaan, salaatt, salaava, salaavat, salaisi, salaisimme, salaisin, salai-  
sit, salaisitte, salaisivat, salanne, salannee, salanneet, salannemme, salannen, salannet, salannette, salannevat,  
salannut, salasi, salasimme, salasin, salasit, salasitte, salasivat, salata, salataan, salatakseen, salaten, salatessa,  
salatkaa, salatkaamme, salatko, salatkoon, salatkoot, salattaessa, salattaisi, salattaisiin, salattako, salattakoon,  
salattaman, salattane, salattaneen, salattava, salattiin, salattu} — (73) 59 words;

{katkeamaton} — (74) 1 word;

{katkea, katkeaa, katkeaisi, katkeaisimme, katkeaisin, katkeaisit, katkeaisitte, katkeaisivat, katkeama, katke-  
maan, katkeamaisillaan, katkeamalla, katkeaman, katkeamassa, katkeamasta, katkeamatta, katkeaminen, katke-  
mista, katteamme, katkean, katkeat, katkeatte, katkeava, katkeavat, katkeisi, katkeisimme, katkeisin, katkeisit,  
katkeisitte, katkeisivat, katkenne, katkennee, katkenneet, katkennemme, katkennen, katkennet, katkennette, kat-  
kennevat, katkennut, katkesi, katkesimme, katkesin, katkesit, katkesitte, katkesivat, katketa, katketaan, katketak-  
seen, katketen, katketessa, katketkaa, katketkaamme, katketko, katketkoon, katketkoot, katkettasa, katkettaisi,  
katkettaisiin, katkettako, katkettakoon, katkettaman, katkettane, katkettaneen, katkettava, katkettiin, katkettu} —  
(74) 66 words;

{selviämätön} — (75) 1 word;

{selvinne, selvinnee, selvinneet, selvinnemme, selvinnen, selvinnet, selvinnette, selvinnevät, selvinnyt, selvisi, sel-  
visimme, selvisin, selvisit, selvisitte, selvisivät, selviten, selvitessä, selvitkää, selvitkäämme, selvitkö, selvitköön,  
selvitkööt, selvittiin, selvitty, selvittäässä, selvittäisi, selvittäisiin, selvittäkö, selvittäköön, selvittämän, selvittäne,  
selvittäneen, selvittävä, selvitä, selvitökseen, selvitään, selviä, selviäisi, selviäsimme, selviäisin, selviäisit, sel-  
viäisitte, selviäisivät, selviäminen, selviämistä, selviämme, selviämä, selviämäillä, selviämällä, selviämän,  
selviämässä, selviämästä, selviämättä, selviämään, selviän, selviät, selviätte, selviävä, selviävät, selviää} — (75)  
60 words;

{taitamaton} — (76) 1 word;

{taida, taidamme, taidan, taidat, taidatte, taideta, taidetaan, taidettaessa, taidettaisi, taidettaisiin, taidettako,  
taidettakoon, taidettaman, taidettane, taidettaneen, taidettava, taidettiin, taidetu, tainne, tainnee, tainneet, tain-  
nemme, tainnen, tainnet, tainnette, tainnevät, tainnut, taisi, taisimme, taisin, taisit, taisitte, taisivat, taitaa, tai-  
taakseen, taitaen, taitaessa, taitaisi, taitaisimme, taitaisin, taitaisit, taitaisitte, taitaisivat, taitakaa, taitakaamme,  
taitako, taitakoon, taitakoot, taitama, taitamaan, taitamaisillaan, taitamalla, taitamaan, taitamassa, taitamasta,  
taitamatta, taitaminen, taitamista, taitane, taitanee, taitaneet, taitanemme, taitanen, taitanet, taitanette, taitane-  
vat, taitanut, taitava, taitavat} — (76) 69 words;

{kumajaa, kumajaisi, kumaji} — (77) 3 words;

{kaikaa, kaikaavat, kaikaisi, kaikaisivat} — (78) 4 words.

Here, words ending in *-maton* and *mätön* are negative participles, which roughly translate into English words with prefixes *il-*, *im-*, *in-*, *ir-*, *non-*, *un-*, or suffix *-less*. To be consistent with our clustering policy for English words, we have separated negative participles from other inflected forms of Finnish verbs.

### 8.1.3 Heuristic detection of compounds

The following algorithm for heuristic detection of Finnish compounds differs from the Danish version (Algorithm 5.12) only in some specific details. To make the context clear, we still state the algorithm in full. (In what follows, the string minus operation  $\hat{\beta} \ominus \hat{\alpha}$  is prescribed by Definition 5.11.)

**Algorithm 8.13** (Heuristic identification of Finnish binary compounds). Let  $\Lambda^{\hat{\rho}} = \{\hat{\rho}_1, \dots, \hat{\rho}_Q\}$  be a list of distinct Finnish essential roots (without vowel blotting) that contain at least one instance of **V** (in either upper or lower case) and DO NOT match the following string patterns:

$$(\mathbf{CVij|VC|ko|mi|te|tu|(ant|ei|ent|nut|nyt|to|vat|väl|vät)}\mathbf{X}).$$

The output of the function  $\text{CpdDet}(\Lambda^{\hat{\rho}})$  is obtained through the following procedures:

- (1) Construct a list  $\{(\hat{\rho}_1, \lambda_1^{\hat{\rho}}), \dots, (\hat{\rho}_Q, \lambda_Q^{\hat{\rho}})\}$  where  $\lambda_q^{\hat{\rho}} = \{\hat{\rho}_{(q,1)}, \dots, \hat{\rho}_{(q,n_q)}\}$  is a subset of  $\Lambda^{\hat{\rho}}$  whose members all match the string pattern  $\hat{\rho}_{q\sim}$ , for  $q \in \mathbb{Z} \cap [1, Q]$ .
- (2) Expand the aforementioned entry  $(\hat{\rho}_q, \lambda_q^{\hat{\rho}})$  into a list of triplets  $\{(\hat{\rho}_{(q,1)}, \hat{\rho}_q, \hat{\rho}_{(q,1)} \ominus \hat{\rho}_q), \dots, (\hat{\rho}_{(q,n_q)}, \hat{\rho}_q, \hat{\rho}_{(q,n_q)} \ominus \hat{\rho}_q)\}$  for every  $q \in \mathbb{Z} \cap [1, Q]$  such that  $\lambda_q^{\hat{\rho}} \neq \emptyset$ . Collect all these triplets as one runs through the list  $\{(\hat{\rho}_1, \lambda_1^{\hat{\rho}}), \dots, (\hat{\rho}_Q, \lambda_Q^{\hat{\rho}})\}$ . The list of these triplets  $\{(\hat{\rho}_{(1)}, \hat{\eta}_{(1)}, \hat{\rho}_{(1)} \ominus \hat{\eta}_{(1)}), \dots, (\hat{\rho}_{(Q')}, \hat{\eta}_{(Q')}, \hat{\rho}_{(Q')} \ominus \hat{\eta}_{(Q')})\}$  contains potentially valid decompositions of compounds.
- (3) Screen the aforementioned list of triplets as follows: for every  $q' \in \mathbb{Z} \cap [1, Q']$ , if  $(\hat{\rho}_{(q')}, \hat{\eta}_{(q')}, \hat{\tau}_{(q')}) = \hat{\rho}_{(q')} \ominus \hat{\eta}_{(q')}$  satisfies

$$\ell(\hat{\rho}_{(q')} \ominus \hat{\eta}_{(q')}) \geq 2 \quad \text{AND} \quad \hat{\rho}_{(q')} \ominus \hat{\eta}_{(q')} = \mathbf{X}_1 \mathbf{V} \mathbf{X}_2,$$

then construct  $\hat{\tau}_{(q')}^*$  by performing  $(e|i)ksi \sim \rightarrow \emptyset$ ,  $(en|En|et|i|ta|ten|tä) \sim \rightarrow \emptyset$ ,  $(i|I|ie)n \sim \rightarrow \emptyset$ ,  $e(L|lt|n|S|T)(a|e|ä) \sim \rightarrow \emptyset$ ,  $i(n|S|T)(a|e|ä) \sim \rightarrow \emptyset$  on  $\hat{\tau}_{(q')}$ , before generating a list  $\lambda_{(q')}^{\hat{\tau}}$  by members of  $\Lambda^{\hat{\rho}}$  that match the pattern  $(\hat{\tau}_{(q')}|\hat{\tau}_{(q')}^*)$ ; otherwise, set  $\lambda_{(q')}^{\hat{\tau}} = \emptyset$ .

- (4) Collect all the triplets  $(\hat{\rho}_{(q')}, \hat{\eta}_{(q')}, \lambda_{(q')}^{\hat{\tau}})$  where  $\lambda_{(q')}^{\hat{\tau}}$  is non-void and  $\hat{\tau}_{(q')}$  DOES NOT match the following string patterns:

$$(Tom|ton|Töm|tön|τντ)\mathbf{X}.$$

Trim the tags in the triplets by doing  $\sim^{\mathbf{X}\epsilon}(\mathbf{VC})\mathbf{V}_m \rightarrow \mathbf{X}$  on the strings. This resulting list of triplets  $\text{CpdDet}(\Lambda^{\hat{\rho}})$  contains the heuristic decompositions of all the identified binary compounds.

**Algorithm 8.14** (Approximate clustering of Finnish words with heuristic detection of compounds). The approximate clustering of a list of Finnish words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  respecting compounding is completed in four stages:

- (1) Do as in Algorithm 8.12(1).
- (2) Do as in Algorithm 8.12(2). Save both the tagged clusters  $\{\Gamma_1 = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, \Gamma_K = \{G_{(K,1)}, \dots, G_{(1,m_K)}\}\}$  and the list of word clusters  $\{\check{\Gamma}_1, \dots, \check{\Gamma}_K\}$  for further use.
- (3) Construct a tagged list of word clusters  $\{(\check{\Gamma}_1, \Lambda_1^{\hat{\rho}}), \dots, (\check{\Gamma}_K, \Lambda_K^{\hat{\rho}})\}$  where  $\Lambda_k^{\hat{\rho}}$  is the union of all Finnish essential roots (without vowel blotting) available to  $\Gamma_k$ , for  $k \in \mathbb{Z} \cap [1, K]$ . Set  $\Lambda^{\hat{\rho}} = \Lambda_1^{\hat{\rho}} \cup \dots \cup \Lambda_K^{\hat{\rho}}$ , and evaluate  $\text{CpdDet}(\Lambda^{\hat{\rho}})$ .
- (4) The first component  $\hat{\rho}_{(q'')}$  of each triplet  $(\hat{\rho}_{(q'')}, \hat{\eta}_{(q'')}, \lambda_{(q'')}^{\hat{\tau}})$  in  $\text{CpdDet}(\Lambda^{\hat{\rho}})$  is called a “dissolvable compound”, the second component  $\hat{\eta}_{(q'')}$  a “heuristic head”, and the first member in the third component  $\lambda_{(q'')}^{\hat{\tau}}$  a “heuristic tail”. In the tagged list  $\{(\check{\Gamma}_1, \Lambda_1^{\hat{\rho}}), \dots, (\check{\Gamma}_K, \Lambda_K^{\hat{\rho}})\}$ , every entry containing a “dissolvable compound” is removed, and regrouped with the entries matching its “heuristic head” and “heuristic tail”. Finally, remove all tags.

*Example 8.14.1.* If we have the inflected compound words from the aforementioned KOTUS-51 sample (containing 27 words), as well as the words *nuori* “young” and *pari* “couple”, then our heuristic compound detection generates the following output:

{*nuoreenpariin*, *nuoreksipariksi*, *nuorellaparilla*, *nuorelleparille*, *nuoreltaparilta*, *nuorenparina*, *nuorenparin*, *nuoressaparissa*, *nuorestaparista*, *nuoretparit*, *nuoretaparitta*, *nuorienparien*, *nuoriinpareihin*, *nuoriksipareiksi*, *nuorinapareina*, *nuorinepareineen*, *nuorinparein*, *nuoripari*, *nuorissapareissa*, *nuoristapareista*, *nuorittapareitta*, *nuortaparia*, *nuortenparien*, *pari*} — 24 words;

{*nuoreenpariin*, *nuoreksipariksi*, *nuorellaparilla*, *nuorelleparille*, *nuoreltaparilta*, *nuorenparina*, *nuorenparin*, *nuoressaparissa*, *nuorestaparista*, *nuoretparit*, *nuoretaparitta*, *nuori*, *nuoriapareja*, *nuorienparien*, *nuoriinpareihin*, *nuoriksipareiksi*, *nuorillapareilla*, *nuorillepareille*, *nuoriltapareilta*, *nuorinapareina*, *nuorinepareineen*, *nuorinparein*, *nuoripari*, *nuorissapareissa*, *nuoristapareista*, *nuorittapareitta*, *nuortaparia*, *nuortenparien*} — 28 words.

This is partially satisfactory: a small fraction of the compounds were not decomposed properly (component *pari* undetected), but were simply regarded as derivatives of the word *nuori*. If we have *nuorin* “youngest” instead of *nuori* “young” in our input, then we get something slightly worse:

{*nuoriapareja*} — 1 word;

{*nuoreenpariin*, *nuoreksipariksi*, *nuorellaparilla*, *nuorelleparille*, *nuoreltaparilta*, *nuorenparina*, *nuorenparin*, *nuoressaparissa*, *nuorestaparista*, *nuoretparit*, *nuorettaparitta*, *nuorienparien*, *nuoriinpareihin*, *nuoriksipareiksi*, *nuorinepareineen*, *nuorinparein*, *nuoripari*, *nuorissapareissa*, *nuoristapareista*, *nuorittapareitta*, *nuortaparia*, *nuortenparien*, *pari*} — 23 words;

{*nuoreenpariin*, *nuoreksipariksi*, *nuorellaparilla*, *nuorelleparille*, *nuoreltaparilta*, *nuorenparina*, *nuorenparin*, *nuoressaparissa*, *nuorestaparista*, *nuoretparit*, *nuorettaparitta*, *nuorienparien*, *nuoriinpareihin*, *nuoriksipareiksi*, *nuorillapareilla*, *nuorillepareille*, *nuoriltapareilta*, *nuorin*, *nuorinapareina*, *nuorinepareineen*, *nuorinparein*, *nuoripari*, *nuorissapareissa*, *nuoristapareista*, *nuorittapareitta*, *nuortaparia*, *nuortenparien*} — 27 words.

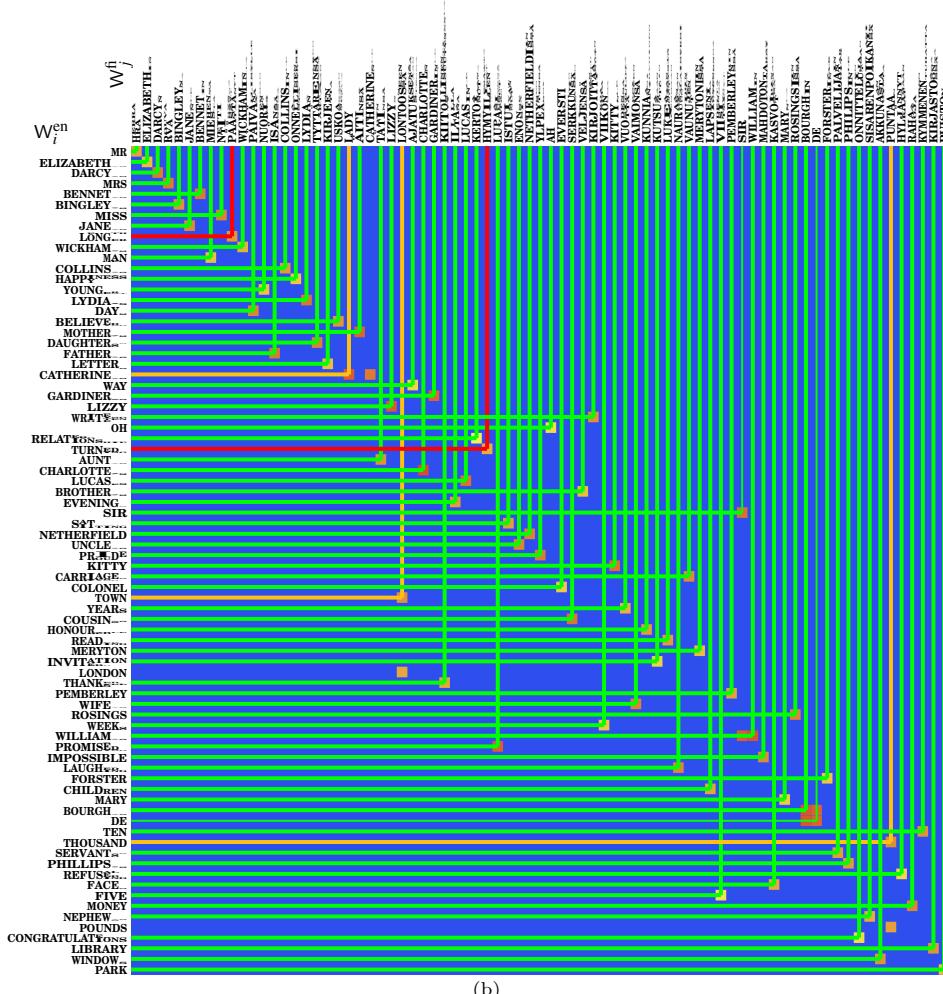
*Example 8.14.2.* In Fig. S12, we further test the aforementioned word clustering algorithm (with heuristic detection of compounds) on topic extraction and machine translation (see Table S1 for text sources).

*Peni(n)kulma* is an archaic measure of length used in Finland. Although this is not numerically equivalent to the English *mile*, we consider these length units exact matches in Fig. S12b' and b''.

In Darwin's *Origin of Species*, the word "fresh" is generally used in the context of *fresh water* and the word "reciprocal" is mostly used in the phrase *reciprocal cross*. Therefore, we consider *suolattoman* (genitive singular of *suolaton* "saltless") an exact match to *fresh* and *vastavuoroinen* "reciprocal" a close match to *graft* in Fig. S12b''. Similarly, we regard *siitoselimiristö* "breeding organ" as a near synonym to *reciprocal*.

HRA~ ELIZABETH~ TULE~ DARCY~ SANOT~ HXYÄ~ TUNTE~ PUHE~  
 KUULLA~ RVA~ BINGLEY~ JANE~ KÄY~ BENNET~ SISARENSX~ NÄHÖ~  
 RAKAS~ MIEHEN~ NÄT~ MIELE~ TIE~ PÄÄSTX~ ASIAS~ PITÄ~ TOS~ AIX~  
 WIICKHAM~ LÄHT~ TARA~ TOIVO~ PÄIVÄ~ VARMASTA~ YSTÄVÄSX~ SUURE~ NUORE~  
 ISANSA~ COLLINS~ ANTA~ OSA~ PUOLEL~ ONNEELINER~ TAHTO~ LYDIA~ HUOMA~  
 TYTTÄRENSX~ SEURA~ VIERA~ KATSEL~ KIRJEEN~ SILMÄ~ KERPA~ LUULI~ VAHÄN~ NAISE~ ARVAAN~  
 USKO~ MIELX~ LADY~ NAYTTA~ VIIME~ PERHEENS~ PER~ SUINKÄÄ~ TYT89~ ÄITINSX~ CATHERINE~ ILO~  
 NAIMISTON~ IHMISSX~ ILM~ PAHA~ ARVO~ PÄÄTKE~ PALAA~ SYDÄMES~ KÄSITÄN~ ODOTTA~ IHAS~  
 TATI~ OT~ KOHTELIA~ LIZZY~ MATKÄ~ HALU~ KAUNIS~ LONTOOSX~ AJATUST~ JATKOT~ SUOSIUS~ PARU~ HUUDAHDI~  
 PIDÄT~ UU~ TIETV~ SUHTEES~ CHARLOTTE~ AJATTEL~ NIME~ LONGBOURNIAS~ MEN~ TALOS~ TANSIA~  
 KUITTOL~ SURF~ LÄHETT~ ALBES~ LAUSU~ SAAPU~ MÄHDOLLIST~ EPÄILY~ TILAISUUVX~ LUONTEEN~ LUCAS~ KOTI~  
 KERTOS~ PITR~ MAINTT~ SALLI~ AIHETTA~ KÄRBY~ HYMYLLX~ VAKUTUS~ LUKU~ MERKILLY~ KÄYTÖ~ LUUPA~ ISTUK~ ENON~  
 YLPEX~ JOUTU~ ELÄX~ SATTU~ AH MUIST~ EVERSTI~ SERKKUN~ KITTY~ YUOW~ TÄYDELL~ VAIMONSA~  
 KYSTY~ PELKÄSE~ LOPUT~ VIHKOS~ KIELTAX~ TUTTAVUUV~ HETKE~ ISÄNTÄ~ PYSY~ LUKU~ KUTSU~ SALA~  
 VASTA~ TARVITSE~ MAINTT~ KUNNIA~ KAUAN~ SAATTAS~ SELITTAS~ JÄÄS~ SELVÄV~  
 NAUR~ ESKY~ TAPARIT~ TIEDUS~ KIPPS~ VAIKEA~ VIMM~ OMIAIS~ VAKUUS~ SÄVY~ KOLME~  
 KULKE~ TÖÖJ~ ROHES~ VÄSTAREN~ KIPPS~ VAIKEA~ VIMM~ OMIAIS~ VAKUUS~ PÄÄTKE~  
 VAL~ TÄR~ AMUN~ DRUK~ VÄRÄVÄ~ VISA~ KÄVY~ VÄRÄVÄ~ SÄVY~ PÄÄTKE~  
 KASVY~ KINTYVÄ~ ERIC~ AIBO~ HAMMASTI~ HAMMAS~ VÄRÄVÄ~ VÄRÄVÄ~ VÄRÄVÄ~  
 KUNGA~ BOURGH~ MELPYVERE~ EHPOYKE~ KOVAK~ ASTUR~ AKHERK~ YMMÄRXT~ DE ASHEM~  
 ERO~ HARTA~ VAKUTU~ SIISZÄTT~ PIENE~ HELLSY~ STVAK~ HASKA~ ANN~ KELLOS~ HUOL~  
 KIRJ~ MOITY~ KESY~ PALVELLAJAC~ OPPRE~ EHDY~ AIDH~ MAINNO~ HURST~ FITZWILLIAM~  
 UUTISYA~ JUOKS~ ARVOSA~ TERVETILMÄ~ HÄPEM~ TUTTAS~ LÄHETT~ KÄRBY~  
 KOR~ YLÄMÄ~ MÄNAS~ LIHTA~ ROSINTASA~ ANSATSV~ LIVOG~ HÄYÄV~ ESTITA~  
 ONNITTEL~ ONNITTELU~ KATKEBE~ HUNSFORDE~ YLSTV~ JALOSUKUUTSENKA~  
 BRIGHTONIAN~ VALSIMÄTTOMV~ PHIRYTA~ OVS~ KALLA~ HEPUS~ VIERAS~ ENEMPA~  
 LAUL~  
 W<sup>f</sup><sub>j</sub>  
 W<sup>en</sup><sub>i</sub>

(a)



(b)

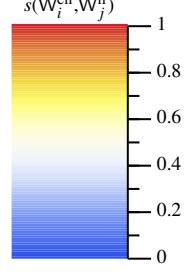


Fig. S12. Text mining in Finnish.  
 (a) Statistically identified topics ( $n_{ii} \geq 20$ ) in a Finnish version of *Pride and Prejudice*, with the same color encoding scheme as Fig. S3. (b) Semantic similarities  $s(W_i^e, W_j^f)$  between selected topics in English and Finnish versions of *Pride and Prejudice*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) ranking of the similarity score for the optimal node. Green (resp. red) cross-hair indicates a correct (resp. incorrect) match. Amber cross-hair marks a link between distinct concepts that share the same hypernyms.

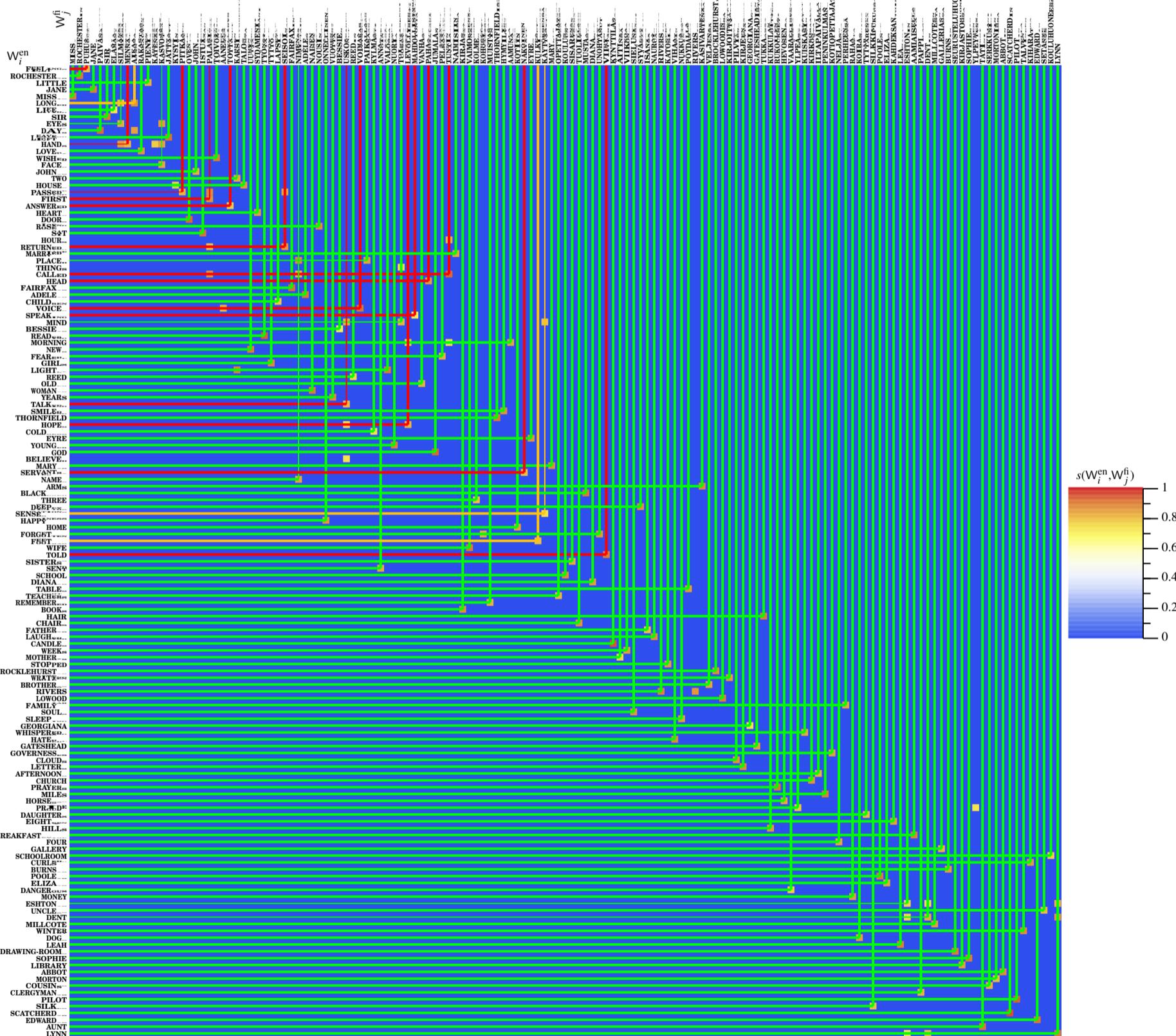
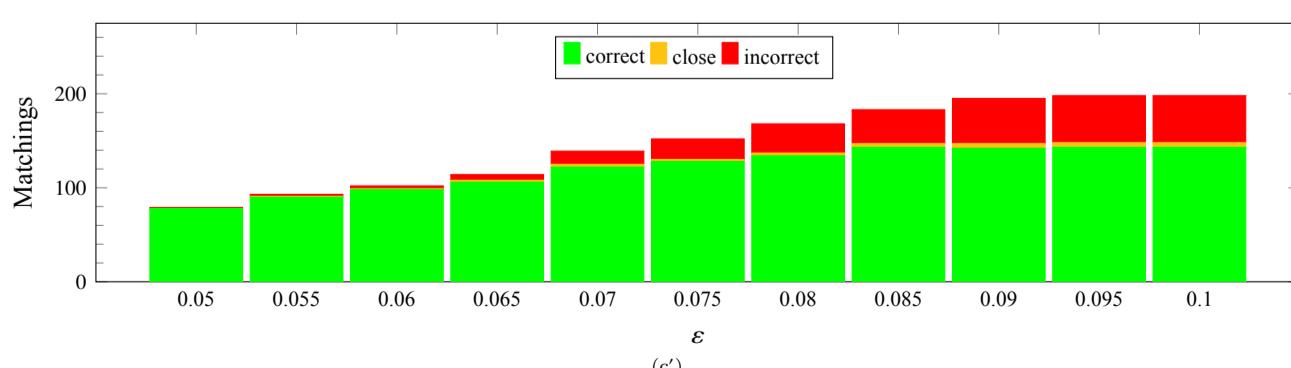


Fig. S12. Text mining in Finnish. (Continued)  
 (a') Statistically identified topics ( $n_{ii} \geq 20$ ) in a Finnish version of *Jane Eyre*, with the same color encoding scheme as Fig. S3. (b') Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{fi}})$  between selected topics in English and Finnish versions of *Jane Eyre*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) rank of the similarity score for the optimal node. Green (resp. red) cross-hair indicates a correct (resp. incorrect) match. Amber cross-hair marks a link between distinct concepts that share the same hypernyms. (c') Results from control experiments with different choices of the  $\varepsilon$ -parameter in ballpark screening criteria (1.13).



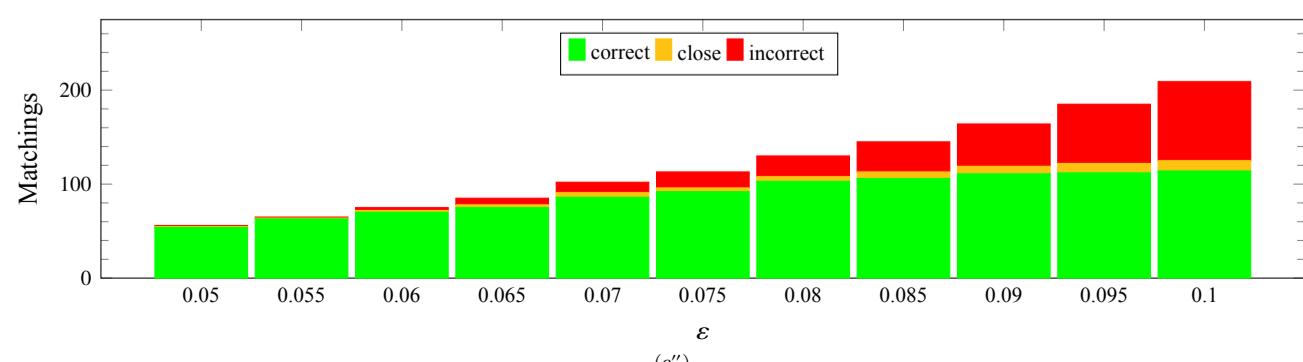
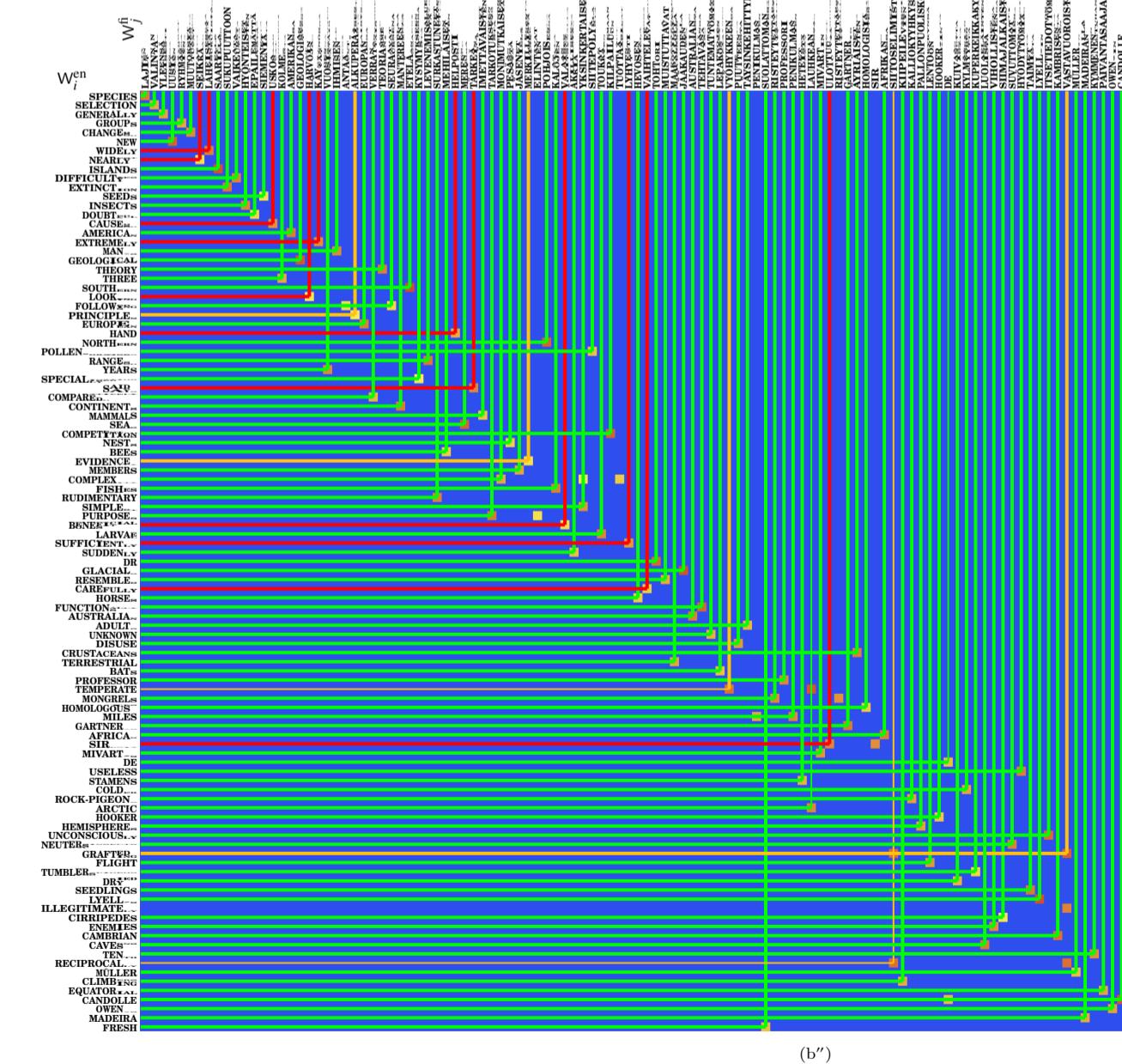


Fig. S12. Text mining in Finnish. (Continued)  
 (a'') Statistically identified topics ( $n_{ii} \geq 20$ ) in a Finnish version of *Origin of Species*, with the same color encoding scheme as Fig. S3. (b'') Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{fi}})$  between selected topics in English and Finnish versions of *Origin of Species*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) rank of the similarity score for the optimal node. Green (resp. red) cross-hair indicates a correct (resp. incorrect) match. Amber cross-hair marks a link between distinct concepts that share the same hypernyms. (c'') Results from control experiments with different choices of the  $\epsilon$ -parameter in ballpark screening criteria (1–13).

## 8.2 Modified Porter stemming algorithm for Hungarian

In our presentation below, we will disable hyphenation for all the Hungarian words. There are two reasons for doing so: (1) The hyphenation rules for Hungarian are both complex and counterintuitive, which may baffle non-Hungarian readers of this document. (2) The Hungarian language option in the L<sup>A</sup>T<sub>E</sub>X *babel* package is incompatible with many other L<sup>A</sup>T<sub>E</sub>X packages, which discourages us from loading it.

**Definition 8.15** (Hungarian stop words). If a word, up to removal of *~-e* (hyphen and the letter *e*, attached at the end<sup>100</sup>), belongs to the following list<sup>101</sup>:

*a, abban, ahoz, ahogy, ahogyan, ahol, ahonnan, ahonnét, ahoa, akár, akárhol, akárhonnan, akárhonnét, akárhova, aki, akik, akkor, aközben, alá, alád, alájuk, alám, alánk, alátok, alatt, alatta, alattad, alattam, alattatok, alattuk, alattunk, alig, aligha, alighanem, alighogy, alól, alóla, alólad, alólám, alólatok, alóluk, alólunk, által, általa, általában, általad, általam, általatok, általuk, általunk, alul, amely, amelyek, amelyekben, amelyeket, amelyet, amelynek, ami, amíg, amikor, amint, amít, amolyan, annak, annyi, annyiba, annyiban, annyiból, annyiért, annyihoz, annyiig, annyként, annyin, annyinak, annyinál, annyira, annyiról, annyiszor, annyit, annyitól, annyivá, annyival, arra, arról, át, az, azaz, azért, azok, azon, azonban, azt, aztán, azután, azzal, azzonal, bár, bárcsak, bárhogyan, bárhol, bárhonnan, bárhonnét, bárhova, bármikor, be, belé, beléd, beléjük, belém, belénk, belétek, beljebb, belőle, belőled, belőlem, belőletek, belőlük, belőlünk, belül, benn, benne, benned, bennem, bennetek, benneteket, bennük, bennünk, bennünket, bent, bizony, cikk, cikkek, cikkeket, csak, csaknem, csinál, csinálandó, csináld, csinálhat, csinálj, csinálja, csináljad, csináljak, csinálják, csináljál, csináljalak, csináljam, csináljanak, csináljatok, csináljátok, csináljon, csináljuk, csináljunk, csinállak, csinálna, csinálnád, csinálnak, csinálnák, csinálnál, csinálnálak, csinálnám, csinálnának, csinálnátok, csinálnék, csinálni, csinálnia, csinálniuk, csinálnod, csinálnom, csinálnotok, csinálnunk, csináló, csinálok, csinálom, csinálsz, csinált, csinálta, csináltad, csináltak, csináltál, csináltalak, csináltam, csináltatok, csináltatók, csináltuk, csináltunk, csinálunk, csinálva, de, dehát, dehogy, e, ebben, eddig, egész, egy, egybe, egyben, egyből, egye, egyéb, egyébként, egyed, egyem, egyen, egyért, egyes, egyet, egyetek, egyetlen, eggé, egygel, egyhez, egyig, egyik, egyikbe, egyikben, egyikből, egyike, egyiked, egyikek, egyikekbe, egyikekben, egyikekből, egyikeken, egyikekért, egyikeket, egyikekhez, egyikekig, egyikekké, egyikekkel, egyikekként, egyikeknak, egyikeknél, egyikekre, egyikekről, egyikektől, egyikem, egyiken, egyikért, egyiket, egyikhez, egyikig, egyikké, egyikkel, egyikként, egyiknek, egyiknél, egyikre, egyikről, egyiktől, egyikük, egyikünk, egyiként, egymás, egymásba, egymásban, egymásból, egymásért, egymáshoz, egymásig, egymásként, egymásnak, egymásón, egymásra, egymásról, egymássá, egymással, egymást, egymástól, egynek, egnél, egyre, egyről, egytől, együk, egyné, együtt, ehhez, ekkor, eközben, el, elé, eléd, elég, eléjük, elém, elénk, elétek, ellen, ellenben, ellené, ellenem, ellenére, ellenetek, ellenük, ellenünk, elő, elől, előle, előled, előlem, előletek, előlük, előlünk, először, előtt, előtte, előtted, előtttem, előttetek, előttük, előttünk, első, emilyen, én, énbelőlem, engem, ennek, éppen, erre, érte, érted, értetek, értük, értünk, és, esetében, esetedben, esetben, esetén, esetekben, esetleg, esetükben, esetünkben, ez, ezek, ezen, ezért, ezt, ezzel, fel, fele, felé, feléd, feléjük, felém, felénk, felétek, felett, felől, felőle, felőled, felőlem, felőletek, felőliuk, felőlünk, felül, fenn, fent, fenti, fentibe, fentiben, fentiből, fentiek, fentiekbe, fentiekben, fentiekből, fentieken, fentiekért, fentieket, fentiekhez, fentiekig, fentiekké, fentiekkel, fentiekként, fentieknek, fentieknel, fentiekre, fentiekről, fentiektől, fentiért, fentihez, fentiig, fentiként, fentin, fentinek, fentinél, fentire, fentiről, fentit, fentitől, fentivé, fentivel, fog, fogja, fogják, fogjátok, fogjuk, fognak, fogod, fogok, fogom, fogsz, fogtok, fogunk, fogva, folyamán, folytán, föl, fölé, föld, föléjük, fölém, fölenk, fölétek, följebb, fölött, fölöttem, fölöttem, fölötterek, fölötük, fölötünk, fölül, fölüle, fölüled, fölulem, fölületek, fölülük, fölülünk, gyanánt, ha, hadd, hanem, hánny, hasonlóan, hát, hátha, helyett, helyette, helyettem, helyettetek, helyettiuk, helyettünk, hiszen, hogy, hogyan, hogyha, hogyne, hol, holott, honnan, honnét, hosszat, hova, hozzá, hozzájuk, hozzáám, hozzáánk, hozzátok, ide, időközben, igen, így, ill, ill., illetve, ily, ilyen, ilyenkor, innen, iránt, iránta, irántad, irántam, irántatok, irántuk, irántunk, is, ismét, ison, itt, jóvoltából, kedvér, kegyed, kegyedben, kegyedből, kegyeddé, kegyeddel, kegyeden, kegyedért, kegyedet, kegyedhez, kegyedig, kegyedként, kegyednek, kegyednél, kegyedre, kegyedről, kegyedtől, kell, kellett, képest, keressünk, keresztül, kezdve, ki, kijebb, kinn, kint, kivéve, kívül, köré, köréd, köréjük, körém, körénk, körétek, körül, körülötte, körülöttem, körülöttek, körülöttük, körülöttünk, következetében, közben, közbeni, közé, közéjük, közédm, közénk, közétek, között, közöttem, közöttem, közöttetek, közöttük, közöttünk, közé, közé, közüle, közüled, közülem, közületek, közülik, közülink, különben, le, leendő, legalább, légy, legyek, legyél, legyen, legyenek, legyetek, legyünk, lehet, lehetett, lejebb, lenn, lenne, lenned, lennem, lennének, lennetek,*

<sup>100</sup>Word final *~-e* occurs in Hungarian tag questions.

<sup>101</sup>Our list of Hungarian stop words is based on <http://snowball.tartarus.org/algorithms/hungarian/stop.txt> (prepared by Anna Tordai), with extensive additions to roughly match their counterparts in English. Although the words *jó* “good”, *jobban* “better”, *jól* “well”, *nagy* “big”, *nagyobb* “bigger”, *új* “new”, *újabb* “newer” showed up in Anna Tordai’s list of Hungarian stop words, they do not appear in our list here, for the sake of consistency with other languages.

or matches one of the following string patterns:

( $\emptyset$ |jö|ön)mag(a|á|u( $\emptyset$ |n)k)(d|m|tok)( $\emptyset$ |at|ba|ban|ból|dá|dal|ért|hoz|ig|ká|kal|ként|má|mal|n|nak|nál|on|ra|ról|t|tó|vá|vá).

(as)f[es]mi|ki|sem[m]i)f[é]l(e|é)(Ø|k)(Ø|be|ben|b[ó]l|en|ént|ért|et|hez|ig|k[é]|kel|k[é]nt|n|ne[k]|né|re|r[ó]l|t|t[ó]l|v[é]|vel).

( $\emptyset|a|akár|bár|más| minden| mint| ne|né| se| vala)ho(gy|l|va|vá).$

(Ø|mind|ugyan)e(bbe|bben|bból|ddig|hhez|kként|nnék|nnél|rre|rról|ttól|z|zek|zekbe|zekben|zekból|zekent|zekért|zeket|zekhez|zekig|zekké|zekkel|zekeként|zeknek|zeknél|zekre|zekról|zektől|zen|zér|zt|zzé|zzel),

(Ø|mind|ugyan)a(bba|bba|bból|ddig|hhoz|kként|nnak|nnál|rra|rról|ttól|z|zér|zok|zokat|zokba|zokban|zokból|zokért|zokhoz|zokig|zokká|zokkal|zokként|zoknak|zoknál|zokon|zokra|zokról|zoktól|zon|zt|zzá|zzal),

(Ø|akárm|am|bárm|m|né|sem|ugyan|valam)(annyi|elyik|ennyi|iłyen)(Ø|e)(Ø|d|je|jei|jeid|jeik|jeim|jeink|jeitek|jük|m|nk|tek)(Ø|e)(Ø|be|ben|ból|dé|del|ért|hez|ig|k|kbe|kben|kból|ké|kel|ken|ként|kért|ker|khez|kig|kké|kkel|kként|knék|knél|kre|król|ktől|mé|mel|n|nek|nél|re|ről|t|től|vél|vel),

(Ø|a|akár|bár|minden|sem|sen|vala)(k|m)i(Ø|d|é|je|jei|jeid|jeik|jeim|jeink|jeitek|jük|m|nk|tek)(Ø|e)(Ø|be|ben|ból|dé|del|ért|hez|ig|k|kbe|kben|kból|ké|kel|ken|ként|kért|ker|khez|kig|kké|kkel|kként|knék|knél|kre|król|ktől|mé|mel|n|nek|nél|re|ről|t|től|vél|vel),

(enyéim|enyém|magáé|maguké|magukéi|miek|mienk|miénk|öné|önéi|önöké|önökéi|övé|övéi|övéik|övék|tied|tiéd|tieid|tieitek|tieitek|tiétek)(Ø|be|ben|ból|dé|del|ek|ekbe|ekben|ekból|eken|ekért|eket|ekhez|ekig|lekké|ekkel|ekként|eknek|eknél|ekre|ekról|ektől|en|ért|et|hez|ig|ké|kel|ként|mé|mel|nek|nél|re|ről|től|vél|vel),

(ahány|akárhány|néhány|sok|valahány)(Ø|ak|akat|akba|akban|akból|akért|akhoz|akig|akká|akkal|akként|aknak|aknál|akon|akra|akról|aktól|at|ba|ban|ból|ért|hoz|ig|ká|kal|ként|nak|nál|on|ra|ról|től),

(kevés|kevesebb|legkevesebb|legtöbb|több)(Ø|ed|em|etek|je|jei|jeid|jeik|jeim|jeink|jeitek|jük|iink)(Ø|d|k|m|v)(Ø|e|é|ek|ekbe|ekben|ekból|eken|ekért|eket|ekhez|ekig|ekké|ekkel|ekként|eknek|eknél|ekre|ekról|ektől|el|en|en|ért|et|hez|ig|ként|nek|nél|ől|re|ről|től),

(Ø|akár|bár|se)(ak|ek|jó|mek)kor(a|ába|ában|ából|áért|ához|áig|ák|ákat|ákba|ákban|ákból|aként|ákért|ákhoz|ákig|ákká|ákkal|ákként|áknak|áknál|ákon|ákra|ákról|áktól|án|ának|ánál|ára|áról|át|ától|ává|ával),

(Ø|a|h(á|a)nyadik(Ø|ak|akat|akba|akban|akból|akért|akhoz|akig|akká|akkal|akként|aknak|aknál|akon|akra|akról|aktól|at|ba|ban|ból|ért|hoz|ig|ká|kal|ként|nak|nál|ok|okat|okba|okban|okból|okért|okhoz|okig|okká|okkal|okként|oknak|oknál|okon|okra|okról|októl|on|ot|ra|ról|től),

(Ø|épp|ugyan)olyan(Ø|ja|jai|jaid|jaik|jaim|jaink|jaitok|juk|od|om|otok|unk)(Ø|ba|ban|ból|dá|dal|ért|hoz|ig|ká|kal|ként|má|mal|ná|nak|nal|nál|ok|okat|okba|okban|okból|okért|okhoz|okig|okká|okkal|okként|oknak|oknál|okon|okra|okról|októl|on|ra|ról|től),

(Ø|a|bár|né|vala)mely(Ø|be|ben|ból|ek|ekbe|ekben|ekból|eken|ekért|eket|ekhez|ekig|ekké|ekkel|ekként|eknek|eknél|ekre|ekról|ektől|en|ért|et|hez|ig|ként|nek|nél|re|ről|től),

then we consider it a Hungarian stop word. All the Hungarian stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list.

The Hungarian verbs for “buy”, “carry”, “drink” and “eat” have highly irregular conjugations (see Example 8.24.2 below for full lists of their conjugated forms), which also resemble other unrelated words (to untrained eyes and naïve algorithms). We are not going to include these verbs in our list of stop words, but will use their full conjugation tables to define word patterns **buyHungarian**, **carryHungarian**, **drinkHungarian** and **eatHungarian** in order to handle them separately in the word clustering algorithm below. □

Up to some modifications, our method is modeled after Anna Tordai’s stemming algorithm for Hungarian (<http://snowball.tartarus.org/algorithms/hungarian/stemmer.html>). Our modifications of Tordai’s stemming algorithm aim at two goals: to fully take care of verb conjugations and to fully accommodate to vowel alternations and fleeting vowels in declensions.

### 8.2.1 Effective spelling and essential root

It is assumed that all Hungarian words are converted to lowercase before going through any of the procedures below.

We note that Hungarian vowel lengths are sometimes used to distinguish words of unrelated etymologies, such as *agy* “brain” vs. *ágy* “bed”, *örül* “rejoices” vs. *őrül* “goes crazy”, and *vall* “confesses” vs. *váll* “shoulder”. Meanwhile, many other words may alter between short *a* and long *á* (or short *ö* and long *ő*) during inflections. We will allow vowel length variation as a general rule, and accommodate to special cases by adding exceptions.

The pronunciations of Hungarian *s* (like English *sh*) and *sz* (like English *s*) are exactly the opposite to the Polish practices. However, for convenience, we will still replace the Hungarian diagraph *sz* by *ś* for clustering purposes. Likewise, in what follows, the replacements of *gy* and *ty* by the Icelandic letters *ð* and *þ* are not motivated by phonological resemblances.

**Definition 8.16** (Hungarian Vowels). Hereafter in §8.2, the symbol  $\mathbf{V}$  stands for any member from the list of Hungarian vowels  $\{a, \acute{a}, e, \acute{e}, i, \acute{i}, o, \acute{o}, \ddot{o}, \ddot{\acute{o}}, u, \acute{u}, \ddot{u}, \ddot{\acute{u}}\}$ . In line with the multiplicity notations introduced in Definition 3.3, the symbol  $\mathbf{V}_m$  stands for a text string formed by consecutive appearance of one or more (not necessarily identical or distinct) members from the set of Hungarian vowels.

Dual to the notations above, the symbol  $\mathbf{C}$  stands for any character that does not belong to the list  $\{a, \acute{a}, e, \acute{e}, i, \acute{i}, o, \acute{o}, \ddot{o}, \ddot{\acute{o}}, u, \acute{u}, \ddot{u}, \ddot{\acute{u}}\}$ , and  $\mathbf{C}_{m_0}$  stands for a text string formed by consecutive appearance of zero or more (not necessarily identical or distinct) characters that do not belong to the same list.  $\square$

The following protected range of Hungarian words is defined according to Anna Tordai's Hungarian stemming algorithm (<http://snowball.tartarus.org/algorithms/hungarian/stemmer.html>). Some subsequent steps in our clustering algorithms are also inspired by hers. We note, however, that Anna Tordai's algorithm only treats the declensions of Hungarian nouns and adjectives (without vowel alternations and fleeting vowels), and does not fully consider the verb conjugations. Some cosmetic augmentations are needed to make the clustering algorithm work for all the test words.

**Definition 8.17** (Hungarian protected range). Let  $\hat{\sigma}$  be a Hungarian text string, then its protected range  $\text{ProtRg}(\hat{\sigma}) = \max\{\lambda_1(\hat{\sigma}), \lambda_2(\hat{\sigma})\}$  is specified through the following procedures:

- Look for the string pattern  $(\emptyset|be|sel|ki|le|meg)\mathbf{V}_m\mathbf{C}\sim$  in the string  $\hat{\sigma}$ ;
- If such a string pattern is found, its last occupying position defines  $\lambda_1(\hat{\sigma})$ ; otherwise, set  $\lambda_1(\hat{\sigma}) = 0$ ;
- Look for the string pattern  $(\emptyset|be|sel|ki|le|meg)\mathbf{C}_m\mathbf{V}\sim$  in the string  $\hat{\sigma}$ ;
- If such a string pattern is found, its last occupying position defines  $\lambda_2(\hat{\sigma})$ ; otherwise, set  $\lambda_2(\hat{\sigma}) = 0$ .  $\square$

**Algorithm 8.18** (Hungarian effective spelling). For a Hungarian word  $\hat{\sigma}$  (converted to lower case form), its effective spelling  $\text{EffSpell}(\hat{\sigma})$  is constructed in sequential steps:

(1) Convert to lowercase, do  $-X \rightarrow \emptyset$  (i.e. delete everything after hyphen),  $dr \sim \rightarrow \delta r$ ,  $tr \sim \rightarrow \tau r$ ,  $vissziik \rightarrow visszik$ ,  $\sim nyig \rightarrow ny$ ,  $\sim r(a|e|o)st(u|\ddot{u})l \rightarrow r$ , before replacing<sup>102</sup>

$(\emptyset be)vall\sim$	$(\emptyset édes)any(\emptyset j)$	$(\emptyset édes)ap(a \acute{a})$	$(\emptyset édes)apj$
$a\delta\mu\tau$	$any$	$x\varphi a\delta\rho$	$x\varphi a\delta\rho$
$(\emptyset leg legesleg)f\acute{in}om\sim$	$(\emptyset leg legesleg)gazdag\sim$	$(\emptyset leg legesleg)hideg\sim$	
$\varphi i\tilde{v}m$	$\rho i\tilde{e}g$	$k\acute{o}l\delta g$	
$(\emptyset leg legesleg)magá(\emptyset n)ny\sim$	$(\emptyset leg legesleg)meleg\sim$	$(\emptyset leg legesleg)messz\sim$	$(\emptyset leg legesleg)sim\sim$
$\pi\rho i\beta$	$\omega a\rho\mu g$	$\varphi a\rho s z$	$\sigma i\mu$
$(\emptyset leg legesleg)szórák\sim$	$(\emptyset leg legesleg)vidám\sim$	$(\emptyset meg)ért\sim$	$(\emptyset meg)eskü(d sz)\sim$
$e\tilde{v}jok$	$\tau\zeta\rho am$	$\tilde{v}a\sigma\tau\tilde{v}\delta$	$\sigma e\omega\rho\tilde{u}d$
$(\emptyset nagy)nén\sim$	$(báty báttý fivér fivérr öcs öccs)\sim$	$(leg legesleg)kisebb\sim$	$agg\sim$
$\alpha u\tilde{v}\tau$	$\beta\rho\atilde{r}t$	$kis$	$\omega\rho i\tilde{x}$
$\acute{a}gy$	$alak\sim$	$\acute{a}ltalános$	$bácsi\sim$
$\beta e\delta gy$	$alakk$	$anyag\sim$	$báj$
$\alpha a\tilde{l}l$	$\gamma e\tilde{v}\epsilon\rho a\tilde{l}s$	$apu\sim$	$bál\sim$
$\mu a\tilde{t}g$	$\mu a\tilde{t}g$	$ár\sim$	$bámul$
$x\varphi a\delta\rho$	$x\varphi a\delta\rho$	$asztal\sim$	$bán$
$\tau a\beta l$	$\tau a\beta l$	$bácsi\sim$	
$\tilde{u}\tilde{n}kli$	$\tilde{u}\tilde{n}kli$	$báj$	
$\tau x\alpha\rho\mu j$	$\tau x\alpha\rho\mu j$	$bál\sim$	
$\beta á\acute{a}\acute{l}$	$\beta á\acute{a}\acute{l}$	$bámul$	
$\mu a\tilde{t}l$	$\mu a\tilde{t}l$	$bán$	
$bát(or ra)(\emptyset bb kX)$	$bennetX\sim$	$charles$	$csoport$
$\beta a\rho\alpha vr$	$\beta e\tilde{v}\tilde{y}$	$család$	$dönt$
$egész\sim$	$eliz(á a)X\sim$	$beteg$	
$i\tilde{v}\tau sz$	$él$	$charles$	
$e\tilde{l}u\zeta$	$é\acute{r}kez$	$család$	
$\acute{e}lő(s)X$	$esküvő\sim$	$csoport$	
$\acute{e}l$	$est\sim$	$dönt$	
$\rho a\beta z$	$est\sim$	$\acute{e}rkez$	
$\omega e\delta ö$	$ezredes$	$charles$	
$e\acute{e}\beta$	$faj$	$család$	
$\omega \beta e\rho\sigma t$	$fej$	$csoport$	
$\sigma e\pi j$	$fej$	$dönt$	
$\omega \beta e\rho\sigma t$	$fejlett$	$\acute{e}rkez$	
$\sigma e\pi j$	$fejlett$	$charles$	
$\omega \beta e\rho\sigma t$	$férj$	$család$	
$\chi u\sigma\beta j$	$férj$	$csoport$	
$fiatal$	$font$	$charles$	
$\gamma ul$	$form(a \acute{a})$	$család$	
$\pi u\tilde{v}\delta t$	$fosszil$	$csoport$	
$\varphi o\rho\mu a$	$gyermek\sim$	$dönt$	
$\varphi o\sigma\tilde{u}l$	$haj$	$\acute{e}rkez$	
$gyerek$	$haj$	$charles$	
$\eta a\mu j$	$hál(a \acute{a})\sim$	$család$	
$\theta a\gamma\tilde{v}ka$	$hasz(\emptyset o)n$	$csoport$	
$\iota\sigma\tilde{e}\pi n$	$hat\sim$	$dönt$	
$\iota\sigma\tilde{e}\pi n$	$hat\sim$	$\acute{e}rkez$	
$\iota\sigma\tilde{e}\pi n$	$hiáb$	$charles$	
$hiány$	$hív$	$család$	
$\lambda a\tilde{c}kny$	$hivat$	$csoport$	
$\mu o\tilde{v}\theta p$	$hónap\sim$	$dönt$	
$\lambda o\tilde{s}$	$hossz$	$\acute{e}rkez$	
$\varphi i\gamma\rho m$	$idom\sim$	$charles$	
$\beta e\gamma\tilde{v}l$	$indul$	$család$	
$\omega o\epsilon j$	$jaj\sim$	$csoport$	
$\tau o\tilde{n}k$	$játék$	$dönt$	
$\gamma j\tilde{o}ne\sigma$	$jones$	$\acute{e}rkez$	
$\chi a\rho\delta ny$	$kemény$	$charles$	
$\beta eyin$	$kezd$	$család$	
$kisasszon(\emptyset n)y$	$kíván$	$csoport$	
$\mu i\sigma ny$	$könnyv$	$dönt$	
$\omega i\sigma\chi on$	$könnyvtár$	$\acute{e}rkez$	
$\beta kuv$	$köszi$	$charles$	
$\lambda i\beta\rho ar$	$lá(s ss t)$	$család$	
$\theta a\tilde{v}k$	$lament\sim$	$csoport$	
$\sigma ee\tilde{t}$	$leah\sim$	$dönt$	
$\lambda a\mu e\tilde{v}\tau$	$leah\sim$	$\acute{e}rkez$	
$\lambda e\alpha\tilde{v}\tau$	$leah\sim$	$charles$	
$\lambda e\alpha\tilde{v}\tau$	$lyhet$	$család$	
$\lambda e\alpha\tilde{v}\tau$	$lyhet$	$dönt$	

<sup>102</sup>We are being Anglocentric here: *báty* “elder brother”, *fivér* “male sibling”, *öcs* “younger brother” are merged into one class; the same applies to *nővéر* “elder sister” and *húg* “younger sister”. This is more a temporary compromise than a generally healthy practice. In fact, many Asian languages (including Chinese, Japanese, Korean, Khmer, Malay, Mongolian, Tamil, Telugu, Thai, Vietnamese etc.) possess distinct words for siblings (of given gender) depending on seniority.

mad(á a)r	madam	mam(a á)	mari~	maso~	más~	méh	mesél~	mrs~	mulat~	munk <sup>Xε</sup> (a á)~	muszl~
βipődr	μαδαμ	anya	μααρι	μασο	οθρς	βεεη	τελ	φυρχ	φυτ	ωωρκX	μισλ
nagybá(csi ty tty)~	nem~	nővér(Ø r)	ook	oor	orr~	ow	őrül	őrül	ősztön(Ø ös)(Ø sV)	ősztönöz	palánt(a á)
uñkli	ñeu	húg	owk	wor	ñeç	ow	mađl	pjoil	iñσtiñ	iñσtiñzi	πλαñta
papír	pár~	poh(á a)r	puszt	ree~	rég~	remé(l ny nny)~	rend	rossz	rovar	sajn~	szabad
παπιρ	πααρ	κοππρ	δεστρ	ρεε	ολδг	χοπλ	ορδ	βαš	ιñσεκτ	σορη	φριδ
száj~	szalon~	szám	szem~	szer	szeré(Ø n)ny~	szív	szob(a á)	tábor	táj~	táv	
μυñđj	ζαλων	száum	eyem	szrer	μοδστт	κορδαв	ρομ	καμπορ	ρεγ्यј	φάρν	
tavasz~	te(a á)~		te(ss tsz tssz)		tekint	terem	terhet	term	terv	tetszés	téved
leþpsz	τζаji		λικγνδ		nézt	τερεм	τερhek	τερм	πλανv	σατφis	εорed
tit(Ø o)k~	tojás	törvény~	vad	vála(szd ssz szt)X		váll	való	vált	változX	város	
σεκрk	eγys	λаwny	wiłd	σελεкт		σοhöл	iñδεεδó	ξañγet	ξañγet	κιτηs	
véd	vétk	világ	vissz~	zavar	zs(á é)g	cs	jane	miss	sir	tv	~kénti
προκτод	σiñik	ωελтег	βк	κοñφσr	z	γυρπt	γοιñvα	μiσx	σiρx	λаwny	∅

(2) Do  $\alpha a\sim \rightarrow \eta\mu t$ ,  $(Ø|el|meg)vár \rightarrow \omega aatp$ , (Ø|meg)eatHungarian  $\rightarrow a\epsilon a\tau$ , drinkHungarian  $\rightarrow e\delta\sigma i\tilde{v}k$ , mé(gy|sz)  $\rightarrow menni$ , ~buyHungarian  $\rightarrow h\beta uñh$ , ~carryHungarian  $\rightarrow h\kappa\sigma taph$

(3) Replace

(Ø fel meg)kér~	(atlan etlen talan telen tlan tlen)	(két kettő)~	ablak	ajt~	álm	almo~										
kér	zτ	2két	ablaκ	ajτ	aulm	almo										
álon	alom~	barát~	beszé(d)l	egy~	el~	év~	gyerek~	hall~	hib	ifjú~	igy~	jó~	kedv	kert~		
aulom	alom	βarát	βeszéλ	zgy	l	yév	gyerek	χalλ	hiβ	ifja	qgy	jo	keðβ	kert		
kicsi~	leány~	lép	lev~	ond	orvos	szebb	szó~	tér~	tiszt <sup>Xε</sup> (a á i)o	tud	új~	úr~	vagyon-	mr	~l(a)e)g	~th
kis	lány	lép	lev	onð	δokt	széppb	szav	tép	tisztX	tuð	wúj	ur	wagyon	xñep	∅	t

(4) Do  $tiszt \rightarrow tiszt\tau$ ,  $(Ø|leg|legesleg)Xbb \rightarrow X'$ , where one constructs  $X'$  by doing  $\sim\acute{a} \rightarrow a$ ,  $\sim\acute{e} \rightarrow e$  on  $X$ .

(5) Replace

ccs	ggy	lly	nny	ssz	tty	zzs
čč	ðð	łł	ññ	śś	þþ	żż

(6) Replace

cs	gy	ly	ny	sz	ty	zs
č	ð	ł	ñ	ś	þ	ż

and call the result  $\hat{\sigma}'$ .

(7) Break down  $\hat{\sigma}' = \hat{\sigma}_1\hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}^{[\text{ProtRg}(\hat{\sigma}')]}$  (see the notation in Definition 3.1) and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma}')$  is equal to the protected range of  $\hat{\sigma}'$ .

(8) On  $\hat{\sigma}_1$ , do  $(be|fel|le|meg) \sim \rightarrow \emptyset$ ,  $ki\hat{x}\notin(s) \sim \rightarrow \hat{x}$ , and call the result  $\hat{\sigma}'_1$ .

(9) On  $\hat{\sigma}_2$ , perform the following substitutions in a sequel:

(9.1) Do  $\sim(ad|al(Ø|o)m|áš|beli|dal(Ø|o)m|del(Ø|e)m|ék(e|o)ñ|el(Ø|e)m|éš|h(a|e)t|ist(a|á)|ít|ñi|omás|s(á|é)g|šerü)X \rightarrow \emptyset$ ,  $\sim(Ø|m|v)(á|é)\ñX \rightarrow \emptyset$ ,  $\sim izá(ció|l)X \rightarrow \emptyset$ ,  $\sim(st(u|ü)l)X \rightarrow \emptyset$ ,  $\sim(anak|enek)j(á|é)k|n(á|é)(k|m)) \rightarrow \emptyset$ ,  $\sim((á|e|é|o)s|k(e|o|ö)(d|z))X \rightarrow \emptyset$ ,  $\sim t(á|am|él|em) \rightarrow \emptyset$ ,  $\sim ni(Ø|a|e) \rightarrow \emptyset$ ,  $\sim \hat{x}\times 2(a|e)l \rightarrow \hat{x}$ ,  $\sim \hat{x}\notin(l|r)h(a|e)t \rightarrow \hat{x}$ ,  $\sim(\Ø|and)ó \rightarrow \emptyset$ ,  $\sim(\Ø|e|en)(Ø|d)ó \rightarrow \emptyset$ ,  $\sim l(a|e)k \rightarrow \emptyset$ ,  $\sim Vtt \rightarrow \emptyset$ .

(9.2) Do  $\sim^X(Ø|á|é)((a|e|o|ö)n(Ø|ként)|(a|e|o|ö)t|b(a|e)(Ø|n)|b(ó|ö)l|ért|h(e|o|ö)z|ig|ként|képp(Ø|en)) \rightarrow X'$ ,  $\sim^X(Ø|á|é)(n|n(a|e)k|n(á|é)l|r(a|e)|(r|t)(ó|ö)l|t|(u|ü)l|v(á|é)|v(a|e)l) \rightarrow X'$ , where one constructs  $X'$  by doing  $\sim\acute{a} \rightarrow a$ ,  $\sim\acute{e} \rightarrow e$  on  $X$ .

(9.3) Do  $\sim án(Ø|ként) \rightarrow a$ ,  $\sim énként \rightarrow e$ ,  $\sim lom \rightarrow lm$ ,  $\sim(n(a|á|ának|ánk|e|é|ének|énk)|t(á|é)k) \rightarrow \emptyset$ .

(9.4)  $\text{Do } \sim\text{ass}\mathbf{X} \rightarrow \emptyset, \sim\ddot{s}\mathbf{X} \rightarrow \check{s}, \sim(a|e|o|\ddot{o})st \rightarrow \emptyset, \sim\acute{a}st \rightarrow a, \sim\acute{e}st \rightarrow e.$

(9.5)  $\text{Do } \sim\hat{\chi}_{\times 2}(\acute{a}|\acute{e}) \rightarrow \hat{\chi}, \sim\mathbf{C}(d|i|j) \rightarrow \mathbf{C}, \sim n(e|o)d \rightarrow \emptyset.$

The result after these five steps of operations is called  $\hat{\sigma}'_2$ .

(10) Concatenate  $\hat{\sigma}'_1$  and  $\hat{\sigma}'_2$ .

**Algorithm 8.19** (Hungarian essential root). Let  $\hat{\sigma}$  be the effective spelling of a Hungarian word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:

(1)  $\text{Do } \sim\hat{\chi}_{\times 2}(\acute{a}|\acute{e})(\emptyset|l) \rightarrow \hat{\chi}, \sim(e|o|\ddot{o})l \rightarrow \emptyset, \text{ and call the result } \hat{\sigma}'.$

(2) Break down  $\hat{\sigma}' = \hat{\sigma}_1\hat{\sigma}_2$  into the concatenation of two strings  $\hat{\sigma}_1 = \hat{\sigma}^{[\text{ProtRg}(\hat{\sigma}')]} (see the notation in Definition 3.1)$  and  $\hat{\sigma}_2$ , where the length of the first string  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma}')$  is equal to the protected range of  $\hat{\sigma}'$ .

(3) On  $\hat{\sigma}_2$ , perform the following substitutions in a sequel:

(3.1)  $\text{Do } \sim dd \rightarrow \emptyset, \sim(\emptyset|a|e|o|\ddot{o})ké \rightarrow \emptyset, \sim\acute{e}(\emptyset|i) \rightarrow \emptyset, \sim\acute{a}(\acute{e}|i|ké) \rightarrow a, \sim\acute{e}(\acute{e}|i|ké) \rightarrow e;$

(3.2)  $\text{Do } \sim(\emptyset|a|e|o|\ddot{o})d \rightarrow \emptyset, \sim(\emptyset|an|n)(\emptyset|a|e|o)m \rightarrow \emptyset, \sim(\emptyset|j)(a|e) \rightarrow \emptyset, \sim(\emptyset|j|ni|t)(u|\ddot{u})k \rightarrow \emptyset, \sim(\emptyset|u|\ddot{u})nk \rightarrow \emptyset, \sim\acute{a}(\emptyset|djuk|m|nk) \rightarrow a, \sim\acute{e}(\emptyset|djük|m|nk) \rightarrow e, \sim o \rightarrow \emptyset;$

(3.3)  $\text{Do } \sim(\emptyset|(\emptyset|j)(ai|o))tok \rightarrow \emptyset, \sim(\emptyset|j)(\emptyset|e)(\emptyset|i)tek \rightarrow \emptyset, \sim(\emptyset|j)(a|e)(i|id|ik|im|ink) \rightarrow \emptyset, \sim\acute{a}(i|id|ik|im|ink|tok) \rightarrow a, \sim\acute{e}(i|id|ik|im|ink|tek) \rightarrow e, \sim tök \rightarrow \emptyset;$

(3.4)  $\text{Do } \sim\acute{a}k \rightarrow a, \sim\acute{e}k \rightarrow e, \sim(a|e|i|o|\ddot{o})k \rightarrow \emptyset.$

The result after these four steps of operations is called  $\hat{\sigma}'_2$ .

(4) Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}'_2$ .

(5)  $\text{Do } \sim tt(\emptyset|\mathbf{V}) \rightarrow \emptyset.$

(6)  $\text{Do } \sim\hat{\chi}(j|n|\check{s}|t|v)(\emptyset|\mathbf{V})(\emptyset|l) \rightarrow \emptyset.$

(7)  $\text{Do } \hat{\chi}_{\times 2} \rightarrow \hat{\chi}, \sim\hat{\chi}^\epsilon(l|r)h \rightarrow h\hat{\chi}.$

(8) Replace

$kis\sim$	$(\acute{d}e \acute{d}er jö jö jöh jöj jö\check{s})$	$\xrightarrow{\mathbf{X}\epsilon(\emptyset h v)i(\emptyset \acute{d} \acute{d}e \acute{d}él h \check{s})}$	$\xrightarrow{\mathbf{X}\epsilon(l v)(\acute{e} e e\acute{d} \acute{e}\acute{d} e\acute{d}\acute{e}l eh)}$	$\xrightarrow{me(\emptyset \acute{d} h)}$	$\xrightarrow{ve(\acute{d}e \check{s})}$
$\kappa i\sigma$	$jöñ$	$\mathbf{Xin}$	$\mathbf{Xen}$	$men$	$ven$

For clustering purposes, we will also define an operation to shorten the last long vowel in the effective spelling of a Hungarian word.

**Algorithm 8.20** (Length reduction of final vowel). Let  $\hat{\sigma}$  be the effective spelling of a Hungarian word, then  $\text{RdFinLV}(\hat{\sigma})$  is defined by the following procedure:

- Look for the pattern  $\mathbf{V}$  in  $\hat{\sigma}$ .
- If the pattern is not found, do nothing. Otherwise, do  $(\acute{a}|o|\ddot{o}) \rightarrow a, \acute{e} \rightarrow e, \ddot{o} \rightarrow \ddot{o}$  on its last occupying position.

### 8.2.2 Admissible mutation and approximate clustering

**Algorithm 8.21** (Simple heredity test). Construct  $\hat{\sigma}'$  by doing  $\sim(d|e|i|n|t|v) \rightarrow \emptyset$  on  $\hat{\sigma}$ . Define  $\hat{\beta}'$  similarly. The Boolean-valued function  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if the lowercase form of  $\hat{\alpha}$  contains at least one instance of  $\{a, \acute{a}, e, \acute{e}, i, \acute{i}, o, \acute{o}, \ddot{o}, u, \acute{u}, \ddot{u}, \acute{\ddot{u}}\}$  AND at least one of the following three conditions holds:<sup>103</sup>

- (i)  $\hat{\alpha} = \hat{\beta};$
- (ii)  $\hat{\beta} = \hat{\alpha}(d|m|k|v);$
- (iii)  $\ell(\hat{\beta}) > \ell(\hat{\alpha}) \geq \frac{\ell(\hat{\beta})}{2}$  AND  $\hat{\beta} = \hat{\alpha}\mathbf{V}\mathbf{X}.$

In what follows, we define  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta})$  and  $\text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta})$  as what is done in the Danish case (Algorithm 5.7), which was also applied to two other Germanic languages treated in §5.

<sup>103</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical **OR**.

**Algorithm 8.22** (Admissible suffix mismatch and vowel alternation). *For two strings  $\hat{\alpha}$  and  $\hat{\beta}$ , the Boolean-valued function*

$$\text{AdmMut}(\text{RootNW}(\hat{\alpha}, \hat{\beta}), \text{SuffixNW}(\hat{\alpha}, \hat{\beta}), \text{NW}^*(\hat{\alpha}, \hat{\beta}))$$

*returns TRUE if at least one of the following two conditions holds:*

- (i)  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta}) = [\emptyset | \mathbf{V}k, \emptyset | \mathbf{V}k] \text{ AND } (\text{NW}^*(\hat{\alpha}, \hat{\beta}) = [\emptyset, e] | [\emptyset, o] | [\emptyset, \ddot{o}] | [a, \acute{a}] | [e, \acute{e}] | [o, \ddot{o}] \text{ OR } \text{NW}^*(\hat{\beta}, \hat{\alpha}) = [\emptyset, e] | [\emptyset, o] | [\emptyset, \ddot{o}] | [a, \acute{a}] | [e, \acute{e}] | [o, \ddot{o}]) \text{ OR } \text{NW}^*(\hat{\beta}, \hat{\alpha}) = [\emptyset, e] | [\emptyset, o] | [\emptyset, \ddot{o}] | [a, \acute{a}] | [e, \acute{e}] | [o, \ddot{o}] \text{ OR } (\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \emptyset \text{ AND the lowercase form of } \text{RootNW}(\hat{\alpha}, \hat{\beta}) \text{ contains at least one instance of } \mathbf{V})$ ;
- (ii)  $\text{NW}^*(\hat{\alpha}, \hat{\beta}) = \emptyset \text{ AND } \text{SuffixNW}'(\hat{\alpha}, \hat{\beta}) = [a, \ddot{o}] | [ava, \ddot{o}] | [\ddot{o}, \ddot{o}] | [\ddot{o}, \ddot{o}]$ , where one obtains  $\text{SuffixNW}'(\hat{\alpha}, \hat{\beta})$  by doing  $\sim(k|st|v) \rightarrow \emptyset$  on both components of  $\text{SuffixNW}(\hat{\alpha}, \hat{\beta})$ .

Similarly, one can evaluate another Boolean-valued function

$$\text{AdmMut}(\text{RootSW}(\hat{\alpha}, \hat{\beta}), \text{SuffixSW}(\hat{\alpha}, \hat{\beta}), \text{SW}^*(\hat{\alpha}, \hat{\beta}))$$

by trading all the occurrences of NW in the statements above with SW.

**Algorithm 8.23** (Heredity test function). *The structure of the Hungarian heredity test function  $\text{HrdTest}(\hat{\alpha}, \hat{\beta})$  is identical to the German version (Algorithm 8.1.2), except that the functions  $\text{SimpHrdTest}$ ,  $\text{RootNW}$ ,  $\text{SuffixNW}$ ,  $\text{NW}^*$ ,  $\text{RootSW}$ ,  $\text{SuffixSW}$ ,  $\text{SW}^*$  must follow the Hungarian rules stated above.*

The structure of the following clustering algorithm differs from the Finnish counterpart only in cosmetic details.

**Algorithm 8.24** (Approximate clustering of Hungarian words). *The approximate clustering of a list of Hungarian words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  is completed in two stages:*

- (1) *We sort the list  $\{(\hat{\alpha}_1, \text{EffSpell}(\hat{\alpha}_1), \text{RdFinLV}(\text{EffSpell}(\hat{\alpha}_1))), \dots, (\hat{\alpha}_N, \text{EffSpell}(\hat{\alpha}_N), \text{RdFinLV}(\text{EffSpell}(\hat{\alpha}_N)))\}$  alphabetically according to the third component (with higher priority) and the second component (with lower priority). If two consecutive neighbors in the alphabetized list  $\{(\hat{\alpha}_{(1)}, \text{EffSpell}(\hat{\alpha}_{(1)}), \text{RdFinLV}(\text{EffSpell}(\hat{\alpha}_{(1)}))), \dots, (\hat{\alpha}_{(N)}, \text{EffSpell}(\hat{\alpha}_{(N)}), \text{RdFinLV}(\text{EffSpell}(\hat{\alpha}_{(N)})))\}$  satisfies both  $\text{HrdTest}(\text{EffSpell}(\hat{\alpha}_{(n)}), \text{EffSpell}(\hat{\alpha}_{(n+1)})) = \text{FALSE}$  and  $\text{HrdTest}(\text{EffSpell}(\hat{\alpha}_{(n+1)}), \text{EffSpell}(\hat{\alpha}_{(n)})) = \text{FALSE}$ , where  $n \in \mathbb{Z} \cap [1, N]$ , then add a demarcation line between these two entries. In this way, the alphabetized list is divided into separate groups of words with tags:  $\{g_1 = \{(\hat{\alpha}_{(1,1)}, \text{EffSpell}(\hat{\alpha}_{(1,1)}), \dots), \dots, (\hat{\alpha}_{(1,n_1)}, \text{EffSpell}(\hat{\alpha}_{(1,n_1)}), \dots)\}, \dots, g_M = \{(\hat{\alpha}_{(M,1)}, \text{EffSpell}(\hat{\alpha}_{(M,1)}), \dots), \dots, (\hat{\alpha}_{(M,n_M)}, \text{EffSpell}(\hat{\alpha}_{(M,n_M)}), \dots)\}\}$ , where each sublist also preserves alphabetic order. In particular, we have  $\hat{\alpha}_{(1,1)} = \hat{\alpha}_{(1)}$  and  $\hat{\alpha}_{(M,n_M)} = \hat{\alpha}_{(N)}$ .*
- (2) *For each group of words  $g_m = \{(\hat{\alpha}_{(m,1)}, \text{EffSpell}(\hat{\alpha}_{(m,1)}), \dots, (\hat{\alpha}_{(m,n_m)}, \text{EffSpell}(\hat{\alpha}_{(m,n_m)}))\}$  where  $m \in \mathbb{Z} \cap [1, M]$ , we augment it into a tagged entry  $G_m = (g_m, \text{EffSpell}(\hat{\alpha}_{(m,1)}), \text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})), T_m)$ , where*

$$T_m = \text{RdFinLV}(\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))).$$

*The list  $\{G_1, \dots, G_M\}$  is sorted alphabetically, with respect to  $T_m$  (with highest priority),  $\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))$  (with medium priority), and  $\text{EffSpell}(\hat{\alpha}_{(m,1)})$  (with lowest priority). If two consecutive neighbors in the alphabetized list  $\{G_{(1)} = (g_{(1)}, \hat{\gamma}'_{(1)}, \hat{\gamma}''_{(1)}, \hat{\gamma}'''_{(1)}), \dots, G_{(M)} = (g_{(M)}, \hat{\gamma}'_{(M)}, \hat{\gamma}''_{(M)}, \hat{\gamma}'''_{(M)})\}$  satisfy*

$$\text{HrdTest}(\hat{\gamma}''_{(m)}, \hat{\gamma}''_{(m+1)}) = \text{FALSE}$$

**AND**

$$\text{HrdTest}(\hat{\gamma}'_{(m)}, \hat{\gamma}'_{(m+1)}) = \text{FALSE}$$

*where  $m \in \mathbb{Z} \cap [1, M]$ , then add a demarcation line between these two entries. In this manner, the alphabetized list  $\{G_{(1)}, \dots, G_{(M)}\}$  is divided into separate groups of tagged clusters  $\{\Gamma_1 = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, \Gamma_K = \{G_{(K,1)}, \dots, G_{(K,m_K)}\}\}$ . Finally, generate a list of word clusters  $\{\check{G}_1 = (\check{G}_{(1,1)}, \dots), \dots, \check{G}_K = (\check{G}_{(K,1)}, \dots)\}$  by discarding all the tags.*

*Example 8.24.1.* The following samples of regular nouns (declined in all cases, in singular and plural) are extracted from [48, Appendix 2].

álmaitok, álmaitokat, álmaitokba, álmaitokban, álmaitokból, álmaitokért, álmaitokhoz, álmaitokig, álmaitokká, álmaitokkal, álmaitoknak, álmaitoknál, álmaitokon, álmaitokra, álmaitokról, álmaitoktól, álmotok, álmotokat, álmotokba, álmotokban, álmotokból, álmotokért, álmotokhoz, álmotokig, álmotokká, álmotokkal, álmotokként, álmotoknak, álmotoknál, álmotokon, álmotokra, álmotokról, álmotoktól, álom — “your dream”;

*bokor, bokorba, bokorban, bokorból, bokorért, bokorhoz, bokorig, bokorként, bokornak, bokornál, bokorra, bokorrá, bokorral, bokorról, bokortól, bokrok, bokrokat, bokrokba, bokrokban, bokrokóból, bokrokért, bokrokhoz, bokrokig, bokrokká, bokrokkal, bokrokként, bokroknak, bokroknál, bokrokon, bokrokra, bokrokrol, bokroktól, bokron, bokronként, bokrostul, bokrot* — “bush”;

*dinnye, dinnyébe, dinnyében, dinnyéből, dinnyéért, dinnyéhez, dinnyéig, dinnyék, dinnyékbe, dinnyékben, dinnyékből, dinnyéken, dinnyeként, dinnyékért, dinnyéket, dinnyékhez, dinnyékig, dinnyékké, dinnyékkel, dinnyékként, dinnyéknak, dinnyéknél, dinnyékre, dinnyékről, dinnyéktől, dinnyén, dinnyének, dinnyénél, dinnyénként, dinnyére, dinnyéről, dinnyéstü, dinnyét, dinnyétől, dinnyévé, dinnyével* — “melon”;

*ház, házak, házakat, házakba, házakban, házakból, házakért, házakhoz, házakig, házakká, házakkal, házakként, házaknak, házaknál, házakon, házakra, házakról, házaktól, házanként, házastul, házat, házba, házban, házból, házért, házhoz, házig, házként, háznak, háznál, házon, házra, házról, háztól, házzá, házzal* — “house”;

*ház, házaim, házaimat, házaimba, házaimban, házaimból, házaimért, házaimhoz, házaimig, házaimmá, házaimmal, házaimnak, házaimnál, házaimon, házaimra, házaimról, házaimtól, házam, házamat, házamba, házamban, házamból, házamért, házamhoz, házamig, házamként, házammá, házammal, házamnak, házamnál, házammon, házamra, házamról, házamtól* — “my house”;

*iker, ikerbe, ikerben, ikerból, ikerért, ikerhez, ikerig, ikerként, ikernek, ikernél, ikerre, ikerré, ikerrel, ikerről, ikertől, ikek, ikekbe, ikekben, ikekből, ikeken, ikekért, ikeket, ikekhez, ikekgig, ikekké, ikekkel, ikeként, ikeknek, ikeknél, ikekre, ikekről, ikektől, ikren, ikenként, ikrestü, ikret* — “twin”;

*jutalmaink, jutalmaintak, jutalmaintakba, jutalmaintakban, jutalmaintakból, jutalmaintakért, jutalmaintakhoz, jutalmaintakig, jutalmaintakká, jutalmaintakkal, jutalmaintaknak, jutalmaintaknál, jutalmaintakon, jutalmaintakra, jutalmaintakról, jutalmaintaktól, jutalmunk, jutalmunkat, jutalmunkba, jutalmunkban, jutalmunkból, jutalmunkért, jutalmunkhoz, jutalmunkig, jutalmunkká, jutalmunkkal, jutalmunkként, jutalmunknak, jutalmunknál, jutalmunkon, jutalmunkra, jutalmunkról, jutalmunktól, jutalom* — “our reward”;

*kép, képeik, képeikbe, képeikben, képeikból, képeiken, képeikért, képeiket, képeikhez, képeikig, képeikké, képeikkel, képeiknek, képeiknél, képeikre, képeikről, képeiktől, képük, képükbe, képükben, képükből, képükért, képük, képükön, képükig, képükkel, képükként, képüknek, képüknel, képükön, képükre, képükör, képükötől* — “their picture”;

*kéz, kézbe, kézben, kézből, kezek, kezekbe, kezekben, kezeből, kezeken, kezekért, kezeket, kezekhez, kezekig, kezék, kezékkel, kezéként, kezeknek, kezknél, kezere, kezkről, kezektől, kezen, kezenként, kezért, kezestü, kezet, kékhez, kékig, kéként, kéknek, kéknel, kékre, kékrol, kékötől, kékkel — “hand”;*

*kő, kőbe, kőben, kőből, kőért, kőhöz, kőig, kőként, kőnek, kőnél, kőre, kőről, kőtől, kővér, kövek, kövekbe, kövekben, kövekből, köveken, kövekért, köveket, kövekhez, kövekig, kövekké, kövekkel, köveknak, köveknél, kövekre, kövekről, kövektől, kővel, kövenként, kövestü, követ, kövön* — “stone”;

*könyv, könyvbe, könyvben, könyvből, könyvek, könyvekbe, könyvekben, könyvekből, könyveken, könyvekért, könyveket, könyvekhez, könyvekig, könyvekké, könyvekkel, könyvekként, könyveknek, könyveknél, könyvekre, könyvekről, könyvektől, könyvenként, könyvér, könyvestü, könyvet, könyvhöz, könyvig, könyvként, könyvnek, könyvnél, könyvön, könyvre, könyvről, könyvtől, könyvvé, könyvvel* — “book”;

*lánya, lányastul, lányba, lányban, lányból, lányért, lányhoz, lányig, lányként, lánynak, lánynál, lánnyal, lányok, lányokat, lányokba, lányokban, lányokból, lányokért, lányokhoz, lányokig, lányokká, lányokkal, lányokként, lányoknak, lányoknál, lányokon, lányokra, lányokról, lányoktól, lányon, lányonként, lányra, lányról, lányt, lánytól* — “girl”;

*levél, leveled, leveledbe, leveledben, leveledből, leveleddé, leveleddel, leveleden, leveledért, leveledet, leveledhez, leveledig, leveledként, levelednek, levelednél, leveledre, leveledről, leveledtől, leveleid, leveleidbe, leveleidben, leveleidből, leveleiddé, leveleiddel, leveleiden, leveleidért, leveleidet, leveleidhez, leveleidig, leveleidnek, leveleidnél, leveleidre, leveleidről, leveleidtől* — “your letter”;

*lova, lovába, lovában, lovából, lováért, lovához, lovai, lovaiba, lovaiban, lovaiból, lovaiért, lováig, lovaihoz, lovaiig, lovaiként, lovain, lovainak, lovainál, lovaira, lovairól, lovait, lovaitól, lovaivá, lovaival, lovaként, lován, lovának, lovánál, lovára, lováról, lovától, lovává, lovával* — “his/her horse”;

őr, őrbe, őrben, őrből, őrért, őrestül, őrhöz, őrig, őrként, őrnek, őrnél, őrök, őrokbe, őrokben, őrókből, őrokért, őróket, őrókhöz, őrokig, őrokék, őrokkel, őrokéként, őroknek, őroknél, őrokön, őrokre, őrokrol, őróktől, őron, őronként, őre, őrré, őrrel, őrről, őrt, őrtől — “guard”;

pohár, poharak, poharakat, poharakba, poharakban, poharakból, poharakért, poharakhoz, poharakig, poharakká, poharakkal, poharakként, poharaknak, poharagnál, poharakon, poharakra, poharakról, poharaktól, poharanként, poharastul, poharat, pohárba, pohárban, poháról, pohárhoz, pohárig, pohárként, pohárnak, pohárnál, poháron, pohárra, pohárral, pohárról, pohártól — “glass”;

szavak, szavakat, szavakba, szavakban, szavakból, szavakért, szavakhoz, szavakig, szavakká, szavakkal, szavakként, szavaknak, szavaknál, szavakon, szavakra, szavakról, szavaktól, szavanként, szavastul, szavon, szó, szóba, szóban, szóból, szóért, szóhoz, szóig, szók, szóként, szónak, szónál, szóra, szóról, szót, szótól, szóvá, szóval — “word”;

szék, székbe, székben, székből, székek, székekbe, székekben, székekkel, székeken, székekért, székekét, székekhez, székekig, székeké, székekkel, székeként, székeknek, székeknél, székekre, székekről, székektől, széken, székenként, székért, székestül, széket, székhez, székg, székké, székként, széknek, széknél, székre, székről, széktől — “chair”;

táska, táskába, táskában, táskából, táskáért, táskához, táskáig, táskák, táskákat, táskákba, táskákban, táskákból, táskaként, táskákért, táskákhoz, táskákig, táskákká, táskákkal, táskákként, táskáknak, táskánál, táskákon, táskákra, táskákról, táskáktól, táskán, táskának, táskánál, táskánként, táskára, táskáról, táskástul, táskát, táskától, táskává, táskával — “bag”;

tavak, tavakat, tavakba, tavakban, tavakból, tavakért, tavakhoz, tavakig, tavakká, tavakkal, tavakként, tavaknak, tavaknál, tavakon, tavakra, tavakról, tavaktól, tavanként, tavastul, tavat, tavon, tó, tóba, tóban, tóból, tóért, tóhoz, tóig, tóként, tónak, tónál, tóra, tóról, tótól, tóvá, tóval — “lake”;

tükör, tükörbe, tükörben, tükörből, tükörért, tükörhöz, tükörig, tükörként, tükörnek, tükörnél, tükörre, tükörré, tükörrel, tükörről, tükörtől, tükörök, tükörökbe, tükörökben, tükörökkel, tükörkért, tükörköt, tükörkhöz, tükörkig, tükörkék, tükörkékel, tükörkéként, tükörknek, tükörknél, tükörön, tükörkére, tükörkötől, tükörkötől, tükörön, tüköröként, tüköröstü, tüköröt — “mirror”.

These regular nouns are clustered correctly by our algorithm.

We then select three nouns listed in [48, Appendix 3, Section “Metathesis”], as their declensions deviate significantly from the regular patterns covered above.

kehely, kehelybe, kehelyben, kehelyből, kehelyért, kehelyhez, kehelyig, kehelyként, kehellyé, kehellyel, kehelynek, kehelynél, kehelyre, kehelyről, kehelytől, kelyhek, kelyhekbe, kelyhekben, kelyhekből, kelyheken, kelyhekért, kelyheket, kelyhekhez, kelyhekig, kelyhekké, kelyhekkel, kelyhekként, kelyheknek, kelyheknél, kelyhekre, kelyhekről, kelyhektől, kelyhen, kelyhet — “chalice”;

pehely, pehelybe, pehelyben, pehelyből, pehelyért, pehelyhez, pehelyig, pehelyként, pehellyé, pehellyel, pehelynek, pehelynél, pehelyre, pehelyről, pehelytől, pelyhek, pelyhekbe, pelyhekben, pelyhekből, pelyheken, pelyhekért, pelyheket, pelyhekhez, pelyhekig, pelyhekké, pelyhekkel, pelyhekként, pelyheknek, pelyheknél, pelyhekre, pelyhekről, pelyhektől, pelyhen, pelyhet — “flake”;

teher, teherbe, teherben, teherből, teherért, teherhez, teherig, teherként, tehernek, tehernél, teherre, teherré, teherrel, teheről, tehertől, terhek, terhekbe, terhekben, terhekből, terheken, terhekért, terheket, terhekhez, terhekip, terhekké, terhekkel, terhekként, terheknek, terheknél, terhekre, terhekről, terhektől, terhen, terhet — “load”.

They are automatically clustered by our algorithm, as expected.

The declension paradigms of Hungarian adjectives are essentially identical to those of Hungarian nouns. In the following, we only list various adjectives in their positive and comparative forms (in the nominative case, singular number), based on [48, §10.3]. Superlatives are not listed below, as they are consistently formed with a *leg~* prefix to the comparative forms.

bátor, bátrabb — “brave”;

bő, bővebb — “abundant”;

derék, derekabb — “decent”;

drága, drágább — “expensive”;

*édes, édesebb* — “sweet”;  
*érthatő, érthatóbb* — “understandable”;  
*fekete, feketébb* — “black”;  
*hosszabb, hosszú* — “long”;  
*hű, hűbb* — “faithful”;  
*ifjabb, ifjú* — “young”;  
*jó, jobb* — “good”;  
*keserű, keserűbb* — “bitter”;  
*kevés, kevesebb* — “few”;  
*kicsi, kisebb* — “small”;  
*könnyebb, könnyű* — “easy”;  
*különös, különösebb* — “special”;  
*lassabb, lassú* — “slow”;  
*nehéz, nehezebb* — “difficult”;  
*olcsó, olcsóbb* — “cheap”;  
*piros, pirosabb* — “red”;  
*régi, régibb* — “old”;  
*szebb, szép* — “beautiful”;  
*szomorú, szomorúbb* — “sad”;  
*szörnyebb, szörnyű* — “awful”.

Our clustering algorithm groups these adjectives correctly.

*Example 8.24.2.* Our sample of regular Hungarian verbs are chosen from a representative subset of [48, §4.1], so as to cover all the possible regular verb typologies.

*cselekedett, cselekedet, cselekedek, cselekedem, cselekedendő, cselekedett, cselekedhet, cselekedeti, cselekedik, cselekeditek, cselekedj, cselekedje, cselekedjed, cselekedjék, cselekedjél, cselekedjelek, cselekedjem, cselekedjen, cselekedjenek, cselekedjetek, cselekedjétek, cselekedjük, cselekedjünk, cselekedlek, cselekedne, cselekedné, cselekedned, cselekednéd, cselekednek, cselekednék, cselekednél, cselekednélek, cselekednem, cselekedném, cselekednének, cselekednénk, cselekednetek, cselekednétek, cselekedni, cselekednie, cselekedniük, cselekednünk, cselekedő, cselekedsz, cselekedte, cselekedted, cselekedtek, cselekedték, cselekedtél, cselekedtelek, cselekedtem, cselekedtetek, cselekedtétek, cselekedtük, cselekedtünk, cselekedünk, cselekedve, cselekszed, cselekszel, cselekszem, cselekszenek, cselekszetek, cselekszi, cselekszik, cselekszíték, cselekszünk, cselekvő* — “act”;

*emlékezett, emlékezel, emlékezett, emlékezhet, emlékezne, emlékezned, emlékeznek, emlékeznék, emlékeznél, emlékez nem, emlékeznének, emlékeznénk, emlékeznetek, emlékeznétek, emlékezni, emlékeznie, emlékeznitük, emlékezünk, emlékező, emlékeztek, emlékeztél, emlékeztem, emlékeztetek, emlékeztetek, emlékeztünk, emlékezünk, emlékezve, emlékezz, emlékezzek, emlékezzél, emlékezzen, emlékezzenek, emlékezzetek, emlékezzünk, emlékszel, emlékszem, emlékszenek, emlékszetek, emlékszik, emlékszünk* — “remember”;

*gyarapodhat, gyarapodik, gyarapodj, gyarapodjak, gyarapodjál, gyarapodjanak, gyarapodjatok, gyarapodjon, gyarapodjunk, gyarapodna, gyarapodnak, gyarapodnál, gyarapodnának, gyarapodnánk, gyarapodnátok, gyarapodnák, gyarapodni, gyarapodnia, gyarapodniuk, gyarapodnod, gyarapodnom, gyarapodnotok, gyarapodnunk, gyarapodó, gyarapodok, gyarapodott, gyarapods, gyarapodtak, gyarapodtál, gyarapodtam, gyarapodtatok, gyarapodtok, gyarapodtunk, gyarapodunk, gyarapodva* — “increase”;

*igyekezek, igyekezel, igyekezett, igyekezhet, igyekezne, igyekezned, igyekeznek, igyekeznél, igyekeznem, igyekeznének, igyekeznénk, igyekeznetek, igyekeznétek, igyekezni, igyekeznie, igyekezniük, igyekeznünk, igyekező, igyekeztek, igyekeztél, igyekeztem, igyekeztetek, igyekeztünk, igyekezünk, igyekezve, igyekezz, igyekezzek, igyekezzél, igyekezzen, igyekezzenek, igyekezzetek, igyekezziük, igyeksznel, igyekszem, igyekszenek, igyeksztek, igyekszik, igyekszünk, igyekvő — “strive”;*

*mosolygandó, mosolygó, mosolygod, mosolygok, mosolygom, mosolygott, mosolygunk, mosolyog, mosolyogd, mosolyoghat, mosolyogj, mosolyogja, mosolyogjad, mosolyogjak, mosolyogják, mosolyogjál, mosolyogjalak, mosolyogjam, mosolyogjanak, mosolyogjatok, mosolyogjátok, mosolyogjon, mosolyogjuk, mosolyogjunk, mosolyoglak, mosolyogna, mosolyogná, mosolyognad, mosolyognak, mosolyognák, mosolyognál, mosolyognálak, mosolyognám, mosolyognának, mosolyognánk, mosolyognátok, mosolyognék, mosolyogni, mosolyognia, mosolyogniuk, mosolyognod, mosolyognom, mosolyognotok, mosolyognunk, mosolyogsz, mosolyogta, mosolyogtat, mosolyogtak, mosolyogták, mosolyogtál, mosolyogtalak, mosolyogtam, mosolyogtatok, mosolyogtok, mosolyogtuk, mosolyogtunk, mosolyogva — “smile”;*

*sző, sződd, szőj, szője, szőjed, szőjek, szőjék, szőjel, szőjek, szőjel, szőjém, szőjenek, szőjetek, szőjétek, szőjön, szőjük, szőjünk, szőlek, szőne, szőné, szőnéd, szőnek, szőnék, szőnél, szőnélek, szőném, szőnének, szőnénk, szőnétek, szőni, szőnie, szőniük, szőnök, szőnöd, szőnöm, szőnötök, szőnünk, szősz, szőtt, szőtte, szőtted, szőttek, szőtték, szőttél, szőttelek, szőttem, szőttetek, szőtték, szőttök, szőttünk, szőve, szövendő, szövi, szövik, szövitek, szövő, szövöd, szövök, szövöm, szövünk — “weave”;*

*utazhat, utazik, utazna, utaznak, utaznál, utaznának, utaznánk, utaznátor, utaznák, utazni, utaznia, utazniuk, utaznod, utaznom, utaznotok, utaznunk, utazó, utazok, utazol, utazom, utazott, utaztak, utaztál, utaztam, utazzatok, utaztok, utaztunk, utazunk, utazva, utazz, utazzak, utazzál, utazzanak, utazzatok, utazzon, utazzunk — “travel”.*

These verbs are clustered correctly by our algorithm.

Our selection of irregular Hungarian verbs are based on [48, Appendix 1].

*edd, egye, egyed, egyek, egyék, egyél, egyelek, egyem, egyen, egyenek, egyetek, egyétek, együk, egylünk, ehet, enne, enné, enned, ennéd, ennék, ennél, ennélek, ennem, enném, ennének, ennénk, ennetek, ennétek, enni, ennie, enniük, ennünk, eszed, eszek, eszel, eszem, eszi, eszik, esztek, eszlek, esznek, esszük, esztek, eszünk, ette, etted, ettek, ették, ettél, ettelek, ettem, ettetek, ették, ettük, ettünk, evet, evő — “eat”;*

*gyere, gyertek, gyerünk, jöhét, jöjj, jöjjek, jöjjel, jöjjene, jöjjetek, jöjjön, jöjjünk, jön, jönne, jönnek, jönnék, jönnél, jönnének, jönnénk, jönnétek, jönni, jönnie, jönnük, jönnöd, jönnöm, jönnötök, jönnünk, jössz, jösztök, jött, jöttek, jöttel, jöttem, jöttetek, jöttök, jöttünk, jöve, jövendő, jövő, jövök, jövünk — “come”;*

*hidd, hiendő, higgy, higgye, higgyed, higgyek, higgyék, higgyél, higgyelek, higgyem, higgyen, higgenek, higgyetek, higgyétek, higgyük, higgyünk, hihet, hinne, hinné, hinned, hinnéd, hinnék, hinnél, hinnélek, hinnem, hinném, hinnének, hinnénk, hinnetek, hinnétek, hinni, hinnie, hinnük, hinnünk, hisz, hiszed, hiszek, hiszel, hiszem, hiszi, hiszik, hisztek, hiszlek, hisznek, hisszük, hisztek, hiszünk, hitt, hitte, hitted, hittek, hitték, hittél, hittelek, hittem, hittetek, hittétek, hittük, hittünk — “believe”;*

*idd, igya, igyad, igyak, igyák, igyál, igyalak, igyam, igyanak, igyatok, igyon, igyük, igyunk, ihat, inná, inná, innád, innák, innál, innálak, innám, innának, innánk, innátok, innék, inni, innia, inniük, innod, innom, innotok, innunk, iszik, iszak, isznak, iszod, iszok, iszol, iszom, issza, isszák, isszátok, isszuk, isztok, iszunk, itta, ittad, ittak, itták, ittalak, ittam, ittatók, ittuk, ittunk, ivandó, ivó, ivott — “drink”;*

*megy, megyek, megünk, lehet, menj, menjek, menjél, menjen, menjenek, menjetek, menjünk, menne, mennen, mennek, mennék, mennél, mennem, mennének, mennénk, mennenek, mennenek, mennyek, mennyünk, menő, ment, mentek, mentél, mentem, mentetek, mentük, menve — “go”;*

*vedd, veendő, végy, vegye, vegyed, vegyek, vegyék, vegyél, vegyelek, vegyem, vegyen, vegyenek, vegyetek, vegyétek, vegyük, vegyünk, vehet, venne, venné, venned, vennék, vennél, vennélek, vennem, venném, vennének, vennénk, vennétek, vennétek, vennétek, venni, vennie, vennük, vennünk, vesz, veszed, veszek, veszem, veszi, veszik, veszitek, veszlek, vesznek, vesszük, vesztek, veszünk, vett, vette, vettered, vettek, vették, vettél, vettelek, vettetem, vettetek, vettétek, vettük, vettünk, véve, vevő — “buy”;*

*vidd, viendő, vigye, vigyed, vigyek, vigyék, vigyél, vigyelek, vigyem, vigyen, vigyenek, vigyetek, vigyétek, vigyük, vigyünk, vihet, vinne, vinné, vinned, vinnék, vinnél, vinnélek, vinnem, vinném, vinnének, vinnénk, vinnetek, vinnétek, vinni, vinnie, vinnük, vinnünk, visz, viszed, viszek, viszel, viszem, viszi, viszik, viszitek, viszlek, visznek, visszük, visztek, viszünk, vitt, vitte, vittered, vittek, vitték, vittél, vittelek, vittem, vittetek, vitték, vittük, vittünk, vive, vivő — “carry”.*

Sending them to our algorithm, we receive the correct clustering results.

### 8.2.3 Heuristic detection of compounds

The following algorithm for heuristic detection of Hungarian compounds differs from the Danish version (Algorithm 5.12) only in some specific details. To make the context clear, we still state the algorithm in full. (In what follows, the string minus operation  $\hat{\beta} \ominus \hat{\alpha}$  is prescribed by Definition 5.11.)

**Algorithm 8.25** (Heuristic identification of Hungarian binary compounds). *Let  $\Lambda^{\hat{\beta}} = \{\hat{\rho}_1, \dots, \hat{\rho}_Q\}$  be a list of distinct Hungarian essential roots (without vowel blotting) that are at two characters long, contain at least one instance of V in the first two letters, and DO NOT match either*

(ag|ar|az|ba|bo|cel|či|de|dö|dul|em|ha|he|hi|ig|is|ka|ke|ko|köt|li|mas|mi|mu|ne|od|ős|rá|re|ta|te|tor|tü|va|vá|vas)

or

(β|ka|(a|e|o|ö)(s|z)|al|am|di|eg|er|ér|fe|fog|hí|kor|mág|még|ol|or|ot|sV|ter|ve)X.

The output of the function  $\text{CpdDet}(\Lambda^{\hat{\beta}})$  is obtained through the following procedures:

- (1) Construct a list  $\{(\hat{\rho}_1, \lambda_1^{\hat{\beta}}), \dots, (\hat{\rho}_Q, \lambda_Q^{\hat{\beta}})\}$  where  $\lambda_q^{\hat{\beta}} = \{\hat{\rho}_{(q,1)}, \dots, \hat{\rho}_{(q,n_q)}\}$  is a subset of  $\Lambda^{\hat{\beta}}$  whose members all match the string pattern  $\hat{\rho}_{q\sim}$ , for  $q \in \mathbb{Z} \cap [1, Q]$ .
- (2) Expand the aforementioned entry  $(\hat{\rho}_q, \lambda_q^{\hat{\beta}})$  into a list of triplets  $\{(\hat{\rho}_{(q,1)}, \hat{\rho}_q, \hat{\rho}_{(q,1)} \ominus \hat{\rho}_q), \dots, (\hat{\rho}_{(q,n_q)}, \hat{\rho}_q, \hat{\rho}_{(q,n_q)} \ominus \hat{\rho}_q)\}$  for every  $q \in \mathbb{Z} \cap [1, Q]$  such that  $\lambda_q^{\hat{\beta}} \neq \emptyset$ . Collect all these triplets as one runs through the list  $\{(\hat{\rho}_1, \lambda_1^{\hat{\beta}}), \dots, (\hat{\rho}_Q, \lambda_Q^{\hat{\beta}})\}$ . The list of these triplets  $\{(\hat{\rho}_{(1)}, \hat{\eta}_{(1)}, \hat{\rho}_{(1)} \ominus \hat{\eta}_{(1)}), \dots, (\hat{\rho}_{(Q')}, \hat{\eta}_{(Q')}, \hat{\rho}_{(Q')} \ominus \hat{\eta}_{(Q')})\}$  contains potentially valid decompositions of compounds.
- (3) Screen the aforementioned list of triplets as follows: for every  $q' \in \mathbb{Z} \cap [1, Q']$ , if  $(\hat{\rho}_{(q')}, \hat{\eta}_{(q')}, \hat{\tau}_{(q')} = \hat{\rho}_{(q')} \ominus \hat{\eta}_{(q')})$  satisfies

$$\ell(\hat{\rho}_{(q')} \ominus \hat{\eta}_{(q')}) \geq 2 \quad \text{AND} \quad \hat{\rho}_{(q')} \ominus \hat{\eta}_{(q')} = \mathbf{X}_1 \mathbf{V} \mathbf{X}_2,$$

then construct  $\hat{\tau}_{(q')}$  by performing  $(\emptyset | a | e) (\emptyset | i | n | i | s | t | u | ü) \sim \rightarrow \emptyset$  and

$(\lambda l(\emptyset | e | ö) | a | á | a | ö | a | l | a | b | e | f | \emptyset | e | e | l | él) | h | á | e | l | ö | z | id | k | e | r | e | l | ö | z | le | \emptyset | n | meg | ne | k | i | o | d | ö | s | r | á | š | é | š | e | r | t | o | t | ü | u | v | é | g | v | i | n) \sim \rightarrow \emptyset$

on  $\hat{\tau}_{(q')}$ , before generating a list  $\lambda_{(q')}^{\hat{\tau}}$  by members of  $\Lambda^{\hat{\beta}}$  that match the pattern  $(\hat{\tau}_{(q')} | \hat{\tau}_{(q')}^*)$ ; otherwise, set  $\lambda_{(q')}^{\hat{\tau}} = \emptyset$ .

- (4) Collect all the triplets  $(\hat{\rho}_{(q')}, \hat{\eta}_{(q')}, \lambda_{(q')}^{\hat{\tau}})$  where  $\lambda_{(q')}^{\hat{\tau}}$  is non-void and  $\hat{\tau}_{(q')}$  DOES NOT match  $z\tau X$ . This list of triplets  $\text{CpdDet}(\Lambda^{\hat{\beta}})$  contains the heuristic decompositions of all the identified binary compounds.

**Algorithm 8.26** (Approximate clustering of Hungarian words with heuristic detection of compounds). *The approximate clustering of a list of Hungarian words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  respecting compounding is completed in four stages:*

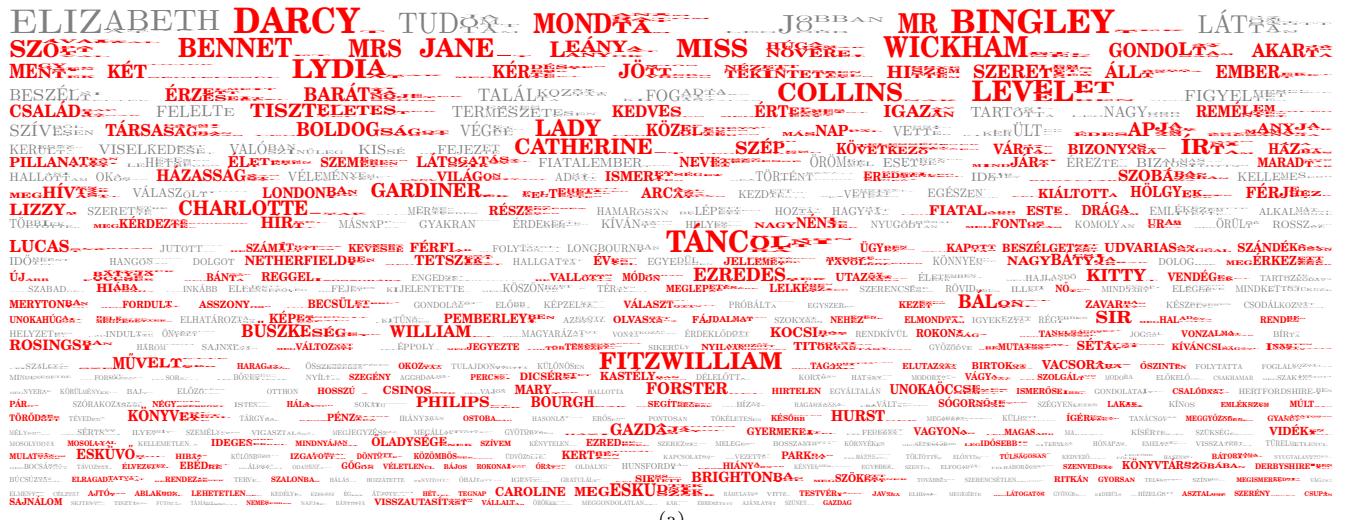
- (1) Do as in Algorithm 8.24(1).
- (2) Do as in Algorithm 8.24(2). Save both the tagged clusters  $\{\Gamma_1 = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, \Gamma_K = \{G_{(K,1)}, \dots, G_{(1,m_K)}\}\}$  and the list of word clusters  $\{\check{\Gamma}_1, \dots, \check{\Gamma}_K\}$  for further use.
- (3) Construct a tagged list of word clusters  $\{(\check{\Gamma}_1, \Lambda_1^{\hat{\beta}}), \dots, (\check{\Gamma}_K, \Lambda_K^{\hat{\beta}})\}$  where  $\Lambda_k^{\hat{\beta}}$  is the union of all Hungarian essential roots (without vowel blotting) available to  $\Gamma_k$ , for  $k \in \mathbb{Z} \cap [1, K]$ . Set  $\Lambda^{\hat{\beta}} = \Lambda_1^{\hat{\beta}} \cup \dots \cup \Lambda_K^{\hat{\beta}}$ , and evaluate  $\text{CpdDet}(\Lambda^{\hat{\beta}})$ .
- (4) The first component  $\hat{\rho}_{(q'')}$  of each triplet  $(\hat{\rho}_{(q'')}, \hat{\eta}_{(q'')}, \lambda_{(q'')}^{\hat{\tau}})$  in  $\text{CpdDet}(\Lambda^{\hat{\beta}})$  is called a “dissolvable compound”, the second component  $\hat{\eta}_{(q'')}$  a “heuristic head”, and the first member in the third component  $\lambda_{(q'')}^{\hat{\tau}}$  a “heuristic tail”. In the tagged list  $\{(\check{\Gamma}_1, \Lambda_1^{\hat{\beta}}), \dots, (\check{\Gamma}_K, \Lambda_K^{\hat{\beta}})\}$ , every entry containing a “dissolvable compound” is removed, and regrouped with the entries matching its “heuristic head” and “heuristic tail”. Finally, remove all tags.

*Example 8.26.1.* In Fig. S13, we apply the aforementioned word clustering algorithm to topic extraction and machine translation (see Table S1 for text sources).

According to Hungarian orthography, vowel lengths (a vs. á, ö vs. ö etc.) are ignored during alphabetization. If implemented as is, stacked vowels in word clusters may become illegible at times. As a compromise, we always sort long vowels after short vowels in Fig. S13.

In Hungarian, *de* “but” is a stop word, so it is understandable that our algorithm does not provide a perfect translation of the proper name *de Bourgh* in Fig. S13b.

In Darwin’s *Origin of Species*, the word “reciprocal” is mostly used in the phrase *reciprocal cross*. Therefore, we consider “reciprocal” and “graft” near synonyms in Fig. S13b''. Furthermore, in Fig. S13b'', sometimes an English word may correspond to a Hungarian phrase or vice versa: we consider English *fresh* an exact match to Hungarian *édesvízi* “fresh water” (literally “sweet water”), and similarly for English *rock-pigeon* and Hungarian *szirti* (the latter being part of the phrase *szirti galamb* that translates *rock-pigeon* into Hungarian).



(a)

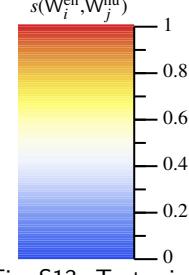
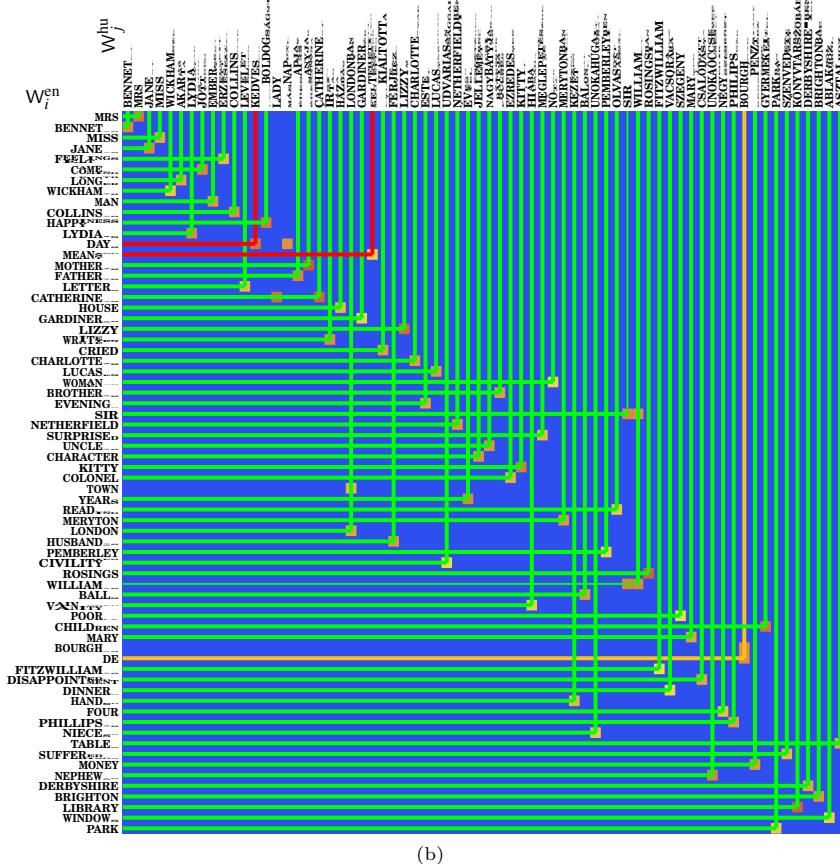


Fig. S13. Text mining in Hungarian. (a) Statistically identified topics ( $n_{ii} \geq 20$ ) in a Hungarian version of *Pride and Prejudice*, with the same color encoding scheme as Fig. S3. (b) Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{hu}})$  between selected topics in English and Hungarian versions of *Pride and Prejudice*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) ranking of the similarity score for the optimal node. Green (resp. red) cross-hair indicates a correct (resp. incorrect) match. Amber cross-hair marks a link between distinct concepts that share the same hypernym.

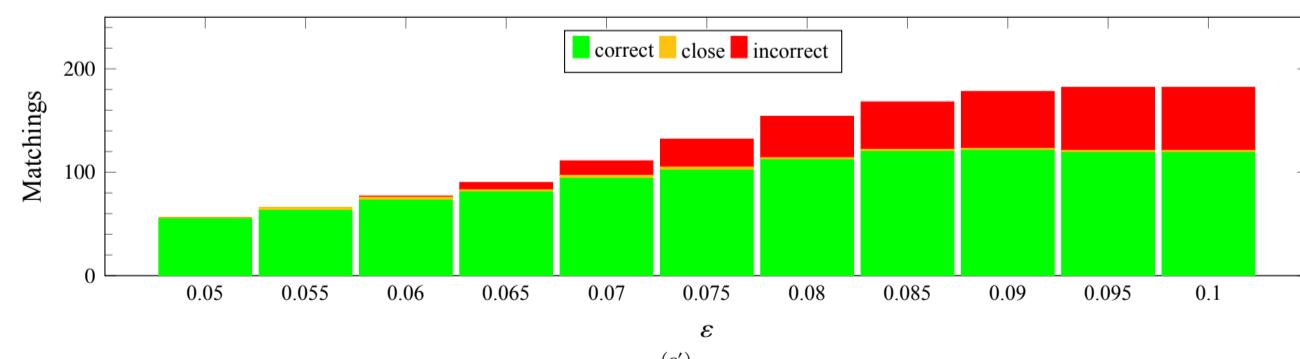
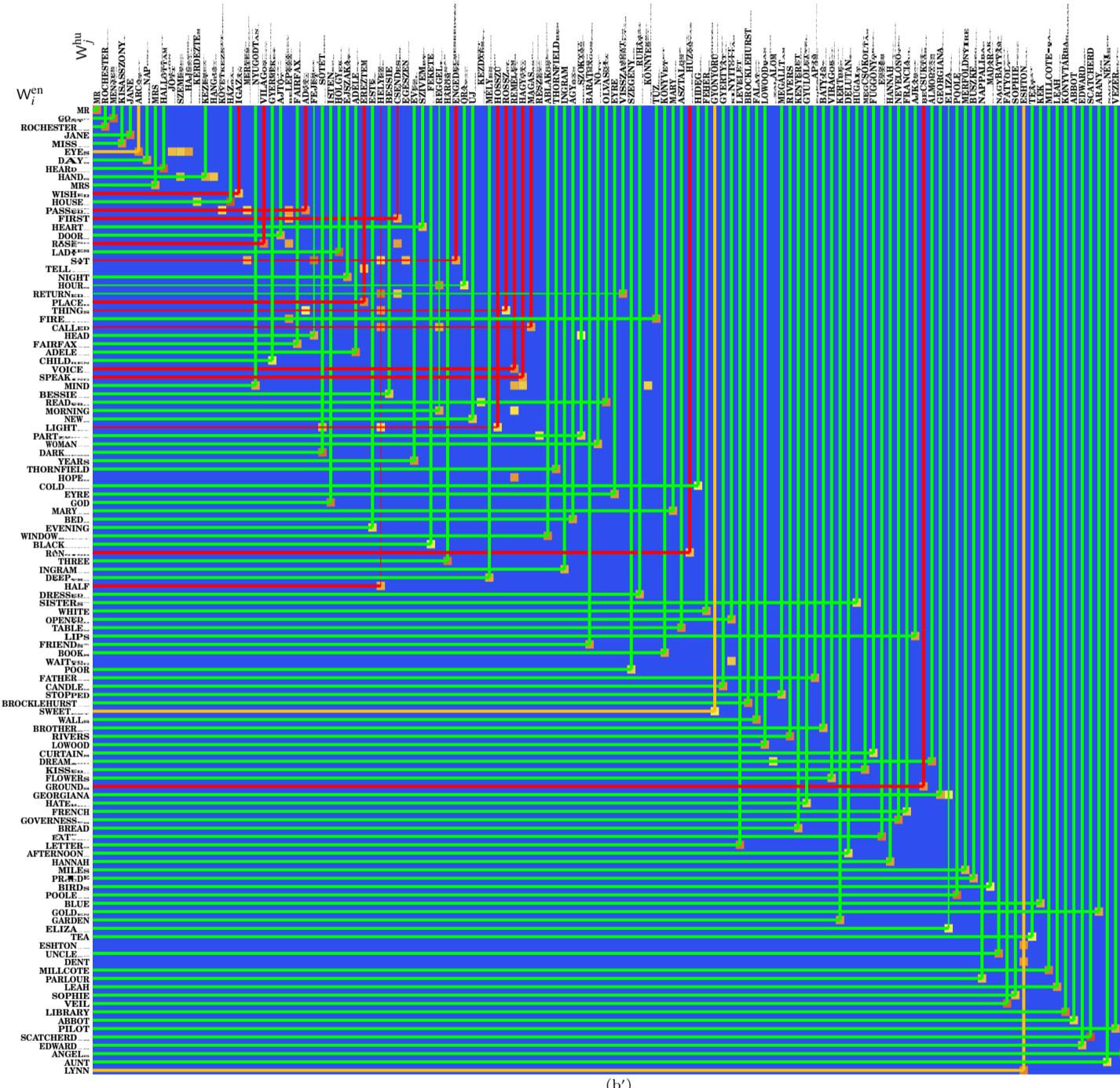
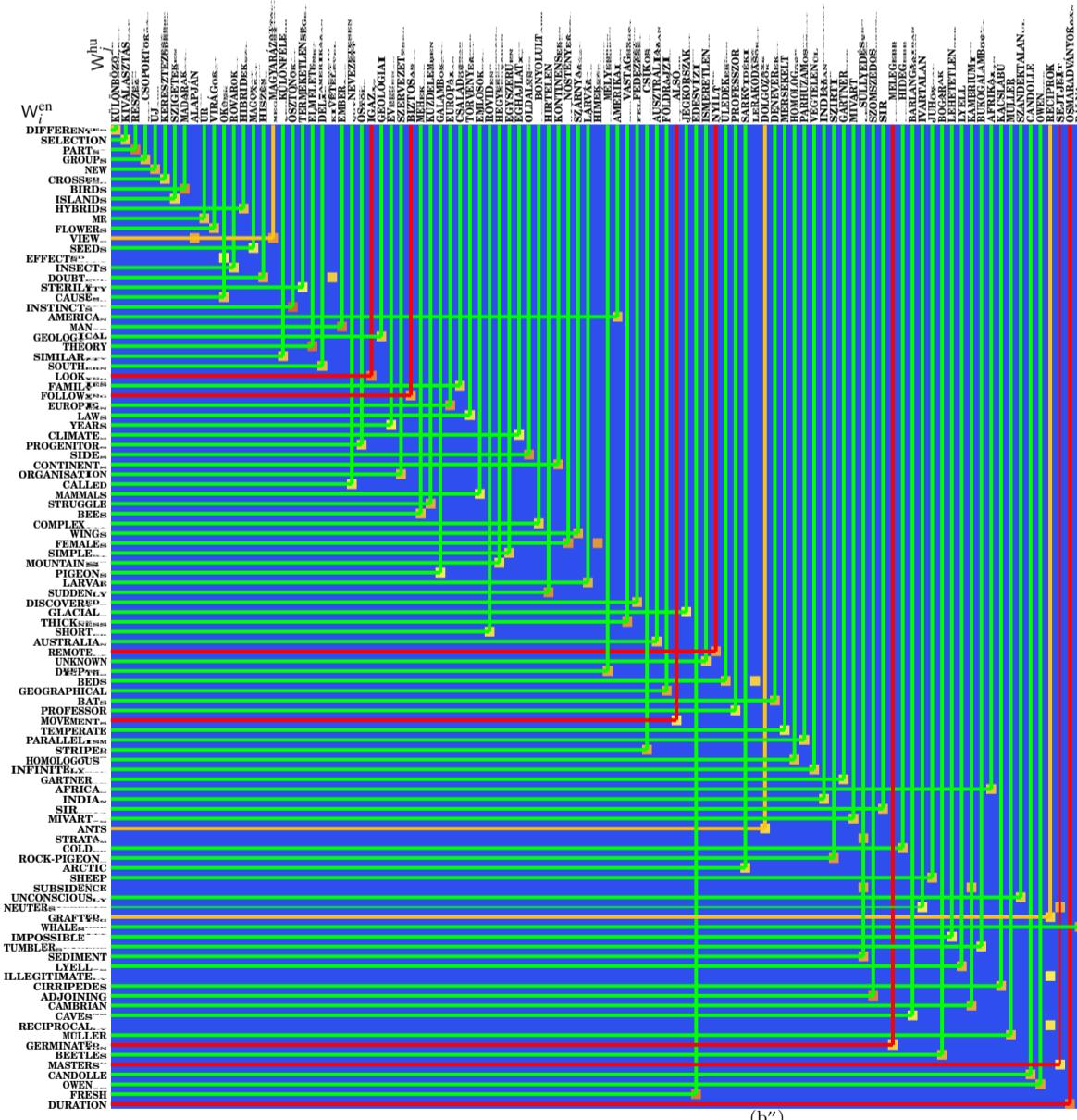
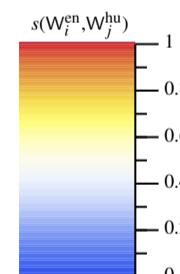


Fig. S13. Text mining in Hungarian. (Continued)  
 (a') Statistically identified topics ( $n_{ii} \geq 20$ ) in a Hungarian version of *Jane Eyre*, with the same color encoding scheme as Fig. S3. (b') Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{hu}})$  between selected topics in English and Hungarian versions of *Jane Eyre*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) rank of the similarity score for the optimal node. Green (resp. red) cross-hair indicates a correct (resp. incorrect) match. Amber cross-hair marks a link between distinct concepts that share the same hypernyms. (c') Results from control experiments with different choices of the  $\varepsilon$ -parameter in ballpark screening criteria (1.13).



(a'')



(b'')

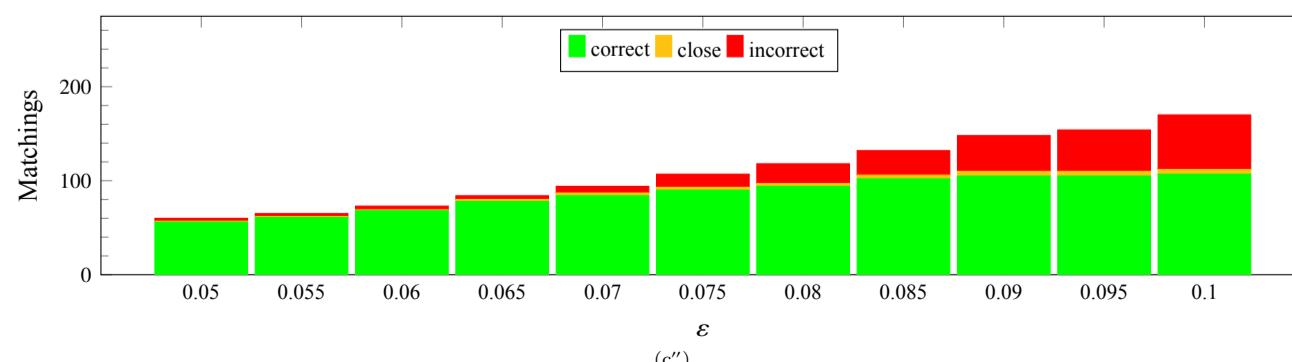


Fig. S13. Text mining in Hungarian. (Continued)  
(a'') Statistically identified topics ( $n_{ii} \geq 20$ ) in a Hungarian version of *Origin of Species*, with the same color encoding scheme as Fig. S3. (b'') Semantic similarities  $s(W_i^en, W_j^hu)$  between selected topics in English and Hungarian versions of *Origin of Species*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) rank of the similarity score for the optimal node. Green (resp. red) cross-hair indicates a correct (resp. incorrect) match. Amber cross-hair marks a link between distinct concepts that share the same hypernyms. (c'') Results from control experiments with different choices of the  $\epsilon$ -parameter in ballpark screening criteria (1.13).

## 9 Approximate word clustering in various agglutinative languages

Basque, Korean and Turkish are three agglutinative SOV languages with different linguistic ancestors,<sup>104</sup> but they do share a few typological and morphological features:

- All these languages are rich in suffixes, but have almost no prefixes. A typical word in these languages assumes the form of



- Typically, the word stem is one or two syllables long.<sup>105</sup> However, due to strong and persisted foreign influences on these languages (French/Spanish on Basque, Chinese/English/Japanese on Korean, Arabic/Persian on Turkish), there are certain exceptions to this estimate. Sometimes, stems of loanwords may be confused with native suffixes, if not handled properly.
- There are hardly any irregular verbs or nouns. Unlike Finnish, the suffixes in these languages (usually) do not cause sound changes upon interactions.
- These languages have more sophisticated systems for kinship terms than typical Indo-European languages. There could be different words for a sibling of a given gender, depending on seniority (Korean) or the gender of the possessor (Basque and Korean); there could be different words for extended family members depending on maternal/paternal sides (Korean and Turkish). Our word clustering algorithms will accommodate to such cultural differences, as much as possible.

---

<sup>104</sup>Some Altaicists believe that both Korean and Turkish descend from the same Proto-Altaic language [49].

<sup>105</sup>Since we do not know, *a priori*, whether the second syllable belongs to the stem or the suffix, a non-trivial stemming algorithm is still necessary.

## 9.1 Modified Porter stemming algorithm for Basque

Our Basque stemming algorithm is adapted from the method of Mikel Otxandorena (<http://snowball.tartarus.org/algorithms/basque/stemmer.html>). Our modest contributions include the following aspects:

- We add special treatments for the only irregular adjective in Basque: *on* “good” [50, §3.1.3.3].
  - We add special routines to cluster finite forms<sup>106</sup> of Basque verbs [50, §3.5.2.1].
  - We carefully handle loanwords that begin with *des~*, *err~*, *irr~* and so on.
  - We carefully treat the negative suffixes *~gabe* and *~gabea*, which carry similar functions as the English prefix *un~* and suffix *~less*.
  - We group Basque kinship terms according to (roughly) their English equivalents.

Basque stop words can carry various agglutinative suffixes. Instead of enumerating all of them in a single list, we will define several string patterns that match the criterion for stop words.

**Definition 9.1** (Basque stop words). Define the pattern **BasqueStopRoot** as any one from the following list:

*aitzin, albo, alde, arte, asko, atze, aurre, azken, azpi, barru, batega, batek, baten, batera, batet, behe, beza, buruzko, dela, denaren, denet, ditu, dizu, eduki, egote, erdi, ere, gain, garen, genieza, genitu, genu, geure, gibel, hain, heure, hieza, hitu, huen, inguru, izan, izate, lizki, neure, nieza, nitu, nuen, ondo, zela, zenieza, zenithu, zenu, zeuen, zeure, zieza, zire, zitu, zizki, zuen.*

Define the pattern **BasqueStopString** as any one from the following list:

<sup>106</sup>All but a handful verbs (like “come”, “go”, “say” and “walk”) in Basque can only exist in non-finite forms, which do not involve personal prefixes.

dakizueke, dakizuen, dakizuke, dakizun, dateke, daude, daudeke, daudela, Dekretu, delako, den, dena, denak, denean, denek, denok, dezagun, dira, dirateke, dire, direla, diren, direnak, dituk, ditun, dizkio, du, dua, duela, dugu, dugun, duk, dun, dut, dute, dutela, duzu, duzuela, duzun, ea, edo, egin, egingo, egiten, egon, egongo, egote, egotea, egoten, elkarrekin, eman, emana, ematen, erdian, eske, esker, eskuratzerakoan, eta, euren, ez, ezean, ezen, ezik, ezin, ezinbesteko, ezta, gabe, gagozkiake, gagozkian, gagozkie, gagozkieke, gagozkiene, gagoz-  
kik, gagozkin, gagozkinake, gagozkinan, gagozko, gagozkiore, gagozkion, gagozkizu, gagozkizue, gagozkizueke, gagozkizuen, gagozkizuke, gagozkizun, gainean, gainera, gaitez, gaitezke, gaituk, gaitun, gakizkiake, gakizkian, gakizkie, gakizkien, gakizkiore, gakizkion, gakizkizueke, gakizkizuen, gakizkizuke, gakizkizun, gara, garateke, gatzaizkiake, gatzaizkie, gatzaizkik, gatzaizkin, gatzaizkinake, gatzaizkio, gatzaizkioke, gatzaizkiz-  
zu, gatzaizkizue, gatzaizkizueke, gatzaizkizuke, gaude, gaudeke, gauden, gehiago, gehiegi, gehien, gengozkiake, gengozkian, gengozkie, gengozkieken, gengozkien, gengozkinake, gengozkinan, gengozkiore, gengozkioken, gengozkion, gengozkizue, gengozkizueken, gengozkizuen, gengozkizuke, gengozkizukeen, gengozkizun, geni-  
tuzke, genituzkeen, genizue, genkizkiake, genkizkian, genkizkie, genkizkieken, genkizkien, genkizkiore, gen-  
kizkioken, genkizkion, genkizkizue, genkizkizueken, genkizkizuen, genkizkizuke, genkizkizukeen, genkizkizun,  
genuke, genukeen, gero, gerok, geroztik, geu, geundeke, geundekeen, geunden, geure, ginateke, ginatekeen, gin-  
duan, gindunan, ginen, gintezen, ginteze, gintezeen, gintaizkiake, gintaizkian, gintaizkie, gintaizkie-  
keen, gintaizkien, gintaizkinake, gintaizkinan, gintaizkiore, gintaizkioken, gintaizkion, gintaizkizue, gintaiz-  
kizueken, gintaizkizuen, gintaizkizuke, gintaizkizukeen, gintaizkizun, gisa, guhaur, gurea, gurekin, gutxi,  
gutxiago, gutxienez, guztiak, guztion, guztiotzat, hadin, hago, hagoke, hagokidake, hagokie, hagokieke, hago-  
kigu, hagokiguke, hagokio, hagokioke, hagokit, hainbeste, haiteke, haiz, haizateke, hakidake, hakidan, hakieke,  
hakien, hakiguke, hakigun, hakioke, hakion, haraino, harantz, haratago, hargatik, hartara, hatzaidake, hatzaie,  
hatzaieke, hatzaigu, hatzaiguke, hatzaio, hatzaioke, hatzait, haukin, hauen, hedin, hengoer, hengoerke, hengo-  
keen, hengokidake, hengokidakeen, hengokidan, hengokieke, hengokieken, hengokien, hengokiguke, hengoki-  
gukeen, hengokigun, hengokioke, hengokiokeen, hengokion, henkidake, henkidakeen, henkidan, henkieke, hen-  
kiekeen, henkien, henkiguke, henkigukeen, henkigun, henkioke, henkiokeen, henkion, herori, heu, heure, hihiar,  
hinteke, hintekeen, hintzaidake, hintzaidakeen, hintzaidan, hintzaieke, hintzaiekeen, hintzaiken, hintzaiguke, hin-  
tzaigukeen, hintzaigun, hintzaioke, hintzaiokeen, hintzaion, hintzateke, hintzatekeen, hintzen, hitzke, hitzkeen,  
hola, honaino, honantz, honen, horraino, horrantz, horrek, horren, hortxe, huke, hukeen, hurrengu, ia, inguru, in-  
guruau, inolako, inoren, izan, izanen, izango, izate, izatea, izaten, kanpo, kanpoan, kontra, landa, laster, legoke,  
legokiate, legokidake, legokieke, legokiguke, legokinake, legokioke, legokizue, legokizuke, legozkiake, legoz-  
kidake, legozkieke, legozkiguke, legozkinake, legozkiore, legozkizue, legozkizuke, lehenago, lekiake, lekidake,  
lekiike, lekiguke, lekiroke, lekizkiake, lekizkidake, lekizkie, lekizkiguke, lekizkiore, lekizkizue, lekizkizuke, le-  
kizueke, lekizuke, leudeke, lirateke, liteke, litzke, litzkeen, litzaiae, litzaidake, litzaiaeke, litzaiguke, lit-  
zainake, litzaioke, litzaizkiake, litzaizkidake, litzaizkie, litzaizkiguke, litzaizkinake, litzaizkiore, litzaizkizue-  
ke, litzaizkizuke, litzaizue, litzaizuke, litzateke, lortu, lortzen, luke, lukeen, maiatza, maiz, mundu, nadin, nago,  
nagoen, nagoke, nagokiate, nagokian, nagokie, nagokieke, nagokien, nagokik, nagokinake, nagokinan,  
nagokio, nagokioke, nagokion, nagokizu, nagokizue, nagokizuen, nagokizuke, nagokizun, nahiko,  
nahikoa, nahiz, naiteke, naiz, naizateke, nake, nakiake, nakian, nakiike, nakien, nakiuke, nakion, nakizueke, na-  
kizuen, nakizuke, nakizun, nan, natzaiake, natzaie, natzaieke, natzaik, natzain, natzainake, natzaio, natzaike,  
natzaizu, natzaizue, natzaizuke, nau, nauk, naun, nekatu, nekez, nendin, nengoer, nengoerke, nen-  
gooken, nengokiate, nengokian, nengokieke, nengokieken, nengokien, nengokinake, nengokinan, nengokioke,  
nengokiokeen, nengokion, nengokizue, nengokizueken, nengokizuen, nengokizuke, nengokizukeen, nengoki-  
zun, nenkiake, nenkian, nenkieke, nenkien, nenkiore, nenkioken, nenkion, nenkizue, nenkizueken,  
nenkizuen, nenkizuke, nenkizukeen, nenkizun, nerau, neu, neure, nihaur, nik, ninduan, nindunan, nintekte, ninte-  
keen, nintzaiake, nintzaian, nintzaie, nintzaieke, nintzaiken, nintzain, nintzainake, nintzainan, nintzaioke, nintzaiokeen,  
nintzaion, nintzaizue, nintzaizueken, nintzaizuen, nintzaizuke, nintzaizukeen, nintzaizun, nintzateke, nintzate-  
keen, nintzen, nitzke, nitzkeen, norakoak, norbaiten, norberaren, norberaren, nori, nuke, nukeen, omen, omeone,  
ondo, ondoan, ondoren, ondoriox, oraindiak, ordea, oso, osoa, ostera, ote, sartu, segidan, sekula, truk, truke, utzi,  
utziz, zagozkidake, zagozkidakete, zagozkidate, zagozkie, zagozkieke, zagozkiekete, zagozkiete, zagozkigu, zagoz-  
kiguke, zagozkigukete, zagozkigute, zagozko, zagozkiore, zagozkiot, zagozkit, zaiake, zaidake,  
zaie, zaieke, zaigu, zaiguke, zaik, zain, zainake, zaio, zaike, zait, zaitez, zaitezke, zaitezke, zaizkiate,  
zaizkidake, zaizkie, zaizkieke, zaizkigu, zaizkiguke, zaizkik, zaizkin, zaizkinake, zaizkio, zaizkioke, zaizkit, zaiz-  
kizu, zaizkizue, zaizkizueke, zaizkizuke, zaizu, zaizue, zaizueke, zaizuke, zakizkidake, zakizkidake, zakizkidan,  
zakizkidaten, zakizkie, zakizkiekete, zakizkien, zakizkieten, zakizkiguke, zakizkigukete, zakizkigun, zakizkiguten,  
zakizkioke, zakizkiokete, zakizkion, zakizkioten, zara, zarateke, zarete, zaretekete, zatekeen, zatia, zatzaizkide,  
zatzaizkidakete, zatzaizkide, zatzaizkie, zatzaizkieke, zatzaizkiekete, zatzaizkiete, zatzaizkigu, zatzaizkiguke, za-  
tzazkigukete, zatzaizkigute, zatzaizko, zatzaizkioke, zatzaizkiokete, zatzaizkioke, zatzaizkit, zaude, zaudeke, zau-  
dekete, zaudete, zazu, ze, zedin, zegooken, zegokian, zegokidakeen, zegokien, ze-

gokigukeen, zegokigun, zegokinan, zegokiokeen, zegokion, zegokizuekeen, zegokizuen, zegokizukeen, zegokizun, zegozkian, zegozkidakeen, zegozkidan, zegozkiekeen, zegozkien, zegozkigueen, zegozkigun, zegozkinan, zegozkiokeen, zegozkion, zegozkizuekeen, zegozkizuen, zegozkizukeen, zegozkizun, zehar, zekian, zekidakeen, zekidan, zekiekeen, zekien, zekigukeen, zekigun, zekiokeen, zekion, zekizkian, zekizkidakeen, zekizkidan, zekizkiekeen, zekizkien, zekizkigueen, zekizkigun, zekizkirokeen, zekizkion, zekizkizuekeen, zekizkizuen, zekizkizukeen, zekizkizun, zekizuekeen, zekizuen, zekizukeen, zekizun, zen, zengozkidake, zengozkidakeen, zengozkidakete, zengozkidaken, zengozkidan, zengozkidaten, zengozkieke, zengozkiekeen, zengozkiekete, zengozkieketen, zengozkien, zengozkieten, zengozkigupe, zengozkigueen, zengozkigukete, zengozkiguketen, zengozkigun, zengozkiguten, zengozkioke, zengozkiokeen, zengozkiokete, zengozkioketen, zengozkion, zengozkioten, zenidate, zenigate, zenituze, zenituzkeen, zenituzkete, zenituzketen, zenkizkidake, zenkizkidakete, zenkizkidaketen, zenkizkidan, zenkizkidaten, zenkizkieke, zenkizkiekeen, zenkizkiekete, zenkizkieketen, zenkizkien, zenkizkieten, zenkizkiguke, zenkizkigueen, zenkizkigukete, zenkizkiguketen, zenkizkigun, zenkizkiguten, zenkizkioke, zenkizkiokeen, zenkizkiokete, zenkizkioketen, zenkizkion, zenkizkioten, zenuke, zenukeen, zenukete, zenuketen, zeren, zergatik, zerok, zerori, zertaz, zeu, zeudekeen, zeuden, zeuak, zeundeke, zeundekeen, zeundeketen, zeunden, zeundeten, zeure, zinateke, zinatekeen, zinatekete, zinateketen, zinen, zineten, zintezene, zintezke, zintezkeen, zintezketen, zintezten, zintzaizkidake, zintzaizkidakeen, zintzaizkidakete, zintzaizkidaketen, zintzaizkidan, zintzaizkidaten, zintzaizkieke, zintzaizkiekeen, zintzaizkiekete, zintzaizkieketen, zintzaizkien, zintzaizkieten, zintzaizkiguke, zintzaizkigukeen, zintzaizkigukete, zintzaizkiguketen, zintzaizkigun, zintzaizkiguten, zintzaizkioke, zintzaizkiokeen, zintzaizkiokete, zintzaizkioketen, zintzaizkion, zintzaizkieten, zitauan, zituan, zituanan, zitzaian, zitzaidakeen, zitzaidan, zitzaikeen, zitzaien, zitzaiukeen, zitzaign, zitzainan, zitzaiokeen, zitzaison, zitzazkian, zitzazkidakeen, zitzazkidan, zitzazkiekeen, zitzazkien, zitzazkigueen, zitzazkigun, zitzazkinan, zitzazkiokeen, zitzazkion, zitzazkizuekeen, zitzazkizuen, zitzazkizukeen, zitzazkizun, zitzazuekeen, zitzazuen, zitzazukeen, zitzazun, zuan, zuhaur, zuihauk, zunan, zurea, zutaz.

Define the pattern **PronBasque** as any one from the following list:

bera, beraiek, beraien, berak, beraren, berau, berauek, bere, beren, berori, beroriek, eurak, ezer, genian, genien, geninan, genion, genizuen, genizun, gu, gure, haiiek, haien, hala, han, handik, hango, hara, haren, hau, hauek, hemen, hemendik, hemengo, hi, hidan, hien, higun, hion, hire, hona, honela, hor, hori, horiek, horko, horra, horrela, hortik, hura, inoiz, inola, inon, inor, ni, nian, nien, ninan, nion, nire, nizuen, nizun, noiz, nozbait, nola, nolabait, non, nonbait, nondik, nongo, nor, nora, norbait, orain, orduan, zein, zenidan, zenidaten, zenien, zenieten, zenigun, zeniguten, zenion, zenioten, zer, zerbait, zian, ziaten, zidan, zidaten, zien, zieten, zigun, ziguten, zinaten, zion, zioten, zizuen, zizueten, zizun, zizuten, zu, zuek, zuen, zure.

Define **tenerBasque** as the following pattern:

( $\emptyset$ |ba|e)( $\emptyset$ |hen|le|nen)( $\emptyset$ |be|da|ga|ge|ha|na|za|ze)( $\emptyset$ |n)ukaX.

If a word exactly matches any one of the following:

(BasqueStopRoot)X,  
BasqueStopString,  
(Ø|edo)PronBasque(Ø|bait|nahi|txe|xe)X,  
tenerBasque,

then we consider it a Basque stop word. All the Basque stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list.  $\square$

### **9.1.1 Effective spelling and essential root**

**Definition 9.2** (Basque protected range). Set  $\mathbf{V} = (a|e|i|o|u)$ ,  $\mathbf{C} = \overline{\mathbf{V}}$ . Search for the pattern  $\mathbf{V}_m \mathbf{C}_m \mathbf{V}_{\sim}$  or

$(\emptyset|ba|bi|ko)(\emptyset|hen|hin|le|li|nen)(\emptyset|be|da|ga|ge|gi|ha|na|ni|za|ze|zi)\mathbf{C}_m\mathbf{V}_m\mathbf{C}\sim$

in the string  $\hat{\sigma}$ . If such a pattern cannot be found, or  $\hat{\sigma}$  does not contain  $\mathbf{V}$ , then define  $\rho_* = 0$ ; otherwise, define  $\rho_*$  as the last position occupied by the aforementioned pattern. If  $\rho_* > 0$ ,  $\ell(\hat{\sigma}) > \rho_*$ ,  $\hat{\sigma}^{[\rho_*]} = r$  (i.e. the  $\rho_*$ -th character is the letter  $r$ ) and  $\hat{\sigma}^{[\rho_*+1]} = \mathbf{C}$  (i.e. the  $(\rho_* + 1)$ -st character is not a vowel), then  $\text{ProtRg}(\hat{\sigma}) = 1 + \rho_*$ ; otherwise, set  $\text{ProtRg}(\hat{\sigma}) = \rho_*$ .

**Algorithm 9.3** (Basque effective spelling). Set  $\mathbf{V} = (a|e|i|o|u)$ ,  $\mathbf{C} = \overline{\mathbf{V}}$ . For a Basque word  $\hat{\sigma}$ , its effective spelling  $\text{EffSpell}(\hat{\sigma})$  is constructed in the following steps:

- (1) Convert to lowercase and do  $\hat{g}_{m, \text{abe}}(\emptyset | a) \sim \rightarrow xx\hat{g}_m$ .

(2) Do  $\text{tx} \rightarrow \check{c}$ .

(3) Do  $(izeb|izeko)\sim \rightarrow \alpha\text{unt}$ ,  $\text{Cy} \rightarrow \text{Ciy}$ ,  $\text{anai}\sim \rightarrow \text{neb}$ ,  $\text{arreb}\sim \rightarrow \text{ahizp}$ ,  $\text{lehengus}(in|u)\text{X}\sim \rightarrow \text{cow}\sigma i\check{v}$ .

(4) Replace

$\text{aban}\sim$	$\text{abe}\sim$	$\text{agir}\sim$	$\text{aldarri}\sim$	$\text{andere}\check{n}\sim$	$\text{ater}\sim$	$\text{bakar}\sim$	$\text{bar}\sim$	$\text{begira}\sim$	$\text{begiru}\sim$	$\text{begitart}\sim$	
$\beta a$	$ab$	$ayir$	$aldrri$	$\mu i\sigma$	$atr$	$\sigma i\check{v}yr$	$\lambda a\varphi$	$\lambda uka$	$\rho e\sigma \pi tu$	$\check{f}a\sigma t$	
$\text{bil}_{\hat{\chi}}(a)\sim$	$\text{bila}\sim$	$\text{burg}\sim$	$\text{buru}\sim$	$\text{carolin}\sim$	$\text{charlotte}\sim$	$\text{dantz}\sim$	$\text{denbor}\sim$	$\text{des}_{\hat{\chi}}(i)\sim$	$\text{desio}\sim$		
$\beta i\lambda\hat{\chi}$	$\beta i\alpha\lambda$	$bupg$	$bru$	$\kappa raln$	$\kappa ralotte$	$\delta ac$	$\tau me\pi$	$\delta ds\hat{\chi}$	$dsir$		
$\text{desi}\sim$	$\text{desk}\sim$	$\text{egona}\sim$	$\text{elizaX}\sim$	$\text{era}_{\hat{\chi}}(m i)\sim$	$\text{err}\sim$	$\text{familia}\sim$	$\text{fitz}\sim$	$\text{gela}\sim$	$\text{gorrot}\sim$	$\text{hertf}\sim$	
$dsi$	$dsk$	$\sigma \tau aya$	$\varepsilon el\zeta$	$exq\hat{\chi}$	$\varepsilon rr$	$\varphi ma\lambda i$	$\varphi \tau \zeta$	$gela$	$\eta at$	$\eta r\tau e\varphi$	
$\text{hobeX}\sim$	$\text{hunsf}\sim$	$\text{ibil}(ald er)\text{X}\sim$		$\text{ideia}\sim$	$\text{igaro}\sim$	$\text{ikar}\sim$	$\text{indi}\sim$	$\text{irr}\sim$	$\text{izen}\sim$	$\text{janeX}\sim$	
$\beta e\sigma\tau$	$hnsf$	$ibili$		$\iota\delta ea$	$\sigma pendo$	$\ckar$	$ndi$	$urr$	$\tilde{v}amu$	$ja\check{v}$	$cant$
$\text{karta}\sim$	$\text{kol}\sim$	$\text{liburuteg}\sim$	$\text{liydia}\sim$	$\text{londres}\sim$	$\text{lots}\sim$	$\text{maria}\sim$	$\text{mariy}\sim$	$\text{meriy}\sim$	$\text{mingarri}\sim$	$\text{musik}\sim$	
$\kappa rta$	$kl$	$\lambda\beta ip$	$\lambda\eta\delta ia$	$lndres$	$\sigma ha\mu$	$\mu ari\alpha$	$\mu ari\eta$	$\mu ry$	$\pi ai\tilde{v}i$	$\mu i\sigma k$	
$\text{negozio}\sim$	$\text{negu}\sim$	$rst$	$\text{senarra}\sim$	$\text{senti}\sim$	$\text{tentel}\sim$	$\text{ugar}\sim$	$\text{zakar}\sim$	$\text{amak}$	$de$	$long$	
$ngoz$	$\tilde{v}egw$	$\rho\sigma\tau$	$hu\sigma\beta a$	$\varphi ee\lambda i$	$\varphi \omega\lambda l$	$ugr$	$\rho u\delta r$	$ama$	$\delta ee$	$\tilde{v}noy$	
										$\text{on}(\emptyset a ak ik)$	
										$\beta e\sigma\tau$	

(5) Replace

$(\emptyset ba bi zi)(\emptyset li)(d g h j n z)(\emptyset in)(\emptyset d)oa(\emptyset k l n t z)\text{X}\sim$									
$\beta go\beta$									
$(\emptyset ba bi zi)(\emptyset li)(d g h j n z)(\emptyset in)(\emptyset d)(\emptyset a)e\text{ki}(\emptyset k l n t z)\text{X}\sim$									
$\beta know\beta$									
$(\emptyset ba e)(\emptyset hen le nen)(\emptyset be da ga ge ha na za ze)(\emptyset n)ka(r z)\text{X}\sim$									
$ekarr$									
$(\emptyset ba e)(\emptyset hen le nen)(\emptyset be da gal ge ha na za ze)(\emptyset n)rama\text{X}\sim$									
$\beta \tau a\kappa\beta$									
$ahazt\sim$	$\text{begiX}\sim$	$\text{ego}\sim$	$\text{eliz}_{\hat{\chi}}(a)\sim$	$\text{elizab}\sim$	$\text{harro}\sim$	$\text{kont}(a u)\text{X}\sim$	$\text{bru}(\emptyset bid\text{X} ko)$		
$sop\gamma t$	$eye$	$eg$	$\varepsilon li\zeta\hat{\chi}$	$\varepsilon lizab\beta$	$\pi i\rho\delta o$	$\kappa o\tilde{v}\tau$	$bruak$		

(6) Call the result so far as  $\hat{\sigma}^\dagger$ , and break it down into  $\hat{\sigma}^\dagger = \hat{\sigma}_1^\dagger \hat{\sigma}_2^\dagger$ , where  $\ell(\hat{\sigma}_1^\dagger) = \text{ProtRg}(\hat{\sigma}_1^\dagger)$ . Work on  $\hat{\sigma}_2^\dagger$  as follows:

(6.1) Do  $\check{c} \rightarrow tx$ .

(6.2) Delete any pattern in the following list:

ago, aldatu, aldi, aldia, an, arazi, ari, atu, atze, bera, bide, bidea, dako, du, dura, ean, era, erreza, erreza, eta, etan, etari, ez, eza, ezin, ezina, gai, gaia, gailu, gailua, gaitza, gale, galea, go, gune, gunea, gura, ide, idea, ka, kaitza, kaitza, kan, kari, karia, karria, kera, keta, ki, kide, kidea, kin, kina, kizun, kizuna, kor, korra, kuna, kunde, kundea, kune, kunea, kuntza, kura, la, lari, le, men, mena, or, orra, pen, pena, pera, pide, pidea, rean, rekin, taile, tailea, taldi, taldea, tarazi, tari, tatu, tezin, tezina, tio, tu, tun, tuna, tura, tzaga, tzaile, tzailea, tzaka, tzake, tzat, tze, tzeke, tzez

and the letters (if any) thereafter.

(6.3) Delete any pattern in the following list:

ada, ail, aizun, ak, alde, aldea, aldi, aldia, anda, anga, antza, ar, ara, ari, aria, aro, aroa, arte, artea, asi, asia, asun, asuna, aurre, aurrea, behar, bera, bizia, burua, dar, dara, degi, degia, denda, di, du, dua, dun, duna, duri, duria, duru, durua, egi, egia, ek, eko, eme, emea, ena, enea, eria, ero, eroa, eroz, eroza, estu, estua, eta, etako, etan, etara, etxe, etxea, ez, eza, ezia, ga, gabe, gabea, gai, gaia, garna, garren, garrena, ge, gei, geia, gela, gerren, gerrena, gibel, gibela, gile, gilea, gintza, gintzo, gintzu, giro, go, goi, gune, gunea, handi, handia, i, ka, kabe, kabea, kada, kail, kaila, kalde, kaldea, kan, kana, kari, karia, kera, keria, ket, keta, ki, kia, kide, kin, kina, kintza, kirri, kirria, ko, koa, koi, koia, koitz, koitza, kondo, kondo, kor, korra, kote, kotea, kume, kumea, kuntza, lari, laria, larria, leku, lekua, liar, liara, mendi, mendia, mendua, mendua, mentua, min, mina, n, na, nahi, ne, nea, ngo, ngoa, no, noa, o, oa, ohi, ohia, oi, oia, ola, ondo, ondoa, ontzi, ontzia, orde, ordea, ordua, oro, oroa, os, osa, oso, osoa,

*oste, ostea, pe, pea, pera, ra, ro, sa, ska, skila, sko, sta, ta, tako, takoa, talde, taldea, taldi, taldia, tan, tar, tara, tari, taria, tarik, tariko, taro, taroa, tasun, tasuna, te, tea, tegi, tegia, teria, ti, tia, tiar, tiara, tila, to, toa, toki, tokia, tra, tsu, tsua, tto, ttoa, tu, tua, tuko, txo, txoa, txu, txua, tz, tzain, tzaina, tzale, tzalea, tzar, tzara, tzarra, tzo, tzoa, tzu, tzua, una, une, unea, urren, urrena, xka, z, za, zain, zaina, zale, zalea, zaro, zaroa, zino, zinoa, zio, zioa, zione, ziona, zonea, zka, zko, zkoa, zp, zto, ztoa, zu, zua*

*and the letters (if any) thereafter.*

(6.4) Delete any pattern in the following list:

*dade, date, era, ero, gi, go, ik, keria, ki, la, lanik, larik, rik, ro, tade, tate, to, ztik*

*and the letters (if any) thereafter.*

(6.5) Call the result so far as  $\hat{\sigma}_2^{\ddagger}$ . Concatenate  $\hat{\sigma}_1^{\ddagger}$  and  $\hat{\sigma}_2^{\ddagger}$ .

(7) *Replace*

$(\emptyset ba e)(\emptyset hen le nen)(\emptyset be da ga ge ha na za ze)(\emptyset n)to(r z)\mathbf{X} \sim$	$etorr$	
$(\emptyset ba i)(\emptyset hen le nen)(\emptyset be da ga ge ha na za ze)(\emptyset n)bil\mathbf{X} \sim$	$(b d gen h n z zen)io(d e g k n s t z)\mathbf{X} \sim$	$one(\emptyset g)$
$\beta walk\beta$	$esa$	$\beta es\tau$

**Algorithm 9.4** (Basque essential root). Let  $\hat{o}$  be the effective spelling of a Basque word, then its corresponding essential root  $\text{EssRoot}(\hat{o})$  is constructed as  $\hat{o}[\text{ProtRg}(\hat{o})]$ .

### 9.1.2 Approximate clustering

There will not be any mutation rules in our algorithm for Basque, so SimpHrdTest = HrdTest.

**Algorithm 9.5** (Simple heredity test). Suppose that  $\hat{\alpha}$  and  $\hat{\beta}$  are both lowercase strings. The Boolean-valued function  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if  $\hat{\alpha}$  contains at least one instance of **V**, **AND** at least one of the following two conditions holds:<sup>107</sup>

- (i)  $\hat{\alpha} = \hat{\beta}$ ;  
(ii)  $\hat{\alpha}' = \hat{\beta}'$ , where the strings with prime result from doing  $\sim^{\mathbf{X}\epsilon}(\mathbf{V}\hat{\chi})(\emptyset|b)(a|e) \rightarrow \mathbf{X}$  on their counterparts without prime.

**Algorithm 9.6** (Approximate clustering of Basque words). *The approximate clustering of a list of Basque words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  is completed in two stages:*

- (1) We sort the list  $\{(\hat{\alpha}_1, \text{EffSpell}(\hat{\alpha}_1)), \dots, (\hat{\alpha}_N, \text{EffSpell}(\hat{\alpha}_N))\}$  alphabetically according to the last component. If two consecutive neighbors in the alphabetized list  $\{(\hat{\alpha}_{(1)}, \text{EffSpell}(\hat{\alpha}_{(1)})), \dots, (\hat{\alpha}_{(N)}, \text{EffSpell}(\hat{\alpha}_{(N)}))\}$  satisfy  $\text{HrdTest}(\text{EffSpell}(\hat{\alpha}_{(n)}), \text{EffSpell}(\hat{\alpha}_{(n+1)})) = \text{FALSE}$ , where  $n \in \mathbb{Z} \cap [1, N]$ , then add a demarcation line between these two entries. In this way, the alphabetized list is divided into separate groups of words tagged with their effective spellings:  $\{g_1 = \{(\hat{\alpha}_{(1,1)}, \text{EffSpell}(\hat{\alpha}_{(1,1)})), \dots, (\hat{\alpha}_{(1,n_1)}, *)\}, \dots, g_M = \{(\hat{\alpha}_{(M,1)}, *), \dots, (\hat{\alpha}_{(M,n_M)}, *)\}\}$ , where each sublist also preserves alphabetic order. In particular, we have  $\hat{\alpha}_{(1,1)} = \hat{\alpha}_{(1)}$  and  $\hat{\alpha}_{(M,n_M)} = \hat{\alpha}_{(N)}$ .
  - (2) For each group of words  $g_m = \{(\hat{\alpha}_{(m,1)}, *), \dots, (\hat{\alpha}_{(m,n_m)}, *)\}$  where  $m \in \mathbb{Z} \cap [1, M]$ , we augment it into a tagged entry  $G_m = (g_m, \text{EffSpell}(\hat{\alpha}_{(m,1)}), \text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})))$ . The list  $\{G_1, \dots, G_M\}$  is sorted alphabetically, with respect to  $\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))$  (with higher priority), and  $\text{EffSpell}(\hat{\alpha}_{(m,1)})$  (with lower priority). If two consecutive neighbors in the alphabetized list  $\{G_{(1)} = (g_{(1)}, \hat{\gamma}'_{(1)}, \hat{\gamma}''_{(1)}), \dots, G_{(M)} = (g_{(M)}, \hat{\gamma}'_{(M)}, \hat{\gamma}''_{(M)})\}$  satisfy

**HrdTest**( $\hat{\gamma}_{(m)}''$ ,  $\hat{\gamma}_{(m+1)}''$ ) = FALSE

AND

**HrdTest( $\hat{\gamma}'_{(m)}, \hat{\gamma}'_{(m+1)}$ ) = FALSE**

where  $m \in \mathbb{Z} \cap [1, M]$ , then add a demarcation line between these two entries. In this manner, the alphabetized list  $\{G_{(1)}, \dots, G_{(M)}\}$  is divided into separate groups of tagged clusters  $\{\Gamma_1 = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, \Gamma_K = \{G_{(K,1)}, \dots, G_{(K,m_K)}\}\}$ . Finally, the output list of word clusters  $\{\tilde{\Gamma}_1, \dots, \tilde{\Gamma}_K\}$  is constructed by discarding all the tags (effective spellings, essential roots) from each  $\Gamma_k$ , where  $k \in \mathbb{Z} \cap [1, K]$ .

*Example 9.6.1.* Our clustering algorithm correctly groups the following families of Basque words:

<sup>107</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical OR.

*ahizpa, ahizpak, ahizparekin, ahizparen, ahizparengan, ahizparengana, ahizparenganaino, ahizparenganako, ahizparenganantz, ahizparengandik, ahizparengatik, ahizparentzat, ahizpari, ahizparik, ahizpatzat, ahizpaz, ahizpei, ahizpek, ahizpekin, ahizpen, ahizpengan, ahizpengana, ahizpenganaino, ahizpenganako, ahizpenganantz, ahizpengandik, ahizpengatik, ahizpentzat, ahizpez, arreba, arrebak, arrebarekin, arrebaren, arrebarengan, arrebarengana, arrebarenganaino, arrebarenganako, arrebarenganantz, arrebarengandik, arrebarengatik, arrebarentzat, arrebari, arrebarik, arrebatzat, arrebaz, arrebei, arrebek, arrebekin, arreben, arrebengan, arrebengana, arrebenganaino, arrebenganako, arrebenganantz, arrebengandik, arrebengatik, arrebentzat, arrebez — “sister”; anaia, anaiaik, anaiarekin, anaiaaren, anaiaarenan, anaiaarenan, anaiaarenanaino, anaiaarenanako, anaiaarenanantz, anaiaengandik, anaiaengatik, anaiarentzat, anaiali, anaiarik, anaiaztat, anaiaz, anaiei, anaiek, anaiekin, anaien, anaient, anaientan, anaientanaino, anaientanako, anaientanantz, anaientangandik, anaientatik, anaientatzat, anaiez, neba, nebak, nebarekin, nebarenen, nebarengan, nebarengana, nebarenganaino, nebarenganako, nebarenganantz, nebarengandik, nebarengatik, nebarentzat, nebari, nebarik, nebatzat, nebaz, nebei, nebek, nebekin, neben, nebengan, nebengana, nebenganaino, nebenganako, nebenganantz, nebengandik, nebengatik, nebentzat, nebez — “brother”;*

*izeba, izeba, izebak, izebak, izebarekin, izebarekin, izebaren, izebarenen, izebarengan, izebarengan, izebarengana, izebarengana, izebarenganaino, izebarenganaino, izebarenganako, izebarenganako, izebarenganantz, izebarenganantz, izebarengandik, izebarengandik, izebarengatik, izebarengatik, izebarentzat, izebarentzat, izebari, izebari, izebarik, izebatzat, izebaz, izebe, izebek, izebekin, izeben, izebengan, izebengan, izebenganaino, izebenganako, izebenganantz, izebengandik, izebengatik, izebentzat, izebez, izeko, izekoak, izekoak, izekoarekin, izekoaren, izekoarengan, izekoarengana, izekoarenganaino, izekoarenganako, izekoarenganantz, izekoarengandik, izekoarengatik, izekoarentzat, izekoari, izekoaz, izekoei, izekoek, izekoekin, izekoen, izekoengan, izekoengana, izekoenganaino, izekoenganako, izekoenganantz, izekoengandik, izekoengatik, izekoenatzat, izekoez, izekok, izekorekin, izekoren, izekorengan, izekorengana, izekorenganaino, izekorenganako, izekorenganantz, izekorengandik, izekorengatik, izekorentzat, izekori, izekorik, izekotzat, izekoz — “aunt”.*

It also correctly cluster each of the following verbs:

*bagenbiltza, bagenbilzkie, bagenbilzkik, bagenbilzkin, bagenbilzkio, bagenbilzkizu, bagenbilzkizue, bahenbilkie, bahenbilkigu, bahenbilkio, bahenbilkit, balebilkie, balebilkigu, balebilkik, balebilkin, balebilkio, balebilkit, balebilkizu, balebilkizue, balebiltza, balebilzkie, balebilzkigu, balebilzkik, balebilzkin, balebilzkio, balebilzkit, balebilzkizu, balebilzkizue, banenbilkie, banenbilzik, banenbilkin, banenbilko, banenbilkizu, banenbilzkizue, bazenbiltza, bazenbiltzate, bazenbilzkidate, bazenbilzkie, bazenbilzkiete, bazenbilzkigu, bazenbilzkigute, bazenbilzkio, bazenbilzkiole, bazenbilzkit, bebiltza, bebilzkie, bebilzkigu, bebilzkik, bebilzkin, bebilzkio, bebilzkit, bebilzkizu, bebilzkizue, benbilkie, benbilkigu, benbilkik, benbilkin, benbilkio, benbilkit, benbilkizu, benbilkizue, dabilela, dabilke, dabilkiake, dabilkidake, dabilkie, dabilkie, dabilkigu, dabilkigute, dabilkik, dabilkin, dabilkinake, dabilkio, dabilkioke, dabilkit, dabilkizu, dabilkizue, dabilkizuke, dabilkizuke, gabiltza, gabiltzan, gabilzke, gabilzkiake, gabilzkian, gabilzkie, gabilzkie, gabilzkien, gabilzkik, gabilzkin, gabilzkinake, gabilzkinan, gabilzkie, gabilzkion, gabilzkizu, gabilzkizue, gabilzkizuke, gabilzkizuen, gabilzkizuke, gabilzkizun, genbiltzan, genbilzke, genbilzkeen, genbilzkiake, genbilzkian, genbilzkie, genbilzkieken, genbilzkien, genbilzkinake, genbilzkinan, genbilzkioke, genbilzkioken, genbilzkion, genbilzkizue, genbilzkizuen, genbilzkiuke, genbilzkizukeen, genbilzkizun, habilke, habilkidake, habilkie, habilkie, habilkigu, habilkiguke, habilkio, habilkioke, habilkit, henbilen, henbilke, henbilkeen, henbilkidake, henbilkidakeen, henbilkidan, henbilkie, henbilkieken, henbilkien, henbilkigu, henbilkiguke, henbilkigukeen, henbilkigun, henbilkioke, henbilkioeken, henbilkion, ibili, ibilik, ibilitze, ibiltzen, lebilke, lebilkiake, lebilkidake, lebilkie, lebilkigu, lebilkinake, lebilkioke, lebilkizue, lebilkizuke, lebilzke, lebilzkiake, lebilzkidake, lebilzkie, lebilzkigu, lebilzkinake, lebilzkioke, lebilzkizue, lebilzkizuke, nabilen, nabilke, nabilkiake, nabilkian, nabilkie, nabilkie, nabilkien, nabilkik, nabilkin, nabilkinake, nabilkinan, nabilkio, nabilkioke, nabilkion, nabilkizu, nabilkizue, nabilkizuen, nabilkizuke, nabilkizun, nenbilen, nenbilke, nenbilkeen, nenbilkiake, nenbilkian, nenbilkie, nenbilkieken, nenbilkien, nenbilkinake, nenbilkinan, nenbilkioke, nenbilkioken, nenbilkion, nenbilkizue, nenbilkizueken, nenbilkizuen, nenbilkizuke, nenbilkizukeen, nenbilkizun, zabilta, zabilzate, zabilzke, zabilzkete, zabilzkidake, zabilzkidakete, zabilzkidate, zabilzkie, zabilzkiake, zabilzkie, zabilzkigu, zabilzkigute, zabilzko, zabilzkioke, zabilzkokete, zabilzkio, zabilzkite, zabilzkidake, zabilzkidakete, zabilzkidate, zabilzkie, zabilzkiakete, zabilzkie, zabilzkigu, zabilzkigute, zabilzkidake, zabilzkidakete, zabilzkidate, zabilzkie, zabilzkiakete, zabilzkiote, zabilzkit, zebilen, zebilkeen, zebilkian, zebilkidakeen, zebilkidan, zebilkieken, zebilkien, zebilkigukeen, zebilkizun, zebilkizun, zebiltzan, zebilzkeen, zebilzkan, zebilzkidakeen, zebilzkidan, zebilzkieken, zebilzkien, zebilzkigukeen, zebilzkigun, zebilzkinan, zebilzkiokeen, zebilzkion, zebilzkizue, zebilzkizuen, zebilzkizukeen, zebilzkizun, zenbiltzan, zenbiltzaten, zenbilzke, zenbilzkeen, zenbilzkete, zenbilzkidake, zenbilzkida-*

*keen, zenbilzkidakete, zenbilzkidaketen, zenbilzkidan, zenbilzkidaten, zenbilzkieke, zenbilzkiekeen, zenbilzkiekete, zenbilzkieketen, zenbilzkien, zenbilzkieten, zenbilzkiguke, zenbilzkigukeen, zenbilzkigukete, zenbilzkiguketen, zenbilzkigun, zenbilzkiguten, zenbilzkioke, zenbilzkiokeen, zenbilzkikete, zenbilzkikoketen, zenbilzkion, zenbilzkioten* — “walk”;

diodaz, diogu, dioguz, diok, diosagu, diosat, diosate, dioska, dioskna, diosku, dioskuk, dioskute, dioskutez, dioskuz, dioskuzak, dioskuzu, dioskuzue, dioskuzez, dioskuzuz, diosnagu, diosnat, diosnate, diost, diostak, dioste, diostate, diostaz, diostazak, diostazu, diostazue, diostazez, diostazuz, diot, diote, diotez, diotse, diotsedaz, diotseg, diotseguz, diotsek, diotset, diotsete, diotsetez, diotsez, diotsezak, diotsezu, diotsezue, diotsezuez, diotsezuz, diotso, diotsodaz, diotsogu, diotsoguz, diotsok, diotsot, diotsote, diotsotez, diotsoz, diotsozak, diotsozu, diotsozue, diotsozuez, diotsozuz, diotsu, diotsudaz, diotsue, diotsuedaz, diotsuegu, diotsueguz, diotsuet, diotsuite, diotsuetez, diotsuez, diotsugu, diotsuguz, diotsut, diotsute, diotsutez, diotsuz, dioz, diozaa, diozaagu, diozaat, diozaate, diozak, diozana, diozanagu, diozanat, diozanate, diozu, diozue, diozuez, diozuz, esan, esango, esate, esaten, genioen, genioezen, geniosan, geniosnan, geniotsen, geniotsezen, geniotson, geniotsozen, geniotsuen, geniotsuezen, geniotsun, geniotsuzen, geniozaa, geniozana, hioen, hioezen, hioskun, hioskuzen, hiostan, hiostazen, hiotsen, hiotsezen, hiotson, hiotsozen, nioen, nioezen, niosan, niosnan, niostsuen, niostsuezen, niotsen, niotsezen, niotson, niotsozen, niotsun, niotsuzen, niozaan, niozanen, zenioen, zenioezen, zenioskun, zenioskuten, zenioskutezen, zenioskuzen, zeniostan, zeniostanen, zeniostatezen, zeniostazen, zeniooten, zenioetezen, zeniootsen, zeniotseten, zeniotsetez, zeniotsezen, zeniotson, zeniootsoten, zeniostsotzen, zeniootszen, zioen, zioezen, ziosan, ziosaten, zioskun, zioskuten, zioskutezen, zioskuzen, ziosnan, ziosnaten, ziostan, ziostaten, ziostatezen, ziostazen, zioten, ziotzezen, ziotsen, ziotseten, ziotsetez, ziotsezen, ziotson, ziotsoten, ziotsotezen, ziotsozen, ziotsuen, ziotsueten, ziotsuetez, ziotsuezen, ziotsun, ziotsuten, ziotsutezen, ziotsuzen, ziozaa, ziozaate, ziozana, ziozanate — “say”;

*bagentoz, bagentozkie, bagentozkik, bagentozkin, bagentozkio, bagentozkizu, bagentozkizue, bahentor, bahentorkie, bahentorkigu, bahentorkio, bahentorkit, baletor, baletorkie, baletorkigu, baletorkik, baletorkin, baletorkio, baletorkit, baletorkizu, baletorkizue, baletoz, baletozkie, baletozkigu, baletozkik, baletozkin, baletozkio, baletozkit, baletozkizu, baletozkizue, banentor, banentorkie, banentorkik, banentorkin, banentorkio, banentorkizu, banentorkizue, bazentoz, bazentozkide, bazentozkie, bazentozkiete, bazentozkigu, bazentozkigute, bazentozkio, bazentozkiote, bazentozkit, bazentozte, bentorkie, bentorkigu, bentorkik, bentorkin, bentorkio, bentorkit, bentorkizu, bentorkizue, betor, betoz, betozkie, betozkigu, betozkik, betozkin, betozkio, betozkit, betozkizu, betozkizue, dator, datorke, datorkiake, datorkidake, datorkie, datorkieke, datorkigu, datorkiguke, datorkik, datorkin, datorkinake, datorkio, datorkioke, datorkit, datorkizu, datorkizue, datorkizuke, datorrela, datozi, datoza, datoze, datozi, datoziak, datoziakate, datoziakie, datoziakieke, datoziakigu, datoziakiguke, datoziakik, datoziakin, datoziakine, datoziakio, datoziakit, datoziaku, datoziakue, datoziakuke, etor, etorri, etorriko, etortze, etortzen, gatoz, gatozen, gatozke, gatoziakie, gatozkian, gatozkie, gatozkieke, gatozkiene, gatozkit, gatozkin, gatozkinak, gatozkinan, gatozki, gatozkiok, gatozkion, gatozkizu, gatozkizue, gatozkizuke, gatozkizuen, gatozkizuke, gatozkizun, gentozen, gentozke, gentozkeen, gentozkiak, gentozkie, gentozkieke, gentozkien, gentozkinak, gentozkinan, gentozkiok, gentozkioken, gentozkizue, gentozkizueken, gentozkizuen, gentozkizuke, gentozkizukeen, gentozkizun, hator, hatorke, hatorkidake, hatorkie, hatorkieke, hatorkigu, hatorkiguke, hatorki, hatorkiok, hatorkit, hentorke, hentorkidake, hentorkidakeen, hentorkidan, hentorkie, hentorkieke, hentorkien, hentorkiguke, hentorkigun, hentorkioke, hentorkiokeen, hentorkion, hentorren, letorke, letorkiak, letorkidake, letorkie, letorkigu, letorkinak, letorkioke, letorkizue, letorkizuke, letozke, letozkiak, letozkidake, letozkie, letozkigu, letozkinak, letozkioke, letozkizue, letozkizuke, nator, natorke, natorkiak, natorkian, natorkie, natorkieke, natorkien, natorkik, natorkin, natorkinak, natorkinan, natorkio, natorkion, natorkizu, natorkizue, natorkizuen, natorkizuke, natorkizun, natorren, nentorke, nentorkeen, nentorkiak, nentorkian, nentorkie, nentorkieke, nentorkien, nentorkinak, nentorkinan, nentorkio, nentorkiokeen, nentorkion, nentorkizue, nentorkizuen, nentorkizuke, nentorkizukeen, nentorkizun, nentorren, zatoz, zatozke, zatozkete, zatozkidake, zatozkidakete, zatozkidate, zatozkie, zatozkieke, zatozkiekete, zatozkite, zatozkigu, zatozkiguke, zatozkigukete, zatozkigute, zatozkio, zatozkiok, zatozkiokete, zatozkiote, zatozkit, zatozte, zentozen, zentozke, zentozkeen, zentozketen, zentozkidake, zentozkidakeen, zentozkidakete, zentozkidaketen, zentozkidan, zentozkidaten, zentozkie, zentozkieke, zentozkiekete, zentozkien, zentozkieten, zentozkigu, zentozkiguke, zentozkigukete, zentozkiguketen, zentozkigun, zentozkiguten, zentozkioke, zentozkiokeen, zentozkiokete, zentozkioketen, zentozkion, zentozkioten, zentozten, zetorkeen, zetorkian, zetorkidakeen, zetorkidan, zetorkie, zetorkien, zetorkigukeen, zetorkigun, zetorkinan, zetorkiokeen, zetorkion, zetorkizue, zetorkizuen, zetorkizukeen, zetorkizun, zetorren, zetozen, zetozkeen, zetozkian, zetozkidakeen, zetozkidan, zetozkie, zetozkien, zetozkigukeen, zetozkigun, zetozkinan, zetozkiokeen, zetozkion, zetozkizue, zetozkizuen, zetozkizukeen, zetozkizun — “come”;*

*bagindoaz, bagindoazkie, bagindoazkik, bagindoazkin, bagindoazkio, bagindoazkizu, bagindoazkizue, bahindoa, bahindoakie, bahindoakigu, bahindoakio, bahindoakit, balihoa, balihoakie, balihoakigu, balihoakik, balihoakin,*

*balihoakio, balihoakit, balihoakizu, balihoakizue, balihoaz, balihoazkie, balihoazkigu, balihoazkik, balihoazkin, balihoazkio, balihoazkit, balihoazkizu, balihoazkizue, banindoa, banindoakie, banindoakin, banindoakio, banindoakizu, banindoakizue, bazindoaz, bazindoazkide, bazindoazkie, bazindoazkiete, bazindoazkigu, bazindoazkigute, bazindoazkio, bazindoazkiote, bazindoazkit, bazindoazte, bihoa, bihoaz, bihoazkie, bihoazkigu, bihoazkik, bihoazkin, bihoazkio, bihoazkit, bihoazkizu, bihoazkizue, bindoakie, bindoakigu, bindoakik, bindoakin, bindoakio, bindoakit, bindoakizu, bindoakizue, doa, doake, doakiakie, doakidake, doakie, doakieke, doakigu, doakiguke, doakik, doakin, doakinake, doakio, doakioke, doakit, doakizu, doakizue, doakizuke, doala, doaz, doazela, doazke, doazkiakie, doazkidake, doazkie, doazkigu, doazkiguke, doazkik, doazkin, doazkinake, doazkio, doazkiokie, doazkit, doazkizu, doazkizue, doazkizuke, gindoazzen, gindoazke, gindoazkeen, gindoazkiakie, gindoazkian, gindoazkie, gindoazkieken, gindoazkien, gindoazkinake, gindoazkinan, gindoazkioke, gindoazkioken, gindoazkion, gindoazkizue, gindoazkizuen, gindoazkizuke, gindoazkizukeen, gindoazkizun, goaz, goazen, goazke, goazkiakie, goazkian, goazkie, goazkien, goazkik, goazkin, goazkinake, goazkinan, goazkio, goazkiokie, goazkion, goazkizu, goazkizue, goazkizuen, goazkizuke, goazkizun, hindoake, hindoakeen, hindoakidake, hindoakidan, hindoakie, hindoakieke, hindoakien, hindoakiguke, hindoakigun, hindoakiokie, hindoakioken, hindoakion, hindoan, hoa, hoake, hoakidake, hoakie, hoakieke, hoakigu, hoakiguke, hoakio, hoakiokie, hoakit, joan, joango, joate, joaten, lihoake, lihoakiakie, lihoakidake, lihoakie, lihoakiguke, lihoakinake, lihoakiokie, lihoakizue, lihoakizuke, lihoazke, lihoazkiakie, lihoazkidake, lihoazkie, lihoazkiguke, lihoazkinake, lihoazkiokie, lihoazkizue, lihoazkizuke, nindoake, nindoakeen, nindoakiakie, nindoakian, nindoakie, nindoakieke, nindoakizue, nindoakizuen, nindoakizuke, nindoakizukeen, nindoakizun, nindoan, noa, noake, noakiakie, noakian, noakie, noakieke, noakien, noakik, noakin, noakinake, noakinan, noakio, noakiokie, noakizun, noakizue, noakizueke, noakizuen, noakizuke, noakizun, noan, zohoakeen, zohoakian, zohoakidakeen, zohoakidan, zohoakie, zohoakien, zohoakiokie, zohoakigun, zohoakinan, zohoakioken, zohoakion, zohoakizue, zohoakizuen, zohoakizuke, zohoakizun, zohoan, zohoazen, zohoazkeen, zohoazkian, zohoazkidakeen, zohoazidan, zohoazkie, zohoazkien, zohoazkigukeen, zohoazkigun, zohoazkinan, zohoazkiokie, zohoazkion, zohoazkizue, zohoazkizuen, zohoazkizuke, zohoazkizun, zindoaz, zindoazken, zindoazke, zindoazketen, zindoazkidake, zindoazkidakeen, zindoazkidakete, zindoazkidaketen, zindoazkidan, zindoazkidaten, zindoazkie, zindoazkieken, zindoazkiekete, zindoazkieketen, zindoazkien, zindoazkieten, zindoazkiguke, zindoazkigukeen, zindoazkigukete, zindoazkiguketen, zindoazkigun, zindoazkiguten, zindoazkiokie, zindoazkioken, zindoazkiokete, zindoazkioketen, zindoazkion, zindoazkioten, zindoazsten, zoaz, zoazke, zoazkete, zoazkidake, zoazkidakete, zoazkide, zoazkie, zoazkieke, zoazkiekete, zoazkiete, zoazkigu, zoazkiguke, zoazkigukete, zoazkigute, zoazkio, zoazkioke, zoazkiokete, zoazkiote, zoazkit, zoazte — “go”.*

*Example 9.6.2.* In Fig. S14, we further apply the aforementioned word clustering algorithm to topic extraction and machine translation (see Table S1 for text source).

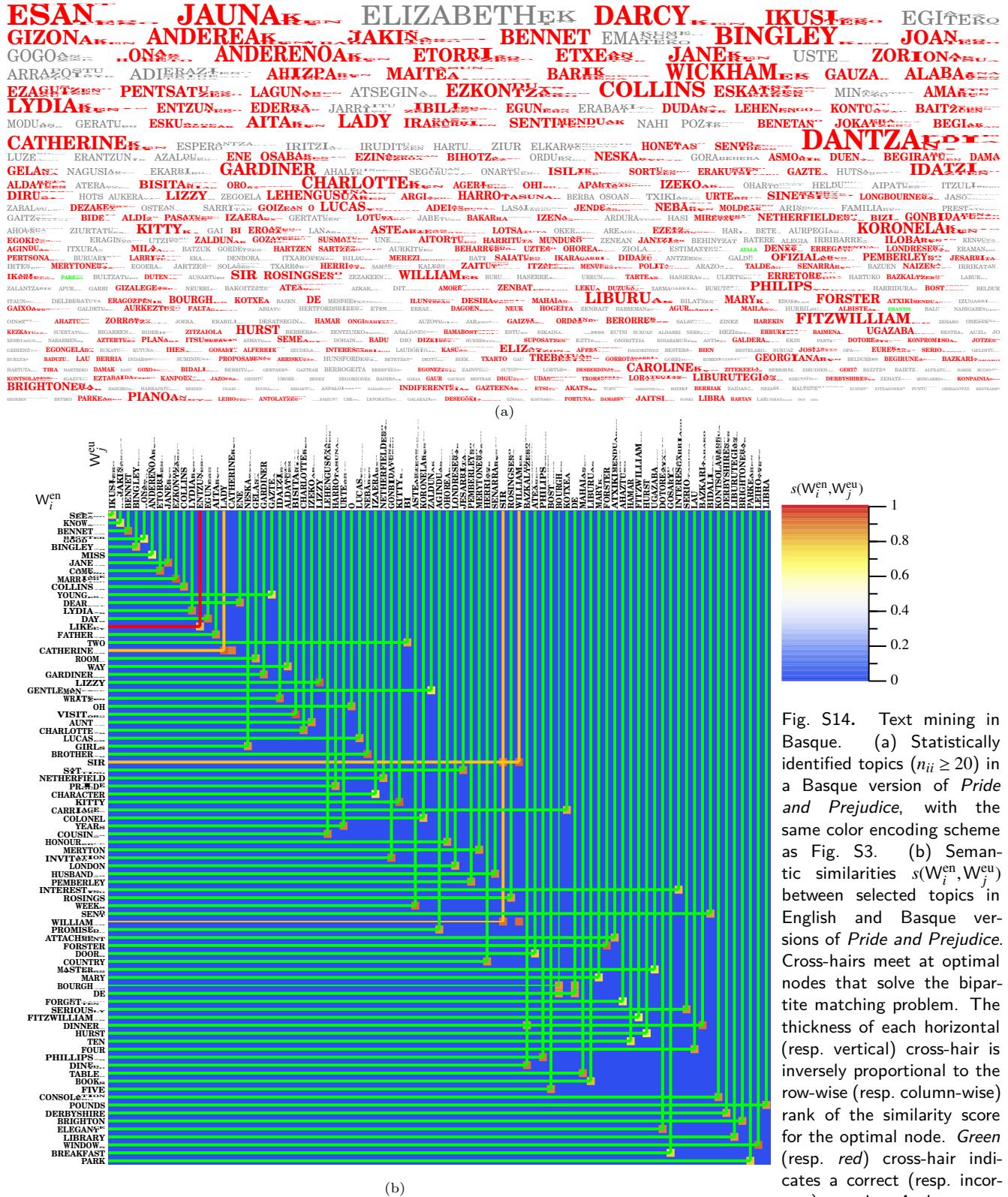


Fig. S14. Text mining in Basque. (a) Statistically identified topics ( $n_{ii} \geq 20$ ) in a Basque version of *Pride and Prejudice*, with the same color encoding scheme as Fig. S3. (b) Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{eu}})$  between selected topics in English and Basque versions of *Pride and Prejudice*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) rank of the similarity score for the optimal node. *Green* (resp. *red*) cross-hair indicates a correct (resp. incorrect) match. *Amber* cross-hair marks a link between distinct concepts that share the same hypernyms.

## 9.2 Modified Porter stemming algorithm for Korean

We will be interested in analyzing Modern Korean literary works, written according to the South Korean standard after 1971. This means that texts are spelt out in the native Korean alphabet (known as *Hangul*), with spaces between words.<sup>108</sup> Each *Hangul* block represents exactly one syllable (see Algorithm 9.8 below for romanization of *Hangul*), hence we can use the “uniform reading speed” approximation (§??).<sup>109</sup> We are not going to deal with an archaic mixed script style (where Sino-Korean vocabulary is spelt out with *Hanja*, the Chinese characters used in the Korean language), which is mostly restricted to academic and legal settings nowadays.<sup>110</sup>

In a Modern Korean translation (see Table S1 for text source) of *Pride and Prejudice*, the opening line reads

재산이 많은 미혼 남성이라면 반드시 아내를 필요로 한다는 말은 널리 인정되는 진리이다.

Its structure can be analyzed as follows:

Hangul:	재산이]	많은	미혼	남성이	라면	반드시	아내를
Mixed:	財產이]	많은	未婚	男性이	라면	반드시	아내를
Romanization:	jaesan-i	manh-eun	mihon	nameseong-ilamyeon	bandeusi	anae-leul	
Gloss:	fortune.NOM	much.PRS.DET	unmarried	male.be.COND	certainly	wife.ACC	
Hangul:	필요로	한다는	말은	널리	인정되는	진리이다	
Mixed:	必要로	한다는	말은	널리	認定되는	眞理이다	
Romanization:	pilyo-lo	handa-neun	mal-eun	neolli	injeong-doeneun	jinli-ida	
Gloss:	necessity.INS	do.IND.NPST.DET	word.TOP	widely	affirm.PAS.PRS.DET	truth.be.IND.NPST	

where word stems are shown in red, and various functions of the suffixes include ACC = accusative, COND = conditional, DET = determiner, INS = instrumental, NOM = nominative, NPST = non-past, PAS = passive, PRS = present, TOP = topic. In principle, we want to automatically remove all kinds of grammatical suffixes from a Korean word, no matter the remaining stem is Sino-Korean (*i.e.* of ultimately Chinese origin) or not. In practice, we must admit to certain limitations in our current technology:

- Native Korean stems are sometimes subject to consonant changes. For example, 널리 neolli “widely”, 넓다 neolbda “to be wide” and 넓비 neobi “width” are supposed to share the same stem, but their *Hangul* spellings in the first syllable have subtle differences. We will try our best to accommodate to such irregular Korean inflections, at the cost of misclassifying certain native Korean words.
- While Sino-Korean stems are immutable in inflections, they suffer from another serious problem: polysemy. Many words with distinguished pronunciations in Chinese become (nearly) homophonic<sup>111</sup> when adapted into Korean, so the disuse of Chinese characters in Korean writing causes ambiguities. Such pairs of *Hangul* homographs [and also (near) homophones] include 인정(認定) injeong “acknowledgement” vs. 인정(人情) injeong “humanity”; 화장(化粧) hwajang “cosmetics” vs. 화장(火葬) hwajang “cremation” etc. We can do nothing to solve such polysemy puzzles in our current work.
- Some native Korean words are spelt in the same way as Sino-Korean words, despite having very different meanings.<sup>112</sup> For example, we have 거리 geoli “road, street” vs. 거리(距離) geoli “distance”; 시장 sijang “hunger” vs. 시장(市場) sijang “market” and 시장(市長) sijang “mayor” etc. Again, we can do nothing right now, as we face Korean documents free from Chinese characters.
- Depending on context (which determines delicate shades of meaning), a given English word might be translated into either a native Korean word or a Sino-Korean word, with drastically different spellings. For example, both 마음 maeum and 심(心) sim mean “heart” in Korean, but there is no simple way to build such links algorithmically.<sup>113</sup> Traditionally, Koreans would learn each *Hanja* by juxtaposing its native Korean gloss with its Sino-Korean pronunciation, in a formula like 心 마음 심. For technical reasons, we have chosen not to feed our computer a long list of *Hanja* readings for rote learning.

<sup>108</sup>This distinguishes Korean texts from Chinese and Japanese ones. In the latter two cases, word segmentation is a non-trivial task.

<sup>109</sup>Vowel length is not marked in Korean spelling. Its effects on reading speed are ignored in our analysis.

<sup>110</sup>When Korean is written in the mixed script, rules for word segmentation are also somewhat different from the pure *Hangul* style. The current North Korean standard does not employ Chinese characters, but its spelling/spacing rules for Sino-Korean vocabulary are different from the practice in South Korea.

<sup>111</sup>For near homophones, we allow differences in vowel lengths that are not marked in *Hangul* spelling.

<sup>112</sup>There is a similar problem for words transliterated from English. For example, 다시 dasi could be a *Hangul* version of “Darcy” or a native Korean word meaning “again”; 제인 jein could be a *Hangul* version of “Jane” or a Sino-Korean word 제인(諸人) “all people”. In the present study, we draw on neighboring words [cf. Algorithm 3.25(4)] to decide whether 다시 means “again” (which is to be removed as a stop word), and regard all words in the form of 제인X as derivatives of “Jane”.

<sup>113</sup>There is a similar situation in Japanese, where the Chinese character 心 can be pronounced either こころ kokoro (native Japanese) or しん shin (Sino-Japanese). In Japanese, “heart” is always written out using the *kanji* 心, irrespective of pronunciation. Thus in text processing, it is trivial to identify native Japanese and Sino-Japanese versions of the same concept. However, even when Korean is written in the mixed script style, only the Sino-Korean vocabulary can be spelt out with *Hanja*, which typographically obscures the connection between native Korean and Sino-Korean counterparts.

Since the phonetic writing system of Modern Korean is a compromise between accuracy and efficiency, so will be our approximate clustering algorithm based on [45].

### 9.2.1 Stop words and transliterations

Korean stop words can carry various agglutinative suffixes. Instead of enumerating all of them in a single list, we will define several string patterns that match the criterion for stop words.

**Definition 9.7** (Korean stop words). Define the pattern **beKorean** as any one from the following list:

계세요, 계셔, 계셔서, 계셔야, 계셔요, 계셨기, 계셨느냐, 계셨다, 계셨습니까, 계셨습니다, 계셨어, 계셨어요, 계셨음, 계시겠다, 계시겠습니다, 계시겠어, 계시겠어요, 계시고, 계시기, 계시느냐, 계시는, 계시는데, 계시니, 계시니까, 계시다, 계시더니, 계시려고, 계시면, 계시지만, 계신, 계신다, 계실, 계심, 계십니까, 계십니다, 계십시오, 아녀, 아녀서, 아녀야, 아녀요, 아녔기, 아녔냐, 아녔다, 아녔습니까, 아녔습니다, 아녔어, 아녔어요, 아녔음, 아니겠다, 아니겠습니까, 아니겠어, 아니겠어요, 아니고, 아니기, 아니냐, 아니니, 아니니까, 아니다, 아니다, 아니더니, 아니면, 아니세요, 아니셔, 아니셔서, 아니셔야, 아니셔요, 아니셨기, 아니셨냐, 아니셨다, 아니셨습니까, 아니셨습니다, 아니셨어, 아니셨어요, 아니셨음, 아니시겠다, 아니시겠습니다, 아니시겠어, 아니시겠어요, 아니시고, 아니시기, 아니시냐, 아니시니, 아니시니까, 아니시다, 아니시더니, 아니시면, 아니시지만, 아니신, 아니신데, 아니실, 아니심, 아니십니까, 아니십니다, 아니어서, 아니어야, 아니어요, 아니었기, 아니었음, 아니지만, 아닌, 아닌데, 아닐, 아님, 아닙니까, 아닙니다, 였기, 였느냐, 였다, 였습니까, 였습니다, 였어, 였어요, 였음, 예요, 이겠다, 이겠습니다, 이겠어, 이겠어요, 이고, 이기, 이냐, 이니, 이니까, 이다, 이다, 이더니, 이면, 이세요, 이셔, 이셔서, 이셔야, 이셔요, 이셨기, 이셨냐, 이셨다, 이셨습니까, 이셨습니다, 이셨어, 이셨어요, 이셨음, 이시겠다, 이시겠습니다, 이시겠어, 이시겠어요, 이시고, 이시기, 이시냐, 이시니, 이시니까, 이시다, 이시더니, 이시면, 이시지만, 이신, 이신데, 이실, 이심, 이십니까, 이십니다, 이어서, 이어야, 이었기, 이었느냐, 이었다, 이었습니다, 이었습니다, 이였어, 이였어요, 이었음, 이에, 이에요, 이지만, 인데, 입니까, 입니다, 있겠다, 있겠습니다, 있겠어, 있겠어요, 있고, 있기, 있느냐, 있는, 있는데, 있다, 있다, 있더니, 있습니까, 있습니다, 있어, 있어서, 있어야, 있어요, 있었기, 있었느냐, 있었다, 있었습니까, 있었습니다, 있었어, 있었어요, 있었음, 있으니, 있으니까, 있으면, 있으세요, 있으셔, 있으셔서, 있으셔야, 있으셔요, 있으셨기, 있으셨냐, 있으셨다, 있으셨습니까, 있으셨습니다, 있으셨어, 있으셨어요, 있으셨음, 있으시겠다, 있으시겠습니다, 있으시겠어, 있으시겠어요, 있으시고, 있으시기, 있으시냐, 있으시니, 있으시니까, 있으시다, 있으시더니, 있으시면, 있으시지만, 있으신, 있으신데, 있으실, 있으심, 있으십니까, 있으십니다, 있을, 있음, 있지만.

Define the pattern **notbeKorean** as any one from the following list:

없겠다, 없겠습니다, 없겠어, 없겠어요, 없고, 없기, 없느냐, 없는, 없는데, 없다, 없다, 없더니, 없습니까, 없습니다, 없어, 없어서, 없어야, 없어요, 없었기, 없었느냐, 없었다, 없었습니다, 없었습니다, 없었어, 없었어요, 없었음, 없으니, 없으니까, 없으면, 없으세요, 없으셔, 없으셔서, 없으셔야, 없으셔요, 없으셔서, 없으셨기, 없으셨어, 없으셨다, 없으셨나, 없으셨습니까, 없으셨습니다, 없으셨어, 있으셨어요, 있으셨음, 있으셨겠다, 없으시겠다, 없으시겠어, 없으시고, 없으시기, 없으시냐, 없으시니, 없으시니까, 없으시다, 없으시더니, 없으시면, 없으시지만, 없으신, 없으신데, 없으실, 없으심, 없으십니까, 없을, 없음, 있지만.

Define the pattern **becomeKorean** as any one from the following list:

돼, 돼라, 돼서, 돼야, 돼요, 됐기, 됐느냐, 됐다, 됐습니까, 됐습니다, 됐어, 됐어요, 됐음, 되, 되겠다, 되겠습니다, 되겠어, 되겠어요, 되고, 되기, 되느냐, 되는, 되는데, 되니, 되니, 되니까, 되다, 되더니, 되려고, 되면, 되세요, 되셔, 되셔서, 되셔야, 되셔요, 되셨기, 되셨느냐, 되셨다, 되셨습니까, 되셨습니다, 되셨어, 되셨어요, 되셨음, 되시겠다, 되시겠습니다, 되시겠어, 되시겠어요, 되시고, 되시기, 되시느냐, 되시는, 되시는데, 되시니, 되시니까, 되시더니, 되시라, 되시려고, 되시면, 되시지만, 되신, 되신다, 되실, 되심, 되십니까, 되십니다, 되십시오, 되어, 되어라, 되어서, 되어야, 되어요, 되었기, 되었음, 되자, 되지만, 된, 된다, 될, 됩니까, 됩니다, 됩시다, 됩시오.

Define the pattern **haveKorean** as any one from the following list:

가져라, 가져서, 가져야, 가져요, 가겼기, 가겼느냐, 가졌다, 가겼습니까, 가겼습니다, 가졌어, 가졌어요, 가졌음, 가지겠다, 가지겠습니다, 가지겠어, 가지겠어요, 가지고, 가지기, 가지느냐,

가지는, 가지는데, 가지니, 가지니까, 가지다, 가지더니, 가지려고, 가지면, 가지세요, 가지셔, 가지셔서, 가지셔야, 가지셔요, 가지셨기, 가지셨느냐, 가지셨다, 가지셨습니까, 가지셨습니다, 가지셨어, 가지셨어요, 가지셨음, 가지시겠다, 가지시겠습니다, 가지시겠어, 가지시겠어요, 가지시고, 가지시기, 가지시느냐, 가지시는, 가지시는데, 가지시니, 가지시니까, 가지시더니, 가지시라, 가지시려고, 가지시면, 가지시지만, 가지신, 가지신다, 가지실, 가지심, 가지십니까, 가지십니다, 가지십시오, 가지어라, 가지어서, 가지어야, 가지어요, 가지었기, 가지었음, 가지자, 가지지만, 가진, 가진다, 가질, 가짐, 가집니까, 가집니다, 가집시다, 가집시오, 갖겠다, 갖겠습니다, 갖겠어, 갖겠어요, 갖고, 갖기, 갖느냐, 갖는, 갖는다, 갖는데, 갖더니, 갖습니까, 갖습니다, 갖아, 갖아라, 갖아서, 갖아야, 갖아요, 갖았기, 갖았느냐, 갖았다, 갖았습니까, 갖았습니다, 갖았어, 갖았어요, 갖았음, 갖으니, 갖으니까, 갖으려고, 갖으면, 갖으세요, 갖으셔, 갖으셔서, 갖으셔야, 갖으셔요, 갖으셨기, 갖으셨느냐, 갖으셨다, 갖으셨습니까, 갖으셨습니다, 갖으셨어, 갖으셨어요, 갖으셨음, 갖으시겠다, 갖으시겠습니다, 갖으시겠어, 갖으시겠어요, 갖으시고, 갖으시기, 갖으시느냐, 갖으시는, 갖으시는데, 갖으시니, 갖으시니까, 갖으시더니, 갖으시라, 갖으시려고, 갖으시면, 갖으시지만, 갖으신, 갖으신다, 갖으실, 갖으심, 갖으십니까, 갖으십니다, 갖으십시오, 갖은, 갖을, 갖음, 갖읍시다, 갖읍시오, 갖자, 갖지만.

Define the pattern **KoreanStopRoot** as any one from the following list:

가운데, 가져, 캐, 거기, 거니, 거라, 거야, 거였, 거예, 거의, 건너편, 겹니, 곳, 귀하, 그러나, 꿔나겠, 꿔나느, 꿔나는, 꿔나니, 꿔나더, 꿔나려, 꿔나세, 꿔나서, 꿔나셨, 꿔나시, 꿔나신, 꿔나십, 꿔나자, 꿔나지, 꿔납니, 꿔납시, 꿔났, 꿔보, 꿔본, 꿔불, 꿔봄, 꿔봄, 꿔봐, 꿔봤, 나네, 나도, 나라면, 나로, 나를, 나만, 나보다, 나중, 나처, 나한테, 난다, 남게, 남는, 남들, 남에, 남은, 남을, 남의, 남한테, 내개, 내겐, 너도, 너를, 너만, 너만큼, 너무, 너에, 너와, 너한, 너희, 누구, 누군가, 다음, 당신, 대개, 대신, 대하여, 덕분, 도중, 되겠, 되나, 되던, 되려, 되리, 되었, 되지, 뒤에, 들끼리, 때까, 때는, 때도, 때때로, 때로, 때마다, 때만, 때만큼, 때면, 때문, 때보다, 때부터, 때에, 때조차, 라고, 마찬가지, 만큼, 많, 몇, 모두, 모든, 못하, 못했, 무렵, 무엇, 무척, 뭐, 뭔, 뭣, 바에, 별로, 본인, 불과, 뿐이, 뿐입, 속으로, 수시로, 스스로, 심지어, 십상이, 아니, 아니요, 아무, 아직, 안에, 안으로, 안은, 안을, 않, 얘들, 얘야, 어느, 어디, 어딘, 어딜, 어때, 어땠, 어띠, 어떤, 어멸, 어멈, 어떻, 어찌, 어쩔, 어쩜, 언제, 얼마, 없게, 없겠, 없었, 없으, 없이, 없지, 여러, 여부, 여전, 였, 예가, 오로지, 왜냐, 외치곤, 우리, 위까, 위로, 위를, 위부터, 위에, 위의, 이라고, 아래, 이랬, 이러, 이런, 이럴, 이럼, 이렇, 이로, 이리, 이만, 이분, 이와, 이외, 이전, 이제, 이쯤, 이후, 일도, 일쑤, 일은, 일을, 일이, 일인, 일하, 있, 자기, 자네, 자신, 자체, 저기, 저까, 저로, 저만큼, 저처럼, 저한테, 저희, 전과, 전까, 전보다, 전부터, 전에, 전으로, 전쯤, 전혀, 절대, 정도, 제가, 제게, 중에, 중이, 지금, 지난, 지마느, 지마세, 지마셨, 지마시, 지마십, 지말겠, 지말더, 지말려, 지말았, 지말지, 지맙니, 지맙시, 차라리, 최소, 텐데, 하거, 하겠, 하고, 하끈, 하기, 하긴, 하나, 하네, 하는, 하니, 하다, 하더, 하던, 하든, 하려, 하며, 하면, 하면서, 하세, 하셨, 하시, 하여, 하지, 한다, 한데, 할지, 함께, 해서, 해야, 해왔, 해준, 해줄, 했, 현재, 혹시, 혹은, 혼자, 후에, 훨씬.

Define the pattern **KoreanStopString** as any one from the following list:

가끔, 가장, 간에, 갖게, 갖은, 같이, 개, 거, 거고, 거나, 거는, 거닌, 거다, 거면, 거야, 거에요, 거예요, 거요, 건, 걸, 걸로, 겁니다, 것이다, 것입니다, 게, 게다가, 결코, 겹, 고로, 곧, 곧잘, 그, 그대, 그래서, 그러나, 그러면, 그리, 그리고, 그리고, 그만, 그만큼, 금방, 기껏해야, 꿔나, 꿔나고, 꿔나기, 꿔나는, 꿔나니, 꿔나라, 꿔나면, 꿔나서, 꿔나셔, 꿔나신, 꿔나실, 꿔나심, 꿔나야, 꿔나요, 꿔나자, 꿔난, 꿔난다, 꿔날, 꿔남, 나, 나는, 나에, 나와, 나의, 난, 남, 내, 내가, 내개, 내내, 내로, 너, 너는, 너무, 너의, 널, 네, 네가, 누가, 늘, 다, 다들, 다소, 대로, 대체로, 댁, 댁은, 더, 더구나, 더욱, 덜, 데, 데가, 데서, 도로, 동안, 되어, 된, 등, 듯, 듯도, 따라, 따라서, 따름, 땐, 때, 맑찢, 또, 라도, 리, 리가, 마, 마나, 마니, 만, 매우, 맨, 먼저, 모든, 못, 무릇, 무슨, 미리, 바람에, 밖에, 번째, 벌써, 뿐, 사이, 사이에, 서로, 서로서로, 수, 수가, 수는, 아니면, 아마, 아마도, 아주, 안, 앞으로, 얘, 어느, 어따, 어떤, 어떻게, 어차피, 언제나, 얼마나, 없, 여기, 옆에, 예, 오히려, 온, 온갖, 왜, 웬, 이, 이는, 이들, 이래, 이로, 이를, 이리, 이미, 이번, 이상, 이상는, 이상에, 이상으로, 이상은, 이상을, 이상의, 이상이, 일, 자, 자넨, 자넬, 자주, 잘, 쟤, 저, 저기, 저나, 저는, 저도, 저를, 저의, 적, 적이, 전, 전에, 전의, 제, 조금, 좀, 좀처럼, 좀체, 줄, 즉, 지, 지가, 지는, 지마는, 지마니, 지마셔, 지마신, 지마실, 지마심, 지만, 지말, 지말고, 지말기, 지말면, 지말아, 지말자, 지맑, 짹이, 쪽, 채, 채로, 퍽, 하고, 하여금, 하지만, 한, 한번, 한번도, 해, 해도, 해라, 해요, 했어, 후, 후로, 혼히.

Define **doKorean** as the following pattern:

(하|하거|하게|하겠다|하겠습니다|하겠어|하겠어요|하고|하곤|하기|하긴|하느냐|하는|하는데|하니|하니  
까|하다|하더니|하던|하든|하려|하려고|하면|하세요|하셔서|하셔야|하셔요|하셨기|하셨느냐|하셨  
다|하셨습니까|하셨습니까|하셨습니다|하셨어요|하셨어요|하시|하시겠다|하시겠습니까|하시겠어  
요|하시고|하시기|하시느냐|하시는|하시는|하는데|하시니|하시니까|하시더니|하시라|하시려고|하시지  
만|하신|하신다|하실|하심|하십니까|하십니다|하십시오|하여라|하여서|하여야|하였기|하였음|하자|하지  
만|한|한대|할|함|합니까|합니다|합시다|합시오|해|해라|해서|해야|해요|했|했기|했느녀|했다|했습니다|했습  
니다|했어요|했어요|(으)까|요)<sub>m</sub>.

Define **DemPronKorean** as the following pattern:

(그|그려|그리|그립|그립|그만|그쳐|그쳤|그친|그칠|그침|그침)|((으)이|이)(것|곳|보다|쪽)|(으)저|(거|건|걸|게|들|때|번))**X**.

If a word exactly matches any one of the following:

(beKorean|becomeKorean|haveKorean|KoreanStopRoot)**X**,  
**KoreanStopString**,  
(으)관|못|인)(**doKorean**)**X**,  
**DemPronKorean**,  
(전|종)**beKorean**,

then we consider it a Korean stop word. All the Korean stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list.  $\square$

To accommodate to consonant changes in Korean inflections, we need to convert between *Hangul* and its romanized form.

**Algorithm 9.8** (Romanization of *Hangul*). *We transliterate each individual Hangul syllable as follows:*<sup>114</sup>

(1) Replace the initial *Hangul* letter by

ㄱ	ㄲ	ㄴ	ㄷ	ㄸ	ㄹ	ㅁ	ㅂ	ㅃ	ㅅ	ㅆ	ㅇ	ㅈ	ㅉ	ㅊ	ㅋ	ㅌ	ㅍ	ㅎ
g	kk	n	d	tt	l	m	b	pp	s	ss	ø	j	jj	ch	k	t	p	h

(2) Replace the medial *Hangul* vowel letter by

ㅏ	ㅐ	ㅑ	ㅒ	ㅓ	ㅕ	ㅕ	ㅕ	ㅗ	ㅘ	ㅙ	ㅚ	ㅛ	ㅜ	ㅞ	ㅚ	ㅟ	ㅢ	ㅡ	ㅔ	ㅣ
a	ae	ya	yae	eo	ye	yeo	ye	o	wa	wae	oe	yo	u	wo	we	wi	yu	eu	ui	i

(3) Replace the final *Hangul* consonant letter (if present) by

ㄱ	ㄲ	ㄴ	ㄷ	ㄹ	ㅁ	ㅂ	ㅅ	ㅆ	ㅈ	ㅊ	ㅋ	ㅌ	ㅍ	ㅎ	ㄺ	ㄻ	ㄻ	ㄻ	ㄻ	
g	kk	n	d	l	m	b	s	ss	ng	j	ch	k	t	p	h	lh	bs	gs	nj	nh

For a string  $\hat{\sigma}$  written in the *Hangul* alphabet, we construct  $HgRom(\hat{\sigma})$  in three steps:

(1') Transliterate each syllable according to the aforementioned rules.

(2') Insert letter Q as the syllable separator.

(3') Convert double letters into their capitalized equivalents.

The reverse operation  $RomHg$  is defined so as to ensure that  $\hat{\sigma} = RomHg(HgRom(\hat{\sigma}))$ .

Example 9.8.1. We have  $HgRom(볶아)$  =  $baKQe$  and  $HgRom(한글)$  =  $hanQgeul$ .

<sup>114</sup>As mentioned before, we concern ourselves with Modern Korean spelling, which eliminates some consonant and vowel forms in Middle Korean: such as ㅠ kki “meal” (from earlier ㅠ pski), ㄸ ttaleum “only, just” (from earlier ㄸ stolom), 땅 ttang “soil, land” (from earlier 땅 sta), ㅕ ttae “time” (from earlier ㅕ psta).

### **9.2.2 Effective spelling and essential root**

Like Japanese, the Korean language has two separate numeral systems, one native and one Chinese, which appear in different contexts. The following algorithm handles this special feature of Korean.

**Algorithm 9.9** (Clustering of Korean numerals). Define **CntKorean** as the following pattern:

(가지|개|대|마리|명|번|분|사람|잔|장|채)(Ø|가|과|와|의|이).

For a string  $\hat{\sigma}$ , we define  $\text{KorNum}(\hat{\sigma})$  by the following procedures:

- (1) Check whether  $\hat{\sigma}$  matches exactly the pattern  $(\emptyset|\text{제})(\mathbf{D}|\text{구|넷|다|섯|두|둘|삼|세|셋|십|아홉|어덟|열|육|이|일|일곱|칠|팔|하나|한})_m(\emptyset|\text{CntKorean}|\text{시}| 째)$ , where  $\mathbf{D} = (0|1|2|3|4|5|6|7|8|9)$ . If no, define  $\text{KorNum}(\hat{\sigma}) = \hat{\sigma}$  and quit; if yes, go to the next step.

- (2) Work on  $\hat{\sigma}$  as follows:

- (2.1)  $\text{Do } \sim(\text{시})|\text{째}|\text{CntKorean} \rightarrow \emptyset$ ,  $\text{제}\sim \rightarrow \emptyset$ .
  - (2.2)  $\text{Do } \sim(\text{구})|\text{아}|\text{춥} \rightarrow '9$ ,  $\sim\text{넷} \rightarrow '4$ ,  $\sim\text{다섯} \rightarrow '5$ ,  $\sim(\text{두})|\text{둘}|^{\circ} \rightarrow '2$ ,  $\sim(\text{삼})|\text{세}|\text{셋} \rightarrow '3$ ,  $\sim(\text{십})|\text{열} \rightarrow '0$ ,  $\sim(\text{어})|\text{육} \rightarrow '6$ ,  $\sim(\text{여덟})|\text{팔} \rightarrow '8$ ,  $\sim(\text{하나})|\text{한} \rightarrow '1$ ,  $\sim(\text{일곱})|\text{칠} \rightarrow '7$ .
  - (2.3) *Repeat the last step.*
  - (2.4)  $\text{Do } ' \rightarrow \emptyset$ .
  - (2.5) *Call the result so far as  $\hat{\sigma}'$ , and define  $\text{KorNum}(\hat{\sigma}') = \sum \hat{\sigma}'_i$ .*

**Algorithm 9.10** (Regularization of Korean verbs). *For a Korean word  $\hat{\sigma}$ , its “regularized verb form”  $\text{RegVb}(\hat{\sigma})$  is defined through the following procedures:*

- (1) If  $\hat{\sigma}$  contains 閃 and does not start with 閃, then replace it by 閃 $\hat{\sigma}$ ; otherwise, leave it intact.

- (2) *Replace*

(∅ 새 시)(파 퍼)(라 란 랑 래 랬 러 런 렁 레 .StatusBadRequest)X~ 파랑	(거 까 꺼 꺼)(마 만 黑色 매 맸 머 먼 黑色 매 ADVERTISEMENT)X~ 까黑恶
(검 검정 깜장)~ 까맣	(노 누)(라 란 랑 랑 래 랬 러 런 렁 렁 레 .StatusBadRequest)X~ 노랑
(발 발 빨 빨)X~ 발강	(발 벌 빨 빨)(가 간 갛 개 갰 거 건 겋 개 갰)~ 빨강
(보 뽀)(야 안 양 애 הייתי)X~ 뽀얗	(부 뿌)(여 연 영 예 הייתי)X~ 뽀얗
(하 허)(야 안 양 양 애 הייתי)X~ 하얗	(하 허)(야 안 양 양 애 הייתי)X~ 하얗
X <sup>ε</sup> (곱 뽑 수줍 씹 입 잡 접 좁)~ X <sup>뾰</sup>	X <sup>ε</sup> (낳 낳 낳 낳 낳)~ X <sup>핥</sup>
X <sup>ε</sup> (닫 돋 쏟 얻)~ X <sup>뜯</sup>	X <sup>ε</sup> (벗 빼 앗 솟 씻)~ X <sup>奚</sup>
그(쳐 쳤 치 친 칠 침 침)~ 그쳤	그만~ 그많
노래(doKorean 가는 로 를 에 와 의)X~ 노랫	노래(doKorean 가는 로 를 에 와 의)X~ 노랫
돌아(보 본 볼 봄 봅 봐 봤)X~ 했了个	만(나 난 날 남 납 났)X~ 罵
불~ 발강	불~ 발강
여(느 는 는데 니 세 셔 셨 시 실 심 십)~ 엽	열(겠 고 기 더 니 면 어 었 지만)~ 엽
일(은 을)~ 일하다	허영심~ 허용씩
(연다 엷) 엽	건(ঠ 다) 결다
이름(ঠ 은 을) 성함	이름(ঠ 은 을) 성함

- (3) Do 마(여|셨|시|신|실|심|심) ~→ 마-swire,  $\hat{x}$ -어(보|본|불|봄|봄|봐|봤) X ~→  $\hat{x}'$ -다, where  $\hat{x}'$  derives from  $\hat{x}$ , after replacing the latter's final 른 (if present) with 니.<sup>115</sup>

**Algorithm 9.11** (Korean normalization). *For a Korean word  $\hat{\sigma}$ , its normalized form  $\text{KN}(\hat{\sigma})$  is constructed in the following steps:*

- (1)  $Do^{X^\epsilon}(\text{가게})|\text{가구}| \text{거실}|\text{구실}|\text{소리}|\text{시시}|\text{시집}|\text{일시}|\text{포도}|\text{포함}) \sim X_{\frac{\epsilon}{1-\epsilon}} \text{ on RegVb}(\hat{\sigma}).^{116}$

<sup>115</sup>In practice, such a replacement needs both HgRom and RomHg to take effect.

<sup>116</sup>The second syllable in many Sino-Korean stems may resemble a Korean suffix. To avoid misidentification of suffixes, we need to protect such stems by a blocker syllable 韓.

## (2) Replace

(∅ 만)(아드님 아들)~ 엎酙	(가르 가른 가를 가름 가瞽 갈라 갈랐)~ 걸톺					
(다 아 어)(보 이) 본 볼 봄 봄 봐 봤(자))X~ 꿰명	(매일 하루)~ 날	(묻 물었 물으)~ 몬뜰				
(벼서 벼야 벼요 벗기 벗나 벗다 벗습 벗어 벗음 빈데 밉니)X~ 꿰툑	(봬 뺐 뵈 빈 벌 쁨 뵙)~ 보					
(이 르 이 른 이 를 이 름 이 릅 일 려 일 렸)X~ 꿰읊		가(려(고) 렸 리 린 릴 립 립)~ 꿰월				
결혼식~ 결혼��	경(doKorean)X~ 경けば	그리워~ ∅	늘어(놓 서 선 설 심 섭)~ 놓늙	다(르 른 를 름 롭)X~ 꿰꿰		
떨어(져 졌 지 진 질 짐 집)X~ 뿔뿔	떨어뜨X~ 뿔	메리턴~ 메꿰팅	목사관~ 목사꽝	목려받~ 잃꿰	배우자~ 배웨찢	부인(doKorean)X~ 꿰꿰
비(겠 고 기 녀 나 나 더 면 세 셔 셨 시 신 실 심 십 어 었 우 운 울 움 웁 위 웠 지)~ 꿰툑	쓸데~ 꿰꿰		쓸쓸(doKorean)X~ 꿰꿰			
어(렵)(려(우 운 울 움 웁 위 웠))~ 어俵	어리석~ 어리썩	일으(커 켰 키 킨 킬 kim kip)X~ 꿰	일이X~ 일하다	조지아나~ 조찢않		
친(doKorean)X~ 친구	침~ 챔	하녀장~ 하녀쨍	햄~ 햄	(벼 비 빈 벌 빔)~ 꿰툑	달(라고 라서 래)~ 주다	

## (3) If the string so far begins with any one of the following list:

가문, 가치, 강요, 거리, 거만, 게임, 겨우, 겨울, 견해, 결과, 결심, 결함, 경고, 경련, 경쟁, 고개, 고기, 고려, 고문, 고민, 고의, 고집, 곰곰, 과거, 과시, 관대, 관련, 관리, 관심, 구름, 구입, 군대, 권리, 기대, 기만, 기분, 기여, 기운, 기울, 기질, 기흔, 끝내, 나름, 나무, 나이, 날아, 노골, 노출, 단언, 달리, 당당, 대가, 대기, 대기, 대단, 대답, 대면, 대지, 도구, 도도, 도리, 도서, 동기, 동네, 동요, 동의, 두려, 둘러, 드려, 떠올, 런던, 리지, 마디, 마련, 마을, 마음, 마치, 맞이, 머리, 메리, 모습, 목숨, 몹시, 무례, 무시, 무심, 물건, 물려, 물려, 물론, 물리, 미래, 바람, 바보, 반대, 방도, 방심, 방어, 방임, 방자, 방치, 방해, 배려, 배신, 변심, 부대, 부리, 부분, 부여, 부자, 분개, 분들, 불리, 비중, 비참, 사건, 사고, 사과, 사라, 사랑, 사실, 사이, 사치, 상대, 서재, 섬세, 성함, 세대, 세련, 세심, 소개, 소리, 소심, 소중, 속도, 손님, 수고, 수입, 수치, 시골, 시기, 시도, 시선, 식구, 신랑, 신분, 신중, 아래, 아첨, 안심, 암시, 압니, 압시, 양도, 양심, 얼굴, 여름, 여보, 여성, 여인, 여자, 역할, 연구, 연기, 연대, 연습, 연인, 연주, 연출, 열심, 열중, 염려, 예의, 오래, 오만, 오해, 외출, 요구, 요리, 용서, 우려, 우아, 우울, 원래, 원인, 위대, 유감, 유리, 유일, 유지, 의도, 의무, 의문, 의미, 의심, 의아, 의자, 의지, 이유, 이의, 이자, 이치, 이해, 일반, 일어, 일치, 임대, 입구, 자랑, 자리, 자립, 자만, 자세, 자연, 자질, 잠시, 장교, 장군, 장담, 장면, 재산, 적당, 전이, 절대, 점심, 점잖, 정당, 정리, 정보, 정신, 정오, 정중, 제인, 조건, 조심, 조지, 조치, 존재, 존중, 주로, 주요, 주의, 주인, 중대, 중심, 중요, 중재, 지나, 지시, 지인, 지참, 지체, 지출, 직접, 진실, 진심, 진지, 질녀, 집중, 차이, 차지, 책임, 처리, 처신, 처치, 천재, 초과, 초대, 초래, 추구, 축하, 충고, 친구, 태도, 편안, 편지, 피로, 하녀, 하인, 한가, 한심, 합의, 합치, 호감, 호기, 화랑, 화려, 화해, 확고, 확보, 확신, 확실, 확인,

then insert 헷 as the third syllable.

## (4) Replace

(∅ 고)남성~ 남甥	(∅ 고)남아~ 뼈값	(∅ 고)남자~ 꿰맞	(∅ 고)남편~ 남 '>'	(∅ 만)(따님 딸)~ 땋톨	(가나)실~ 쌩깎
(계수 동서 새언니 시누이 울케 제수 처제 처형 형수)~ 켁 '>'					(고모 숙모 이모 작은어머니 큰어머니)X~ 恭敬
(금일 오늘)~ 금일	(길겠 길고 길기 길대 길더 길면 길어 길었 길지 깁니)~ 롱롤				

(남동생 브라더 아우 오라버님 오빠 형님 형제)X~ 쁠遞	(내 명)일~ 묽遞	(년 단 란 잔)말(씀 이)~ 옛遞
(누나 누이 시스터 언니 여동생 자매)X~ 찢遞		(도움 돕)X~ 훑遞
(동서 매부 시동생 시매부 시숙 시아주버니 제부 처남 형부)~ ἴլլ		(드 울)(려 렸 리 린 릴 림 립)~ 주
(마누라 와이프 집사람)~ 아내	(먼저 우선 처음 첫)~ 휑뚯	(모르 모른 모를 모름 모릅 몰라 몰랐)~ 핥닭
(바로 바르 바른 바를 바름 바롭 발라 발랐)X~ 쾰הלך	(봄 봐 봤)~ 보	(부르 부른 부를 부름 부릅 불리 불렀)~ 쾰닭
(부친 아버지 아빠 어비)~ 뻔啄	(숙부 아저씨 외삼촌 친삼촌)X~ 꾹죽	(씨 썼 쓴 쓸 씁)X~ 쓰다
(안다 알)~ 핥닭	(애 이 야)기~ 훑클	(아이 애들 자식)~ 깐똸
(자르 자른 자를 자름 자롭 잘라 잘랐)~ 쫙쫙	(어머니 어미 엄마)~ 모친	(이어 이었 이은 이을 이음 잇)~ 잇렁
(추우 추운 추울 추움 추워 추웠 춥겠 춥고 춥기 춥다 춥더 춥습 춥지)~ 깐깐		
X <sup>e</sup> (의 대)답~ X牒	가르(쳐 쳤 치 친 칠 침 침)X~ 툇톨	갈(겠 고 더 려 면 아 았 지만)~ 챘
강(doKorean)X~ 강장	걸렸X~ 걸리다	검~ 검컴
고말(doKorean)X~ 말하다	구(doKorean)X~ 굵큐	겁(니 시)X~ 행거걸다
기(녀 니 세 셔 셨 시 신 실 심 심 십)~ 롱률	길이~ 롱률	까(보 본 볼 봄 봄 봐 봤)~ 쌩깎
나(애 았 으 은 을 음 읍)X~ 냇	나(온 을 옴 읍 와 왔)X~ 解釋	놀(라 란 랄 람 람 랐)X~ 짢랠
느(꺼 깼 끼 낀 낄 낌 찝)~ 느끼	다가X~ 앓הייתי	다시X~ 탈찌
달(도 만 에 은 이 째 씀)~ 월	달라(요)X~ 액즉	달(doKorean)X~ 튕脿
돌아(가 간 갈 감 갑 갔)X~ 뭣깎다	돌아(오 온 올 옴 읍 와 왔)X~ 뭣왁라	드(느 는 니 세 셔 셨 시 신 실 심 십)X~ 뜰뜰다
드(세 셔 셨 시 신 실 심 십)X~ 먹다	드리(나 난 날 남 남 났)X~ 엷툠	들(리 렸 르 른 를 름 롭)X~ 짢뼉
들(렸 리 린 릴 림 립)X~ 듣다	들려(오 온 올 옴 읍 와 왔)X~ 듣다	들어(가 간 갈 감 갑 갔)X~ 慝깎다
들어(도 주 준 줄 줄 줍 줘)X~ 듣다		들어(라 서 야 요)~ 듣다
들어(서 선 설 섬 섭)X~ 앴서다	들어(오 온 올 옴 읍 와 왔)X~ 였워라	들었X~ 듣다
듭(니 시)X~ 뜰뜰다	떠(나 난 날 났)~ 쾰툠	뜻~ ("'", 넝)
만(doKorean)X~ 핥축	말(려 렸 리 린 릴 림 립)~ 마르	리디아~ 리떡얄
모(아 았 여 였 으 은 을 음 읍 이 인 일 임 입)~ 깰툠		마리아~ 꺌툠
		맘~ 마핥
		몸~ 뿔遞

吳~	무도회~	미움X~	바(라 란 랄 랍 팠)~	빈~	빈아들~	밤~	법~	벗~
𦥑	무또회	밉	핥	빨꼴	빨해들	냉훑	렁깛	߻
보잘~	비(doKorean)X~	상(doKorean)X~	샘~	생각X~	세(우 운 울 움 웁 워 웠)X~			
裴	鄙覩	상행	샘	생깎	세			
손(실 해)~	숨~	싫~	심(doKorean)X~			쓸(데 모)X~		
𦥑	숨	싫길	ｾﾝ			斧斧		
아(느 는 니 제 셨 시 신 실 심 십)X~	앓~	어(례 덮 리 린 릴 림 립)~	어(째 깼 찌 찐 쩔 쩜 찝)~	업~				
𦥑	앓	앓	어래	어챘	업			
엿~	오(가 간 같 감 갑 갓)X~	읊~	읊~	웃~	웃~	의(의 거)(doKorean)X~		
𦥑	읊	읊풀	랠랠	윙笯	락뿔	暉		
이(루 문 률 룸 룹 뤼)~	일(doKorean)X~	일라이자~	잃~	입(에 은 을 이)X~				
앳쳤	일이라는	엘리자베스	읊	읊쯤	읊쯤			
자(라 란 랄 람 랍 랐)~	잡수X~	적(혀 혔 히 힌 힐 힘 힙)X~	전(doKorean)X~	젊~				
끓풀	먹다	쓰다	휑	휑	휑	.�		
정(doKorean)X~	좋아(요)X~	주기X~	주시(doKorean)X~	준다X~	줘~	즐거X~	집~	
뒹딩	뭐	주끼	좆꺾	주	주	즐겁다	집	
차(례 덮 리 린 릴 림 립)~	참~	청(doKorean)X~	춤~	크리스마스~	탓~			
찾횧	핥	.�	뗐	.�	.�	.�	.�	
꾀(doKorean)X~	합~	해리엇~	화가Ӄ~	힐~	간	(긴 진 데 길)		
엎_fn	ܚ	해렐엘	화깎Ӄ~	.�	가다	롱롤		
(의 매)주(의 까지 나 만큼 부터 을) ((에 였 일)X)	(거 걸 짚)	(둔 둔)	(둔다 들 둘)	(본 본다 불 봄)				
왁	.�	.�	.�	.�				
(안 앓)	달	드	들려(의 여 요)	들어	쓸	아	오(게 지)	일
�	월	를	듣다	듣다	쓰다	황황	오면	.�
						양		준

and do (두겠|두고|두기|두느|두는|두니|두데|두려|두면|두세|두셔|두셨|두시|두신|두실|두심|두십|두에|두었|두자|두지|둔다|듭니|듭시)~→�, 나오(의|겠다|겠습니|다|겠어|겠어|요|고|기|너|라|느|나|는|는데|나|니까|더|니|려고|면|세|요|셔|셔|셔|야|셔|요|겼|기|겼|느|나|겼|다|겼|습|니|까|겼|습|니|다|겼|어|요|겼|어|요|겼|읍|시|겼|다|시|겼|습|니|다|시|겼|어|시|겼|어|요|시|고|시|기|시|느|나|시|는|시|는|데|시|니|시|니까|시|더|내|시|래|시|려|고|시|면|시|지|만|신|신|다|실|심|십|니|까|십|나|십|시|자|지|만)X~→�, 매X<sub>1</sub>=(년|달|변|월|일)(의|까지|나|만큼|부터|을|이|을)((에|일)X)→X<sub>1</sub>', where X<sub>1</sub>' results from doing 년→연, 달→월 on X<sub>1</sub>.

(5) Do X<sub>ε</sub>(들어|부당|부인|부추|상당|상세|상실|상심|상의|상인)~→X<sub>학</sub>.

(6) Do (가거라|가려고|가면|가서|가서|가서|가서야|가서요|가시느|녀|가시|려|고|가|시|면|가|신|가|신|다|가|실|가|심|간|다)X~→가다, (거|녀|거|느|거|는|거|니|거|세|거|서|거|셨|거|시|거|신|거|심|거|십|건|데|걸|겠|걸|고|걸|기|걸|다|걸|더|걸|려|걸|면|걸|자|걸|지|겁|니|겁|시)X~→행|거|걸|다, (낫|종|최|선)~→�, (커|컸|크|큰|클|큼|큽)~→�, 고|말|(의|야)~→고|.�, .�(추|춘|출|축|취|쳤)~→�, 남~→�, 달(려|덮|리|린|릴|림)~→�, 들(겠|려|고|면|어)X~→�, 들|들|다, 들(에|으|은|을|음)X~→�, 듣|다, 말~→�, 오(겠|고|기|너|라|느|나|는|니|더|니|도|록|라|려|고|면|세|요|셔|셔|겼|시|신|다|실|심|십|니|까|십|나|십|시|자|지|만)X~→와|라, 이(어|었)~→�, 입~→�, 점~→�, 주(계|까)X~→주, (나|가|나|간)→�, (온|온|다|울|옴|옵|니|까|옵|나|다|옵|시|다|옵|시|오|와|와|서)→와|라.

**Algorithm 9.12** (Korean effective spelling). Set  $\mathbf{V} = (a|e|i|o|u|y)$ . For a Korean word  $\hat{\sigma}$ , its effective spelling EffSpell( $\hat{\sigma}$ ) is constructed in the following steps:

(1) Do  $\hat{x}$ 한테X→ $\hat{x}$  before applying Korean normalization.

(2) Do 리학→�학, 하학→황학, 한학→�학, 할학→�학, 함학→�학, 해학→�학.

(3) Do  $\hat{x}$ (beKorean|becomeKorean|(의|당)doKorean)X→ $\hat{x}$ .

(4) Define **KoreanNounSuff** as the pattern (의|건|걸|것|게|쪽)(의|들|뿐|뿐만)(의|가|파|랑|로|로|부|터|로|서|로|씨|를|보|다|아|에|에|게|에|개|서|에|다|에|다|가|에|서|와|으로|으로|부|터|으로|서|로|씨|을|의|이|이|랑|하고)(의|같|이|까지|까지|나|는|대|로|도|따|라|마|다|마|저|만|만|큼|보|다|보|다|도|부|터|씩|야|야|마|로|은|이|나|조|차|조|차|도|쯤|처|럼). Do  $\hat{x}\hat{x}$ 적X→ $\hat{x}\hat{x}$  and  $\sim\hat{x}\hat{x}$ (의|가|감|객|경|게|과|관|권|금|기|력|료|방|범|별|복|부|비|사|상|생|서|석|성|세|소|수|식|실|아|어|용|원|인|자|장|지|직|초|탕|판|편|품|학|형|화|회)KoreanNounSuff(의|커|녕)→ $\hat{x}\hat{x}$ .

(5) *Do*  $\sim\hat{\chi}$ ( $\emptyset|\text{감|개|계|껏|꾸|러|기|꾼|님|동|이|맹|이|보|어|치|오|우|음|이|쟁|이|질|짜|리|째|투|성|이|히|}$ ) **KoreanNounSuff**( $\emptyset|\text{다|는|데|듯|이|라|는|데|라|셔|려|무|나|련|마|는|렵|망|정|수|록|야|말|고|자|마|지|언|정|커|녕|테|니|테|니|까|텐|데|}$ )  $\rightarrow \hat{\chi}$ .

(6) *Apply KorNum.*

(7) *Do the following in a sequel, searching for the longest match in each step:*

(7.1) *Remove one of these word-final patterns if it follows a character:*

길래, 끼리, 네요, 다며, 다시피, 담, 되며, 라도, 라며, 라면서, 만, 만은, 며, 면서, 어요, 에서, 예선, 요, 이란, 이며, 이요, 인, 하며.

(7.2) *Remove one of these word-final patterns if it follows a character:*

내, 는, 대, 라면, 래, 를, 은, 을, 이군, 이라뇨, 이라니, 채

(7.3) *Remove one of these word-final patterns if it follows a character:*

가, 까, 끔, 늘, 니와, 더러, 도, 든, 라, 뿐, 아, 아라, 야, 어라, 애, 오, 자

and *do*  $\sim\hat{\chi}$ ( $\emptyset|\text{년|단|란|잔|}$ )  $\text{다} \rightarrow \hat{\chi}$ ,  $\sim\hat{\chi}$ ( $\emptyset|\text{났|댔|했|겠|}$ )( $\emptyset|\text{어|}$ )  $\rightarrow \hat{\chi}$  in the same sweep.

(7.4) *Define KoreanVerbSuff as any one of the following patterns:*

$\emptyset$ , 가, 같, 거, 겠, 고, 고말, 구나, 구만, 구먼, 군요, 기, 기만, 기엔, 까지, 깨, 는, 다, 다니, 더니, 데, 도록, 든지, 들, 따라, 려고, 로,로부터,로서,로써,를,마다,마저,만큼,면,보다,부터,서,서부터,셔,씩,야,어,에,에게,에게서,에겐,에다,에다가,에서,에서부터,으로,으로부터,으로서,으로써,은,은데,을,음,의,이,자마자,져서,조차,중,지만,쯤,처럼,처럼,하곤,해선,( $\emptyset|\text{는|은|}$ )커녕, ( $\emptyset|\text{이|}$ )나, ( $\emptyset|\text{랑|와|이랑|하고|}$ )( $\emptyset|\text{는|만|은|}$ ), ( $\emptyset|\text{남|다|답|라|립|자|잡|고|}$ ) $\text{고|니|}$ .

*Do*  $\hat{\chi}$ (**KoreanVerbSuff**) $_m \rightarrow \hat{\chi}$ .

(7.5) *Remove one of these word-final patterns if it follows a character:*

게, 네, 는, 니, 더냐, 더라, 던, 데, 디, 세, 시, 어라, 자.

(7.6) *Do*  $\sim\hat{\chi}$ ( $\emptyset|\text{으|}$ )( $\emptyset|\text{셨|시|신|실|심|십|}$ )( $\emptyset|\text{겠|았|았|겠|었|었|겠|}$ )( $\emptyset|\text{습|읍|}$ )  $\rightarrow \hat{\chi}$ .

(7.7) *Remove one of these word-final patterns if it follows a character:*

겨, 겠, 구, 군, 굴, 굽, 귀, 꿔, 기, 긴, 길, 김, 김, 급, 려, 련, 렸, 리, 릴, 림, 립, 반, 여, 였, 우, 운, 울, 움, 웃, 워, 웠, 이, 인, 일, 임, 입, 쪘, 지, 진, 질, 짐, 짐, 추, 춘, 출, 춤, 축, 춰, 췔, 혀, 헌, 헐, 험, 험, 험, 험.

(7.8) *Do*  $\sim\hat{\chi}$ ( $\emptyset|\text{대|로|}($ **겹|답|롭|맞|스|러|스|럽|)** $)$ ( $\emptyset|\text{겨|우|겨|운|겨|울|겨|움|겨|워|겨|웠|다|우|다|운|다|울|다|움|다|위|다|웠|로|우|로|운|로|울|로|음|로|웁|로|워|로|웠|우|운|울|음|웁|위|웠|})$   $\rightarrow \hat{\chi}$ .

(8) *Do*  $\sim\hat{\chi}$ (**려|러|서|셔|서|**)  $\rightarrow \hat{\chi}$ ,  $\hat{\chi}$ (**빼|뺐|빠|뿔|뿜|뿜|**) $X \rightarrow \hat{\chi}$ **빠**,  $\hat{\chi}$ (**려|리|린|림|**) $X \rightarrow \hat{\chi}$ **리**.

(9) *Apply HgRom.*

(10) *Do*  $nolQlX \sim \sim nolQ$ ,  $Q(\emptyset|g|l|s)V_mSX \rightarrow \emptyset$ ,  $QsiQkX \rightarrow \emptyset$ ,  $\sim(\emptyset|b)Q(\emptyset|eub|seub)(\emptyset|Q)(\emptyset|de|deo|di|ge|go|ja|n|ne|neu|ni|sel|si)(\emptyset|l|m|n) \rightarrow \emptyset$ .

(11) *Apply RomHg and do*  $\sim\hat{\chi}$ ( $\emptyset|\text{나|시|으|시|}$ )  $\rightarrow \hat{\chi}$ .<sup>117</sup>

(12) *Do*  $\hat{\chi}$ ( $\text{느|으|}$ ) $X \rightarrow \hat{\chi}$ .

**Algorithm 9.13** (Korean essential root). *Let  $\hat{\sigma}$  be the effective spelling of a Korean word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:*

(1) *Do*  $\sim\hat{\chi}$ (**리|으|아|서|어|어|서|음|이|**)  $\rightarrow \hat{\chi}$ .

(2) *Apply HgRom and do*  $\sim(b|h|m|s|S) \rightarrow \emptyset$ ,  $\sim(d|Qleu(\emptyset|b|h|m|n|Qleo|s|S)) \rightarrow l$ .

(3) *Apply RomHg and do*  $\sim\omega \rightarrow \emptyset$ .

<sup>117</sup>To save execution time, one could have manipulated only the romanized strings, without converting back and forth between *Hangul* and its transliteration. However, we choose to state our algorithm using multiple applications of *HgRom* and *RomHg*, to enhance the readability of our substitution rules.

### 9.2.3 Approximate clustering

There will not be any mutation rules in our algorithm for Korean, so  $\text{SimpHrdTest} = \text{HrdTest}$ .

**Algorithm 9.14** (Simple heredity test). *Suppose that  $\hat{\alpha}$  and  $\hat{\beta}$  are both Hangul strings. The Boolean-valued function  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if at least one of the following five conditions holds:*<sup>118</sup>

- (i)  $\hat{\alpha} = \hat{\beta}$ ;
- (ii)  $\hat{\beta} = \hat{\alpha} \sqcup \cdot$ ;
- (iii)  $\text{HgRom}(\hat{\alpha})l = \text{HgRom}(\hat{\beta})$ ;
- (iv)  $\min\{\ell(\hat{\alpha}), \ell(\hat{\beta})\} \geq 2$ , AND  $\text{HgRom}(\hat{\alpha})' = \text{HgRom}(\hat{\beta})'$ , where the prime denotes an operation  $\sim \mathbf{V}_m(\emptyset | n) \rightarrow \emptyset$ ;
- (v)  $\min\{\ell(\hat{\alpha}), \ell(\hat{\beta})\} \geq 2$ , AND  $\text{HgRom}(\hat{\alpha})^* = \text{HgRom}(\hat{\beta})^*$ , where the asterisk denotes an operation  $Q(\text{jeog}|\text{man}|\text{Pun}|\text{seong})\mathbf{X} \rightarrow \emptyset$ .

**Algorithm 9.15** (Approximate clustering of Korean words). *The approximate clustering of a list of Korean words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  is completed in two stages:*

- (1) *We sort the list  $\{(\hat{\alpha}_1, \text{HgRom}(\text{EffSpell}(\hat{\alpha}_1))), \dots, (\hat{\alpha}_N, \text{HgRom}(\text{EffSpell}(\hat{\alpha}_N)))\}$  alphabetically according to the last component. If two consecutive neighbors in the alphabetized list  $\{(\hat{\alpha}_{(1)}, \text{HgRom}(\text{EffSpell}(\hat{\alpha}_1))), \dots, (\hat{\alpha}_{(N)}, \text{HgRom}(\text{EffSpell}(\hat{\alpha}_N)))\}$  satisfy  $\text{HrdTest}(\text{EffSpell}(\hat{\alpha}_{(n)}), \text{EffSpell}(\hat{\alpha}_{(n+1)})) = \text{FALSE}$ , where  $n \in \mathbb{Z} \cap [1, N]$ , then add a demarcation line between these two entries. In this way, the alphabetized list is divided into separate groups of words tagged with their effective spellings:  $\{g_1 = \{(\hat{\alpha}_{(1,1)}, \text{EffSpell}(\hat{\alpha}_{(1,1)})), \dots, (\hat{\alpha}_{(1,n_1)}, *)\}, \dots, g_M = \{(\hat{\alpha}_{(M,1)}, *), \dots, (\hat{\alpha}_{(M,n_M)}, *)\}\}$ , where each sublist also preserves alphabetic order. In particular, we have  $\hat{\alpha}_{(1,1)} = \hat{\alpha}_{(1)}$  and  $\hat{\alpha}_{(M,n_M)} = \hat{\alpha}_{(N)}$ .*
- (2) *For each group of words  $g_m = \{(\hat{\alpha}_{(m,1)}, *), \dots, (\hat{\alpha}_{(m,n_m)}, *)\}$  where  $m \in \mathbb{Z} \cap [1, M]$ , we augment it into a tagged entry  $G_m = (g_m, \text{HgRom}(\text{EffSpell}(\hat{\alpha}_{(m,1)})), \text{HgRom}(\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))))$ . The list  $\{G_1, \dots, G_M\}$  is sorted alphabetically, with respect to  $\text{HgRom}(\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})))$  (with higher priority), and  $\text{HgRom}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))$  (with lower priority). If two consecutive neighbors in the alphabetized list  $\{G_{(1)} = (g_{(1)}, \hat{\gamma}'_{(1)}, \hat{\gamma}''_{(1)}), \dots, G_{(M)} = (g_{(M)}, \hat{\gamma}'_{(M)}, \hat{\gamma}''_{(M)})\}$  satisfy*

$$\text{HrdTest}(\hat{\gamma}''_{(m)}, \hat{\gamma}''_{(m+1)}) = \text{FALSE}$$

AND

$$\text{HrdTest}(\hat{\gamma}'_{(m)}, \hat{\gamma}'_{(m+1)}) = \text{FALSE}$$

where  $m \in \mathbb{Z} \cap [1, M]$ , then add a demarcation line between these two entries. In this manner, the alphabetized list  $\{G_{(1)}, \dots, G_{(M)}\}$  is divided into separate groups of tagged clusters  $\{\Gamma_1 = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, \Gamma_K = \{G_{(K,1)}, \dots, G_{(K,m_K)}\}\}$ . Finally, the output list of word clusters  $\{\check{\Gamma}_1, \dots, \check{\Gamma}_K\}$  is constructed by discarding all the tags (effective spellings, essential roots) from each  $\Gamma_k$ , where  $k \in \mathbb{Z} \cap [1, K]$ .

*Example 9.15.1.* Our clustering algorithm correctly groups the following families of Korean words:

걷겠다, 걷겠습니다, 걷겠어, 걷겠어요, 걷고, 걷기, 걷느냐, 걷는, 걷는다, 걷는데, 걷더니, 걷습니까, 걷습니다, 걷자, 걷지만, 걸어, 걸어라, 걸어서, 걸어야, 걸어요, 걸었기, 걸었느냐, 걸었다, 걸었습니다, 걸었어, 걸었어요, 걸었음, 걸으니, 걸으니까, 걸으려고, 걸으면, 걸으세요, 걸으셔, 걸으셔서, 걸으셔야, 걸으세요, 걸으셨기, 걸으셨느냐, 걸으셨다, 걸으셨습니까, 걸으셨습니다, 걸으셨어, 걸으셨어요, 걸으셨음, 걸으시겠다, 걸으시겠습니다, 걸으시겠어, 걸으시겠어요, 걸으시고, 걸으시기, 걸으시느냐, 걸으시는, 걸으시는데, 걸으시니, 걸으시니까, 걸으시더니, 걸으시라, 걸으시려고, 걸으시면, 걸으시지만, 걸으신, 걸으신다, 걸으실, 걸으심, 걸으십니까, 걸으십니다, 걸으십시오, 걸은, 걸을, 걸음, 걸읍시다, 걸읍시오 — “walk”;

넓겠다, 넓겠습니다, 넓겠어, 넓겠어요, 넓고, 넓기, 넓다, 넓더니, 넓습니까, 넓습니다, 넓어, 넓어서, 넓어야, 넓어요, 넓었기, 넓었냐, 넓었다, 넓었습니다, 넓었어, 넓었어요, 넓었음, 넓으냐, 넓으니, 넓으니까, 넓으면, 넓으세요, 넓으셔, 넓으셔서, 넓으셔야, 넓으셔요, 넓으셨기, 넓으셨냐, 넓으셨다, 넓으셨습니까, 넓으셨습니다, 넓으셨어, 넓으셨어요, 넓으셨음, 넓으시겠다, 넓으시겠습니다, 넓으시겠어, 넓으시겠어요, 넓으시고, 넓으시기, 넓으시나, 넓으시니, 넓으시니까, 넓으시다, 넓으시더니, 넓으시면, 넓으시지만, 넓으신, 넓으신데, 넓으실, 넓으심, 넓으십니까, 넓으십니다, 넓은, 넓은데, 넓을, 넓음, 넓지만 — “wide”;

<sup>118</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical OR.

오겠다, 오겠습니다, 오겠어, 오겠어요, 오고, 오기, 오너라, 오느냐, 오는, 오는데, 오니, 오니까, 오더니, 오려고, 오면, 오세요, 오셔, 오셔서, 오셔야, 오셔요, 오셨기, 오셨느냐, 오셨다, 오셨습니까, 오셨습니다, 오셨어, 오셨어요, 오셨음, 오시겠다, 오시겠습니다, 오시겠어, 오시고, 오시기, 오시느냐, 오시는, 오시는데, 오시니, 오시니까, 오시더니, 오시라, 오시려고, 오시면, 오시지만, 오신, 오신다, 오실, 오심, 오십니까, 오십니다, 오십시오, 오자, 오지만, 온, 온다, 올, 옴, 옵니까, 옵니다, 옵시다, 옵시오, 와, 와라, 와서, 와야, 와요, 왔기, 왔느냐, 왔다, 왔습니까, 왔습니다, 왔어, 왔어요, 왔음 — “come”;

높겠다, 높겠습니다, 높겠어, 높겠어요, 높고, 높기, 높다, 높더니, 높습니까, 높습니다, 높아, 높아서, 높아야, 높아요, 높았기, 높았냐, 높았다, 높았습니까, 높았습니다, 높았어, 높았어요, 높았음, 높으냐, 높으니, 높으니까, 높으면, 높으세요, 높으셔, 높으셔서, 높으셔야, 높으셔요, 높으셨기, 높으셨냐, 높으셨다, 높으셨습니까, 높으셨습니다, 높으셨어, 높으셨어요, 높으셨음, 높으시겠다, 높으시겠습니다, 높으시겠어, 높으시겠어요, 높으시기, 높으시냐, 높으시니, 높으시니, 높으시니까, 높으시더니, 높으시라, 높으려고, 높이면, 높이지만, 높신, 높신다, 높실, 높심, 높십니까, 높십니다, 높십시오, 높아, 높아요, 높아라, 높어서, 높어야, 높어요, 높었기, 높었느냐, 높었다, 높았습니까, 높았습니다, 높았어, 높았어요, 높았음, 높으냐, 높으니, 높으니까, 높으면, 높으세요, 높으셔, 높으셔서, 높으셔야, 높으셔요, 높으셨기, 높으셨냐, 높으셨다, 높으셨습니까, 높으셨습니다, 높으셨어, 높으셨어요, 높으셨음, 높으시겠다, 높으시겠습니다, 높으시겠어, 높으시겠어요, 높으시기, 높으시냐, 높으시니, 높으시니까, 높으시더니, 높으시라, 높으려고, 높이면, 높이지만, 높신, 높신다, 높실, 높심, 높십니까, 높십니다, 높십시오, 높으십니까, 높으십니다, 높은, 높은데, 높을, 높음, 높지만 — “tall”;

가거라, 가겠다, 가겠습니다, 가겠어, 가겠어요, 가고, 가기, 가느냐, 가는, 가는데, 가니, 가니까, 가더니, 가라, 가려고, 가면, 가서, 가세요, 가셔, 가셔서, 가셔야, 가셔요, 가셨기, 가셨느냐, 가셨다, 가셨습니까, 가셨습니다, 가셨어, 가셨어요, 가셨음, 가시겠다, 가시겠습니다, 가시겠어, 가시겠어요, 가시고, 가시기, 가시느냐, 가시는, 가시는데, 가시니, 가시니까, 가시더니, 가시라, 가시려고, 가시면, 가시지만, 가신, 가신다, 가실, 가심, 가십니까, 가십니다, 가십시오, 가야, 가요, 가자, 가지만, 간다, 갑니까, 갑니다, 갑시다, 갑시오, 갔기, 갔느냐, 갔다, 갔습니까, 갔습니다, 갔어, 갔어요, 갔음 — “go”;

벗겠다, 벗겠습니다, 벗겠어, 벗겠어요, 벗고, 벗기, 벗느냐, 벗는, 벗는데, 벗더니, 벗습니까, 벗습니다, 벗어, 벗어라, 벗어서, 벗어야, 벗어요, 벗었기, 벗었느냐, 벗었다, 벗었습니다, 벗었습니다, 벗었어, 벗었어요, 벗었음, 벗으니, 벗으니까, 벗으려고, 벗으면, 벗으세요, 벗으셔, 벗으셔서, 벗으셔야, 벗으셔요, 벗으셨기, 벗으셨느냐, 벗으셨다, 벗으셨습니까, 벗으셨습니다, 벗으셨어, 벗으셨어요, 벗으셨어요, 벗으셨어요, 벗으셨음, 벗으시겠다, 벗으시겠습니다, 벗으시겠어, 벗으시겠어요, 벗으시고, 벗으시기, 벗으시느냐, 벗으시는, 벗으시는데, 벗으시니, 벗으시니까, 벗으시더니, 벗으시라, 벗으시려고, 벗으시면, 벗으시지만, 벗으신, 벗으신다, 벗으실, 벗으심, 벗으십니까, 벗으십니다, 벗으십시오, 벗은, 벗을, 벗음, 벗읍시다, 벗읍시오, 벗자, 벗지만 — “comb”;

나아, 나아라, 나아서, 나아야, 나아요, 나았기, 나았느냐, 나았다, 나았습니까, 나았습니다, 나았어, 나았어요, 나았음, 나으니, 나으니까, 나으려고, 나으면, 나주세요, 나으셔, 나으셔서, 나으셔야, 나으셔요, 나으셨기, 나으셨느냐, 나으셨다, 나으셨습니까, 나으셨습니다, 나으셨어, 나으셨어요, 나으셨음, 나으시겠다, 나으시겠습니다, 나으시겠어, 나으시겠어요, 나으시고, 나으시기, 나으시느냐, 나으시는, 나으시는데, 나으시니, 나으시니까, 나으시더니, 나으시라, 나으시려고, 나으시면, 나으시지만, 나으신, 나으신다, 나으실, 나으심, 나으십니까, 나으십니다, 나으십시오, 나은, 나을, 나음, 나웁시다, 나웁시오, 낫겠다, 낫겠습니다, 낫겠어, 낫겠어요, 낫고, 낫기, 낫느냐, 낫는, 낫는다, 낫는데, 낫더니, 낫습니까, 낫습니다, 낫지만 — “recover”;

흐르겠다, 흐르겠습니다, 흐르겠어, 흐르겠어요, 흐르고, 흐르기, 흐르느냐, 흐르는, 흐르는데, 흐르니, 흐르니까, 흐르더니, 흐르려고, 흐르면, 흐르세요, 흐르셔, 흐르셔서, 흐르셔야, 흐르셔요, 흐르셨기, 흐르셨느냐, 흐르셨다, 흐르셨습니까, 흐르셨습니다, 흐르셨어, 흐르셨어요, 흐르셨음, 흐르시겠다, 흐르시겠습니다, 흐르시겠어, 흐르시겠어요, 흐르시고, 흐르시기, 흐르시느냐, 흐르시는, 흐르시는데, 흐르시니, 흐르시니까, 흐르시더니, 흐르시라, 흐르시려고, 흐르시면, 흐르시지만, 흐르신, 흐르신다, 흐르실, 흐르심, 흐르십니까, 흐르십니다, 흐르십시오, 흐르자, 흐르지만, 흐른, 흐른다, 흐를, 흐름, 흐릅니까, 흐릅니다, 흐릅시다, 흐릅시오, 훌러, 훌러라, 훌러서, 훌러야, 훌러요, 훌렀기, 훌렀느냐, 훌렀다, 훌렀습니까, 훌렀습니다, 훌렀어, 훌렀어요, 훌렀음 — “flow”.

Moreover, the words 널리 “widely” and 너비 “width” also cluster correctly with the inflected forms of “wide”.

*Example 9.15.2.* In Fig. S15, we further apply the aforementioned word clustering algorithm to topic extraction and machine translation (see Table S1 for text source).

Following the standard practice in Korean dictionaries, we have inserted etymological information for loanwords from Chinese<sup>119</sup> and English, in the parentheses. This was done by using a manually built look-up table. The etyma in Chinese

<sup>119</sup> Due to partial conflations of native Korean and Sino-Korean vocabulary (as discussed earlier), a *Hanja* form may not apply to all the words in a certain cluster. Furthermore, since we have merged 숙모(叔母) “paternal aunt” and 이모(姨母) “maternal aunt” into the same cluster, but have not built a stacking mechanism for Chinese characters (for legibility reasons), the *Hanja* forms of such word stems are not displayed.

or English were not directly inferable from the Korean text under investigation, nor were they an integral part of our word clustering algorithm: we track the etymologies of the word stems so that we can assess the reliability of our experimental methods.

The Korean case shown here offers a telling example that meanings of words are determined by their numerical fingerprints, rather than their etymological sources. Korean words of native, Chinese and English origins all behave in a consistent way in the displayed results.

We caution our readers, however, that the performance of our numerical algorithm is indeed affected by mixed etymological sources of a given concept in Korean. From the control experiments in Fig. S15c, we see that certain correct matching pairs have similarity scores below the 0.7 threshold applicable to Fig. S15b. Part of this can be explained by the coexistence of native Korean and Sino-Korean stems, the latter usually reserved for abstract, formal and polite settings. See the (partial) list below:

English gloss	native Korean stem	Sino-Korean stem
“amiable”	상냥 sangnyang	친절(親切) chinjeol
“feel”	느끼 neukki	감정(感情) gamjeong, 기분(氣分) kibun
“forget”	잊 ij	망각(忘却) manggag
“love”	사랑 salang	애정(愛情) aejeong
“pride”	자랑 jalang	오만(傲慢) oman, 자부심(自負心) jabusim, 거만(倨慢) geoman
“read”	읽 ilg	독서(讀書) dogseo
“walk”	걷 geod	산책(散策) sanchaeg
“wish”	실 sip	원(願) won

Had we forcibly merged these native Korean stems with their Sino-Korean counterparts, before performing the numerical experiments in Fig. S15b, we would have achieved a slightly better yield for correct matchings with similarity scores above 0.7.

말	사람	씨(氏)	엘리자베스(ELIZABETH)
생각	불	할	대(對)
액기	죽	베넷(BENNET)	다시(DARCY)
빙리(BINGLEY)	같	좋	양(娘)
두	리디아(LYDIA)	결혼	부인(夫人)
편지(便紙)	쓰	방	제인(JANE)
친(親)	아버지	여성(女性)	콜린스(COLLINS)
사랑	행복(幸福)	늘	~고 싶
감정(感情)	일	캐서린(CATHERINE)	땅
확신(確信)	성격(性格)	집	사실(事實)
아	진짜	출신	집에
루커스(LUCAS)	마차(馬車)	방문(訪問)	만나
설명(說明)	kitti(KITTY)	점	집에
가능(可能)	청혼(請婚)	나	~고 싶
판금(苦痛)	청찬(稱讚)	나를(我)	만나
표정(表情)	친밀(親切)	나를(我)	만나
사춘(四寸)	초대(招待)	나를(我)	만나
찾았	비단(非難)	나를(我)	만나
책(冊)	포스터(FORSTER)	나를(我)	만나
관계(關係)	포스터(FORSTER)	나를(我)	만나
시선(視線)	주	나를(我)	만나
어려움(困難)	주	나를(我)	만나
오만(傲慢)	주	나를(我)	만나
연주(演奏)	피츠윌리엄(FITZWILLIAM)	나를(我)	만나
피아노(PIANO)	카드(CARD)	나를(我)	만나
성직(聖職)	교구(教區)	나를(我)	만나
조지아나(GEORGIANA)	교양(教養)	나를(我)	만나

엘리자베스(ELIZABETH) 부인(夫人) 제인(JANE)

다시(DARCY) 다른 갈 마음 정(正) 말

언니(JANE) 위컴(WICKHAM)

제인(JANE) 다른 갈 마음 정(正) 말

위컴(WICKHAM) 편지(便紙) 쓰

고 싶 땅 사실(事實) 집

기 땅 만날 대답(對答) 바

문제(問題) 좋아하 행동(行動) 풍

란던(LONDON) 엘

가리너(GARDINER) 가리너(GARDINER)

찰럿(CHARLOTTE) 점(點) 앉 앉 시작(始作)했다 이유(理由) 같

식사(食事) 풍봉(LONGBOURN) 꽝코 만한

오빠 바로 학실(雅實) 화상(狀況) 화 나아오 ~

로징스(ROSINGS) 점(點) 편(便) 대령(大領) 편

찰럿(CHARLOTTE) 살롱(梳妝台) 신사(紳士) 대령(大領) 편

찰럿(CHARLOTTE) 살롱(梳妝台) 신당(淸堂) 편(便) 대령(大領) 편

찰箩(CHARLOTTE) 살롱(梳妝台) 신당(淸堂) 편(便) 대령(大領) 편

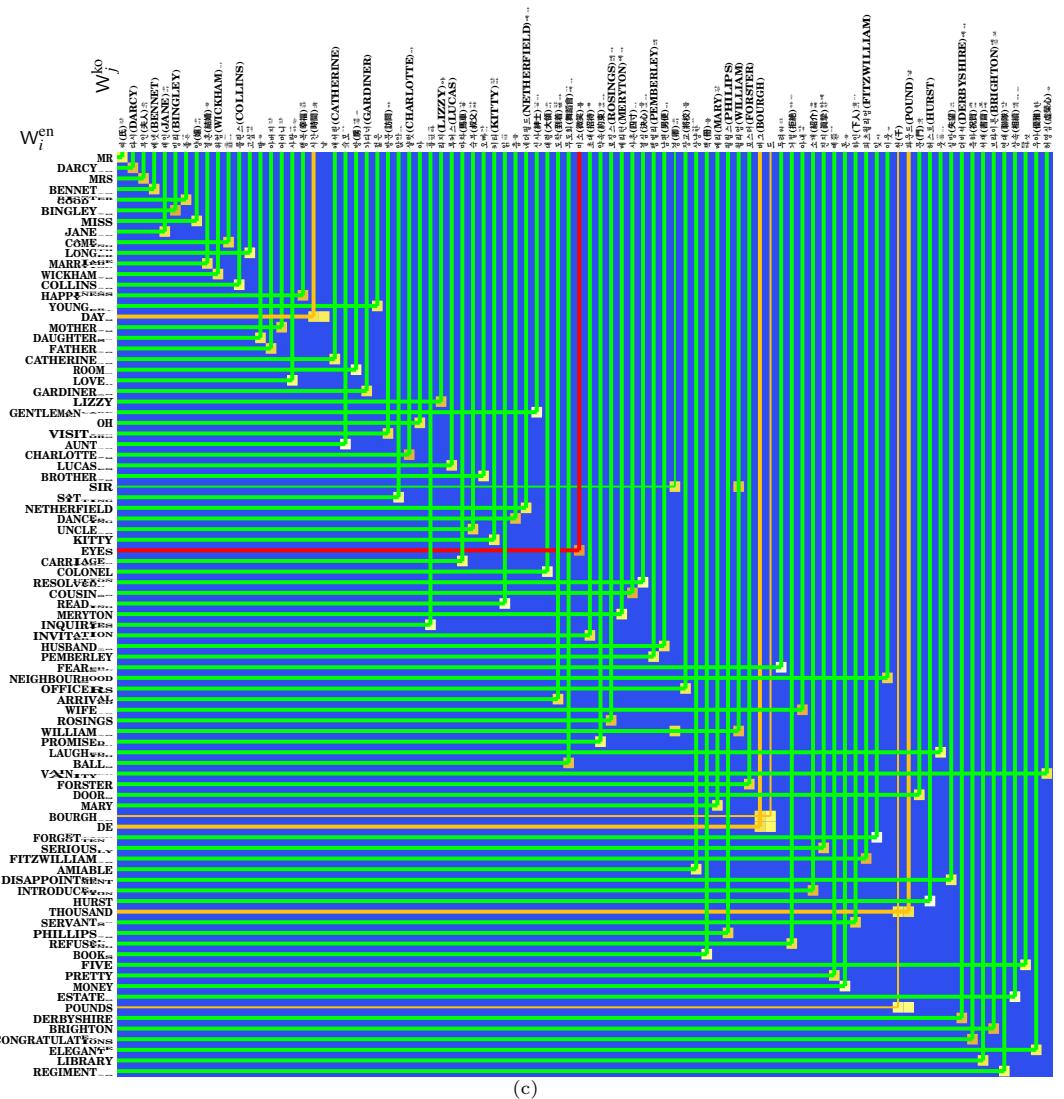


Fig. S15. Text mining in Korean. (Continued) (c) A control experiment in the spirit of Fig. S8c. The censorship  $s(\mathbf{v}_i^{\text{en}}, \mathbf{v}_j^{\text{ko}}) \geq 0.7$  is not imposed.

### 9.3 Modified Porter stemming algorithm for Turkish

Our Turkish stemming algorithm is adapted from the method of Gülsen Eryiğit and Eşref Adalı [51], which was previously implemented by Evren (Kapusuz) Çilden (<http://snowball.tartarus.org/algorithms/turkish/stemmer.html>).

Our algorithm works on Modern Turkish (after the Atatürk Reform of 1928), which uses (a modified version of) the Latin alphabet, and phases out most grammatical forms of Arabic and Persian origins. Since Atatürk's Language Reform, a significant amount of Arabic and Persian loanwords have become obsolete, even if they are apparently more concise (using fewer syllables) than their counterparts in Modern Turkish, such as

English gloss	via Proto-Turkic	via Arabic	via Persian
"friend"	<i>arkadaş</i>	—	<i>dost</i>
"war"	<i>savaş</i>	<i>harp</i>	<i>cenk</i>

While the native Turkish forms *arkadaş* and *savaş* contain more syllables than their Arabic/Persian equivalents, they are free from awkward consonant clusters.<sup>120</sup> Some Arabic/Persian loanwords still survive in Modern Turkish, as they are phonologically compatible with native Turkish words, such as *saat* "hour" (from Arabic) and *zaman* "time, epoch" (from either Arabic or Persian). It is also possible that a native Turkish word coexists with an Arabic/Persian loan for the same concept, with slightly different connotations, such as *ak* "white" (native Turkish) vs. *beyaz* "white" (via Arabic); *kara* "black" (native Turkish) vs. *siyah* "black" (via Persian). The last situation is similar to the coexistence of native Korean and Sino-Korean vocabularies in Modern Korean (see §9.2).

While we design our stemming algorithm for Modern Turkish, we take into account its phonological restrictions on syllable structure. Like what we did for Korean, we are not going to merge native and non-native etyma into the same word cluster.

Turkish stop words can carry various agglutinative suffixes. Instead of enumerating all of them in a single list, we will define several string patterns that match the criterion for stop words.

**Definition 9.16** (Turkish stop words). Define the pattern **TurkishStopRoot** as any one from the following list:

*aksi, altında, ama,anca,ancak,andan,alarında,arasında,artık,aşağı,aşağıda,aynı,ayı,ayırica,ayıryeten,ayıryeten,az,bana,başına,başka,bazen,bazı,bazılıları,bazılıarda,bazılılarına,bazılılarından,bazılılarını,bazılırinin,belgeyle,belki,ben,bende,benden,beni,benim,beri,bile,bir,biraz,biri,birlileri,birisi,birkaç,bırlikte,biz,bizde,bizden,bize,bizi,bizim,bizler,boylece,böyle,böylece,bu,buna,bunda,bundan,bunlar,bunlara,bunlarda,bunlardan,bunları,bunların,bunu,bunun,bununla,bura,burada,buradaki,buradan,buralar,buralara,buralarda,buralardan,buraları,buraların,buranın,buraya,burayı,bütün,civarda,çoğu,çoğunlukla,çok,çünkü,da,daha,dahası,dair,de,dek,dendenle,dışarı,dışında,diger,digerleri,diye,dolayı,dolayısıyla,durumda,e,eden,eder,ederim,ederiz,ederler,edersin,edersiniz,edildi,ediyor,ediyordu,ediyorduk,ediyordum,ediyordun,ediyordunuz,ediyorlar,ediyorlardı,ediyorsun,ediyorsunuz,ediyorum,ediyoruz,eğer,en,etmedi,etmedik,etmediler,etmedim,etmedin,etmediniz,etmem,etmeyiz,etmez,etmezler,etmezsin,etmezsiniz,etmiyor,etmiyordu,etmiyordum,etmiyordun,etmiyordunuz,etmiyorlar,etmiyordurdı,etmiyorsun,etmiyorsunuz,etmiyorum,etmiyorum,etrafında,etsen,etti,ettik,ettiler,ettim,ettin,ettiniz,ettirdi,ettirdik,ettirdiler,ettirdim,ettirdin,ettirdiniz,ettireceğim,ettireceğiz,ettirecek,ettirecekler,ettireceksin,ettireceksiniz,ettirir,ettiririm,ettiririz,ettirirler,ettirirsın,ettirirsınız,ettiriyor,ettiriyordu,ettiriyorduk,ettiriyordum,ettiriyordun,ettiriyordunuz,ettiriyorlar,ettiriyorlardı,ettiriyorsun,ettiriyorsunuz,ettiriyorum,ettiriyoruz,ettirmedı,ettirmedik,ettirmediler,ettirmedim,ettirmedin,ettirmediniz,ettirmem,ettirmeyeceğim,ettirmeyeceğiz,ettirmeyecek,ettirmeyecekler,ettirmeyeceksin,ettirmeyeceksiniz,ettirmeyiz,ettirmez,ettirmezler,ettirmezsin,ettirmezsiniz,ettirmiyor,ettirmiyordu,ettirmiyorduk,ettirmiyordular,ettirmiyordum,ettirmiyordun,ettirmiyordunuz,ettirmiyorlar,ettirmiyorsun,ettirmiyorsunuz,ettirmiyorum,ettirmiyoruz,evet,evvel,fakat,gayet,gelen,gelişme,gereçler,gerek,gerektiği,geride,gibi,göre,günü,hadi,hakkında,halâ,halbuki,halde,haline,hangi,hangileri,hangile-rine,hangilerini,hangilerinin,hangisi,hangisinde,hangisinden,hangisine,hangisini,hangisinin,hangsilerinde,hangsilerinden,hani,hariç,hatta,hayır,hem,hemen,henüz,hepsi,hepsinde,hepsine,hepsini,hepsinin,her,herhangi,herkes,herkesin,hiç,hiçbir,icar,icerede,icin,icinde,icinden,icine,ile,ilgili,irade,ise,ister,işte,iyi,izin,kadar,kâh,kapalı,karşı,karşın,karşısında,kazanılmış,kenara,keşke,kez,ki,kim,kimde,kimden,kime,kimi,kimin,kimisi,kimler,kimlerde,kimlerden,kimlere,kimleri,kimlerin,kimse,kişinin,kişkiye,konum,lâkin,madem,mademki,meğer,meğerki,meğerse,meselâ,mi,misin,misiniz,miydi,miyim,miyiz,mi,misin,misiniz,miyim,miyiz,mi,musun,musunuz,muydu,muyduk,muydum,muydun,muydunuz,muyum,mu-yuz,mü,nasıl,ne,nede,nedenle,neler,nelerde,nelerden,nelerle,neleri,nelerin,nere,nerede,nereden,neredeysse,nereler,nerelerde,nerelerden,nerelere,nereleri,nereye,nereyin,neye,neyi,neyin,neyse,niçin,o,ol,olacağım,olacağız,olacak,olacaklar,olacaksın,olacaksınız,olan,olanları,olarak,oldu,ol-duğu,olduk,oldukça,oldular,oldum,oldun,oldunuz,olma,olmadan,olmadı,olmadık,olmadıkça,olmadılar,*

<sup>120</sup>A typical syllable in a native Turkish word takes the form (C<sub>1</sub>)V(C<sub>2</sub>), containing a single vowel flanked by (optionally present) single consonants at both ends.

*olmadım, olmadın, olmadınız, olmak, olmalı, olmam, olmamalı, olmanın, olmayacağım, olmayacağız, olmayacak, olmayacaklar, olmayacaksın, olmayacaksınız, olmayız, olmaz, olmazlar, olmazsınız, olmazsınız, olmuştur, olmuyor, olmuyordu, olmuyordular, olmuyordum, olmuyordun, olmuyordunuz, olmuyollar, olmuyorsun, olmuyorsunuz, olmuyorum, olmuyoruz, olsa, olsun, olup, olur, olurken, olurlar, olursa, olursun, olursunuz, olurum, oluruz, oluyor, oluyordu, oluyordum, oluyordun, oluyordunuz, oluyorlar, oluyorlardı, oluyorsun, oluyorsunuz, oluyorum, oluyoruz, ona, onda, ondan, onlar, onlara, onlarda, onlardan, onları, onların, onu, onun, ora, orada, oradan, oralar, oralara, oralarda, oralardan, oraları, oraların, oranın, oraya, orayı, ordan, ortasında, oysa, ön, önce, onde, örneğin, ötesinde, ötürü, pek, rağmen, sadece, sahip, sahiptir, salt, sana, sayede, sebeple, sen, sende, senden, seni, senin, sık, sırasında, sırt, siz, sizde, sizden, size, sizi, sizin, son, sonra, sonradan, sonraki, sözgelimi, şayet, şekilde, şey, şimdiye, şu, şuna, şunda, şundan, şunlar, şunlara, şunlarda, şunlardan, şunları, şunların, şunu, şunun, şura, şurada, şuradan, şuralar, şuralara, şuralarda, şuralardan, şuraları, şuraların, şuranın, şuraya, şurayı, tabii, takdirde, tam, tarafından, tek, tekrar, uzak, üstelik, üzere, üzerinde, üzerine, vakte, var,vardı, ve, verdim, vesile, veya, veyahut, vs, ya, yahut, yakında, yaklaşık, yanında, yani, yapacağım, yapacağız, yapacak, yapacaklar, yapacaksin, yapacaksiniz, yapamadım, yapamaz, yapar, yaparım, yaparız, yaparlar, yaparsın, yaparsınız, yapıyor, yapıyordu, yapıyordum, yapıyordun, yapıyordunuz, yapıyorlar, yapıyordular, yapıyorsun, yapıyorsunuz, yapıyorum, yapıyorum, yapıyorum, yapma, yapmadı, yapmadık, yapmadılar, yapmadım, yapmadın, yapmadınız, yapmam, yapmayacağım, yapmayacağız, yapmayacak, yapmayacaklar, yapmayacaksin, yapmayacaksiniz, yapmayız, yapmaz, yapmazlar, yapmazsin, yapmazsiniz, yapmiyor, yapmiyordu, yapmiyorduk, yapmiyordular, yapmiyordum, yapmiyordun, yapmiyordunuz, yapmiyorlar, yapmiyorsun, yapmiyorsunuz, yapmiyorum, yapmiyoruz, yaptı, yaptık, yaptılar, yaptım, yaptın, yaptınız, yerine, yeterince, yine, yoksa, yoluyla, yukarı, yukarıdaki, zaman, zamanda, zaten, zira.*

If a word exactly matches the string pattern **TurkishStopRoot**( $\emptyset|ki$ ) or any one of the following:<sup>121</sup>

*(b(en|iz)im|on(lar|u)n|s(en|iz)in)( $\emptyset|ki(\emptyset|n)$ )( $\emptyset|de(\emptyset|n)|e|i(\emptyset|n)$ ),  
(değil|edebil|öyle|tüm|yok)X,  
edece(ğ|k)X,  
etm(e|is)X,  
kendi(leri|(m|n)( $\emptyset|iz$ )|si)( $\emptyset|n$ )( $\emptyset|de(\emptyset|n)|e|i(\emptyset|n)$ ),  
ol(abil|ma|u(n|r|y))X,  
ol(aca|du)(ğ|k)X,*

then we consider it a Turkish stop word. All the Turkish stop words that appear in a particular document need to be ignored before we perform word clustering on the rest of the vocabulary list.  $\square$

### 9.3.1 Effective spelling and essential root

**Definition 9.17** (Turkish protected range). Set  $\mathbf{V} = (a|\hat{a}|e|i|\hat{i}|o|\hat{o}|u|\hat{u}|i\ddot{u})$ . If  $\hat{\sigma}$  does not contain  $\mathbf{V}$ , define  $\text{ProtRg}(\hat{\sigma}) = 0$ ; otherwise, define  $\text{ProtRg}(\hat{\sigma})$  as the first position occupied by a letter in  $\mathbf{V}$ .

**Algorithm 9.18** (Turkish effective spelling). Set  $\mathbf{V} = (a|\hat{a}|e|i|\hat{i}|o|\hat{o}|u|\hat{u}|i\ddot{u})$ ,  $\mathbf{C} = \bar{\mathbf{V}}$ . For a Turkish word  $\hat{\sigma}$ , its effective spelling  $\text{EffSpell}(\hat{\sigma})$  is constructed in the following steps:

- (1) Do  $I \rightarrow i$ ,  $\hat{I} \rightarrow i$ , before converting to lowercase.
- (2) Look for the pattern  $\mathbf{CV}(m|p|r|s)\mathbf{CV}\sim$ . If the two occurrences of  $\mathbf{CV}$  represent the same consonant-vowel combination, reduce the aforementioned pattern to just  $\mathbf{CV}$ ; otherwise, leave as is.<sup>122</sup>
- (3) Do  $ad\sim \rightarrow na\mu$ ,  $apaci\sim \rightarrow acı$ ,  $emel\sim \rightarrow eμeλ$ ,  $emin\mathbf{X}\sim \rightarrow σuρe$ ,  $ipislak\sim \rightarrow islak$ ,  $uzun\mathbf{X}\sim \rightarrow λoñy$ ,  $mr \rightarrow ηeρρ$ ,  $mrs \rightarrow φρau$ .
- (4) Call the result so far as  $\hat{\sigma}'$ , and break it down into  $\hat{\sigma}' = \hat{\sigma}_1\hat{\sigma}_2$ , where  $\ell(\hat{\sigma}_1) = \text{ProtRg}(\hat{\sigma}')$ . Work on  $\hat{\sigma}_2$  as follows:
  - (4.1) Do  $\sim(\emptyset|u)ms(i|i|u) \rightarrow \emptyset$ ,  $\sim(\emptyset|u)mt(i|i)rak \rightarrow \emptyset$ ,  $\sim(\emptyset|y)(i|i|u|i\ddot{u})p \rightarrow \emptyset$ ,  $\sim(c|ç)(a|e) \rightarrow \emptyset$ ,  $\sim(c|ç)(i|i)k \rightarrow \emptyset$ ;
  - (4.2) Call the result so far as  $\hat{\sigma}'_2$ . If  $\hat{\sigma}'_2$  contains  $s(i|i|u|i\ddot{u})z$ , then define  $\hat{\sigma}''_2 = \tilde{v}\hat{\sigma}^*_2$ , where  $\hat{\sigma}^*_2$  results from doing  $s(i|i|u|i\ddot{u})z\mathbf{X} \rightarrow \emptyset$  on  $\hat{\sigma}'_2$ ; otherwise define  $\hat{\sigma}''_2 = \hat{\sigma}'_2$ .

Concatenate  $\hat{\sigma}_1$  and  $\hat{\sigma}''_2$ .

<sup>121</sup>The Turkish letter  $\mathfrak{s}$  (LATIN SMALL LETTER S WITH CEDILLA) should not be confused with the Romanian letter  $\mathfrak{s}$  (LATIN SMALL LETTER S WITH COMMA BELOW), as they occupy different positions in Unicode.

<sup>122</sup>This step treats the intensified form of certain Turkish adjectives, such as *bembeyaz* “totally white” (cf. *beyaz* “white”) and *bomboş* “totally empty” (cf. *bos* “empty”.)

(5) Do  $(amca|dayı)\mathbf{X} \sim \rightarrow \omega\tilde{\nu}\kappa\lambda e$ ,  $(bibi|eme|hala|teyze)\mathbf{X} \sim \rightarrow \alpha u\tilde{\nu}\tau$ .<sup>123</sup>

(6) Replace

$anla\mathbf{X} \sim$	$der(\emptyset iz) \sim$	$di(\emptyset me)y(ece or)\mathbf{X} \sim$	$dişün\mathbf{X} \sim$	$eleştir\mathbf{X} \sim$	$getir\mathbf{X} \sim$	$gid(e i) \sim$	$götür\mathbf{X} \sim$	$gül \sim$
$\kappa\mu\pi\rho$	$de$	$de$	$\theta\tilde{i}\tilde{\nu}\kappa$	$\kappa\mu\iota\tau r$	$\beta\mu\iota\gamma$	$gitmek$	$\kappa\mu\mu\mu\mu$	$\gamma\mu\mu\mu\mu$
$hanım \sim$	$hatır\mathbf{X} \sim$	$his \sim$	$ileri\mathbf{X} \sim$	$itiraf\mathbf{X} \sim$	$itiraz\mathbf{X} \sim$	$izl\mathbf{X} \sim$	$kismet\mathbf{X} \sim$	$kızar \sim$
$\lambda\mu\delta$	$\mu\tilde{i}\tilde{\nu}\delta$	$hisset$	$\varphi\mu\o\tilde{\nu}\tau$	$\kappa\o\tilde{\nu}\sigma$	$\omega\beta\mu\kappa\tau$	$\varphi\o\lambda\omega$	$\varphi\mu\tau\epsilon$	$\kappa\mu\zeta\mu$
$su(lar yu)\mathbf{X} \sim$	$ulaş\mathbf{X} \sim$	$yemin\mathbf{X} \sim$	$yer(\emptyset iz) \sim$	$yi(\emptyset me)y(ece or)\mathbf{X} \sim$	$anda$	$su(\emptyset da dan dur sun sunuz uz ya)$		
$hho$	$peach$	$vow$	$ye$	$ye$	$now$			$hho$

(7) Call the result so far as  $\hat{\sigma}^\dagger$ , and break it down into  $\hat{\sigma}^\dagger = \hat{\sigma}_1^\dagger \hat{\sigma}_2^\dagger$ , where  $\ell(\hat{\sigma}_1^\dagger) = \text{ProtRg}(\hat{\sigma}^\dagger)$ . Work on  $\hat{\sigma}_2^\dagger$  as follows:

(7.1) Replace

$\sharp\mathbf{X}$	$(abil ebil mak mek)\mathbf{X}$	$(aca ece)(\check{g} k)\mathbf{X}$	$c(a e)s(i i)n(a e)\mathbf{X}$	$l(a e)(r s)\mathbf{X}$
$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
$m(i i u ü)\$X$	$r \sim$	$yor\mathbf{X}$	$\sim(\emptyset n)(d t)(a e)(\emptyset n)$	$Ta$
$\emptyset$	$\rho$	$\emptyset$		
$\sim(d t)(i i u ü)(\emptyset k l(a e)r) m n(\emptyset (i i u ü)z) r\mathbf{X}$	$Tu$	$\sim ki(\emptyset n)(de den e i in)$	$\sim s(i i u ü)n(\emptyset (i i u ü)z)$	$\emptyset$

(7.2) Do  $(d|t)(i|i|u|ü)(\check{g}|k)\mathbf{X} \rightarrow \emptyset$ ,  $Vym(i|i|u|ü)\$X \rightarrow V$ ,  $m(a|e|i|i|u|ü)\mathbf{X} \rightarrow \emptyset$ ,  $\sim(i|i|u|ü)(m|n)(\emptyset|(i|i|u|ü)z) \rightarrow \emptyset$ ,  $\sim V(y(\emptyset|l)(a|e) \rightarrow V$ ,  $\sim V(y(\emptyset|Tu) \rightarrow V$ ,  $\sim V(y(i|i|u|ü)(\emptyset|m|z) \rightarrow V$ ,  $\sim Vyken \rightarrow V$ ,  $\sim Vys(a|e) \rightarrow V$ .

(7.3) Do  $\sim(\emptyset|i|i|u|ü)(m|n) \rightarrow \emptyset$ .

(7.4) Call the result so far as  $\hat{\sigma}_2^\ddagger$ . Concatenate  $\hat{\sigma}_1^\dagger$  and  $\hat{\sigma}_2^\ddagger$ .

(8) Do  $T\mathbf{X} \rightarrow \emptyset$ ,  $\sim r(\emptyset|(i|i|u|ü)z) \rightarrow \emptyset$ .

**Algorithm 9.19** (Turkish essential root). Let  $\hat{\sigma}$  be the effective spelling of a Turkish word, then its corresponding essential root  $\text{EssRoot}(\hat{\sigma})$  is constructed in the following steps:

- (1) Search for the pattern  $\mathbf{C}_{m_0} \mathbf{V} \mathbf{C}_{m_0} \mathbf{V}_{m_0} (\emptyset|\overline{(n|s)}) \sim$ . If this search is successful, define  $\hat{\sigma}^\flat$  as the result from doing  $\sim(i|i|u|ü) \rightarrow \emptyset$  on such a pattern; otherwise, define  $\hat{\sigma}^\flat$  as  $\hat{\sigma}$ .
- (2) On  $\hat{\sigma}^\flat$ , do  $\sim b \rightarrow p$ ,  $\sim c \rightarrow \zeta$ ,  $\sim d \rightarrow t$ ,  $\sim \check{g} \rightarrow k$ .

### 9.3.2 Approximate clustering

There will not be any mutation rules in our algorithm for Turkish, so  $\text{SimpHrdTest} = \text{HrdTest}$ .

**Algorithm 9.20** (Simple heredity test). Suppose that  $\hat{\alpha}$  and  $\hat{\beta}$  are both lowercase strings. The Boolean-valued function  $\text{SimpHrdTest}(\hat{\alpha}, \hat{\beta})$  returns **TRUE** if  $\hat{\alpha}$  contains at least one instance of  $\mathbf{V}$ , AND at least one of the following two conditions holds:<sup>124</sup>

- (i)  $\hat{\alpha} = \hat{\beta}$ ;
- (ii)  $\hat{\alpha}' = \hat{\beta}'$ , where the strings with prime result from doing  $\sim^{\mathbf{X}^{\infty}}(\mathbf{V}\hat{\chi})(\emptyset|l)\mathbf{V} \rightarrow \mathbf{X}$  on their counterparts without prime.

**Algorithm 9.21** (Approximate clustering of Turkish words). The approximate clustering of a list of Turkish words  $\{\hat{\alpha}_1, \dots, \hat{\alpha}_N\}$  is completed in two stages:

- (1) We sort the list  $\{(\hat{\alpha}_1, \text{EffSpell}(\hat{\alpha}_1)), \dots, (\hat{\alpha}_N, \text{EffSpell}(\hat{\alpha}_N))\}$  alphabetically according to the last component. If two consecutive neighbors in the alphabetized list  $\{(\hat{\alpha}_{(1)}, \text{EffSpell}(\hat{\alpha}_{(1)})), \dots, (\hat{\alpha}_{(N)}, \text{EffSpell}(\hat{\alpha}_{(N)}))\}$  satisfy  $\text{HrdTest}(\text{EffSpell}(\hat{\alpha}_{(n)}), \text{EffSpell}(\hat{\alpha}_{(n+1)})) = \text{FALSE}$ , where  $n \in \mathbb{Z} \cap [1, N]$ , then add a demarcation line between these two entries. In this way, the alphabetized list is divided into separate groups of words tagged with their effective spellings:  $\{g_1 = \{(\hat{\alpha}_{(1,1)}, \text{EffSpell}(\hat{\alpha}_{(1,1)})), \dots, (\hat{\alpha}_{(1,n_1)}, *)\}, \dots, g_M = \{(\hat{\alpha}_{(M,1)}, *), \dots, (\hat{\alpha}_{(M,n_M)}, *)\}\}$ , where each sublist also preserves alphabetic order. In particular, we have  $\hat{\alpha}_{(1,1)} = \hat{\alpha}_{(1)}$  and  $\hat{\alpha}_{(M,n_M)} = \hat{\alpha}_{(N)}$ .

<sup>123</sup>In this step, we group certain kinship terms in Turkish so as to match their English counterparts. However, since the Turkish word *kardeş* means both “brother” and “sister”, we cannot do anything to compensate for this.

<sup>124</sup>As a general rule in this document, the truth values of items labeled with Roman numerals are connected to each other with logical OR.

- (2) For each group of words  $g_m = \{(\hat{\alpha}_{(m,1)}, *), \dots, (\hat{\alpha}_{(m,n_m)}, *)\}$  where  $m \in \mathbb{Z} \cap [1, M]$ , we augment it into a tagged entry  $G_m = (g_m, \text{EffSpell}(\hat{\alpha}_{(m,1)}), \text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)})))$ . The list  $\{G_1, \dots, G_M\}$  is sorted alphabetically, with respect to  $\text{EssRoot}(\text{EffSpell}(\hat{\alpha}_{(m,1)}))$  (with higher priority), and  $\text{EffSpell}(\hat{\alpha}_{(m,1)})$  (with lower priority). If two consecutive neighbors in the alphabetized list  $\{G_{(1)} = (g_{(1)}, \hat{\gamma}'_{(1)}, \hat{\gamma}''_{(1)}), \dots, G_{(M)} = (g_{(M)}, \hat{\gamma}'_{(M)}, \hat{\gamma}''_{(M)})\}$  satisfy

$$\text{HrdTest}(\hat{\gamma}''_{(m)}, \hat{\gamma}''_{(m+1)}) = \text{FALSE}$$

AND

$$\text{HrdTest}(\hat{\gamma}'_{(m)}, \hat{\gamma}'_{(m+1)}) = \text{FALSE}$$

where  $m \in \mathbb{Z} \cap [1, M]$ , then add a demarcation line between these two entries. In this manner, the alphabetized list  $\{G_{(1)}, \dots, G_{(M)}\}$  is divided into separate groups of tagged clusters  $\{\Gamma_1 = \{G_{(1,1)}, \dots, G_{(1,m_1)}\}, \dots, \Gamma_K = \{G_{(K,1)}, \dots, G_{(K,m_K)}\}\}$ . Finally, the output list of word clusters  $\{\check{\Gamma}_1, \dots, \check{\Gamma}_K\}$  is constructed by discarding all the tags (effective spellings, essential roots) from each  $\Gamma_k$ , where  $k \in \mathbb{Z} \cap [1, K]$ .

*Example 9.21.1.* Our clustering algorithm correctly groups the following families of Turkish words:

ev, evde, evden, eve, evi, evin, evler, evlerde, evlerden, evlere, evleri, evlerin — “house”;  
 dedi, dedik, dediler, dedim, dedin, dediniz, demedi, demedik, demediler, demedim, demedeniz, demem, demeyeceğim, demeyeceğiz, demeyecek, demeyecekler, demeyeceksin, demeyeceksiniz, demeyiz, demez, demezler, demezsin, demezsiniz, demiyor, demiyordu, demiyorduk, demiyordular, demiyordum, demiyordun, demiyordunuz, demiyorlar, demiyorsun, demiyorsunuz, demiyorum, demiyoruz, der, derim, deriz, derler, dersin, dersiniz, diyecek, diyecekler, diyeceğim, diyeceğiz, diyecek, diyecekler, diyeceksin, diyeceksiniz, diyor, diyordu, diyorduk, diyordum, diyordun, diyordunuz, diyorlar, diyorsun, diyorsunuz, diyorum, diyorus — “say”;  
 al, alacağım, alacağız, alacak, alacaklar, alacaksın, alacaksınız, aldı, aldık, aldılar, aldım, aldın, aldınız, alır, alırırm, alırız, alırlar, alırsın, alırsınız, alıyor, aliyordu, aliyorduk, aliyordum, aliyordun, aliyordunuz, aliyorlar, aliyorlardı, aliyorsun, aliyorsunuz, aliyorum, aliyoruz, almadi, almadık, almadılar, almadım, almadın, almadınız, almam, almayacağım, almayacağız, almayacak, almayacaklar, almayacaksın, almayacaksınız, almayız, almaz, almazlar, almazsin, almazsiniz, almiyor, almiyordu, almiyorduk, almiyordular, almiyordum, almiyordun, almiyordunuz, almiyorlar, almiyorsun, almiyorsunuz, almiyorum, almiyoruz — “buy”;  
 okudu, okuduk, okudular, okudum, okudunuz, okumadı, okumadılar, okumadım, okumadın, okumadınız, okumam, okumayacağım, okumayacağız, okumayacak, okumayacaklar, okumayacaksın, okumayaçaksınız, okumayız, okumaz, okumazlar, okumazsin, okumazsınız, okumuyor, okumuyordu, okumuyorduk, okumuyordular, okumuyordum, okumuyordun, okumuyordunuz, okumuyorlar, okumuyorsun, okumuyorsunuz, okumuyorum, okumuyoruz, okur, okurlar, okursun, okursunuz, okurum, okuruz, okuyacağım, okuyacağız, okuyacak, okuyacaklar, okuyacaksın, okuyacaksınız, okuyor, okuyordu, okuyorduk, okuyordum, okuyordun, okuyordunuz, okuyorlar, okuyorlardı, okuyorsun, okuyorsunuz, okuyorum, okuyoruz — “read”;  
 kitaba, kitabı, kitabin, kitap, kitaplar, kitaplara, kitaplarda, kitaplardan, kitapları, kitapların, kitapta, kitaptan — “book”.

*Example 9.21.2.* In Fig. S16, we further apply the aforementioned word clustering algorithm to topic extraction and machine translation (see Table S1 for text source).

Note that the Turkish translation for *disappointment* consists of two words *hayal kırıklığı* (literally “dream brokenness”), where *hayal* means “dream”. This explains why there are two hot spots in the row for *disappointment* in Fig. S16b.

In Turkish, *de* “also, too” is a stop word.

ELIZABETH - MR Darcy - NİN DÜŞÜNCELERİ - GELİŞİMİ - DEDİ  
GÖRME - VERİ - BINGLEY - İN SÖYLEŞE! KIZ - GIT - MRS BENNET -  
ANLAT - JANE - İN KALOĞ - ADAM'IN İKİ - MISS KARDEŞİ - BÜYÜK SEYİ  
HİSLER - WICKHAM - İN KONUSMAK - LYDIA - NİN DUYURU - CIKMAK - GENÇ ARKADAŞI - ETİĞİ - GERÇEKTE ON  
BİLMİYOR - KENDİ BULŞESEN - ANNESİN - KONUŞMAK - BABASINI - GÜZEL - AİLESİN - GEREKLİYİ - SEBEPI - DAVRANIŞLARINI  
SÖZLER - İSTERİĞİN - İN - DUYURU - İN - DUYURU - İN - DUYURU - İN - DUYURU - İN - DUYURU - İN - DUYURU - İN - DUYURU - İN  
İNANCAK - GÜN - HANIMLAR - GETİRİ - ODASINA - EMİN - ARA - LADY MEKTUBU - MESALET - HUSUSUNA  
DOĞRU - BOYLESİ - İÇERİSİNDE - ANA - FAZLA - EVLENME - OLAMAMI - SAMİMİ - KABUL YAKIN - CATHERINE - İN  
BİRİST - GÖRÜNCÜĞÜ - CEVAP - YANPAN - YAZ - GEÇT - ÜMIT - AH - BAHSİSE - BEKLƏNME - YASAM - EDİLEBİYİ  
İNSAN - BASLAŞTI - GIR - ABLAŞISI - SAAP - DÖNMESİ - BİLİYOR - CHARLOTTE - ON KARAR - GARDINER - İLK  
TANTAN - YÜZÜNDEN - KÜÇÜK AKLIŞ - SEVGİLİ KAHNE - HOS - EVS - YENİDEN YANIT - LADI CALIŞTA - GOSTERR - İN - UZUN BIRAKYI - HATIRI  
MEMNUN AYRILMASI - ILGI - EVLİLİK - OTE - LIZZY - SORSE - TEŞEKKUR - ANDA OGRENME - ZİYARET - KIZI - BAKTI - KARAKTERİNE - AKSAM

(a)

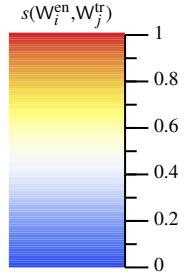


Fig. S16. Text mining in Turkish. (a) Statistically identified topics ( $n_{ii} \geq 20$ ) in a Turkish version of *Pride and Prejudice*, with the same color encoding scheme as Fig. S3. (b) Semantic similarities  $s(W_i^{\text{en}}, W_j^{\text{tr}})$  between selected topics in English and Turkish versions of *Pride and Prejudice*. Cross-hairs meet at optimal nodes that solve the bipartite matching problem. The thickness of each horizontal (resp. vertical) cross-hair is inversely proportional to the row-wise (resp. column-wise) rank of the similarity score for the optimal node. Green (resp. amber) cross-hair indicates an exact (resp. a close but non-exact) match.

- [1] W. E and Y. Zhou, “A mathematical model for universal semantics,” 2020, arXiv:1907.12293v6 [cs.CL].
- [2] W. A. Gale, K. W. Church, and D. Yarowsky, “A method for disambiguating word senses in a large corpus,” *Comput. Humanities*, vol. 26, no. 5-6, pp. 415–439, 1992.
- [3] D. M. Blei, A. Y. Ng, and M. I. Jordan, “Latent Dirichlet allocation,” *J. Mach. Learn. Res.*, vol. 3, no. Jan, pp. 993–1022, 2003.
- [4] P. D. Turney and P. Pantel, “From frequency to meaning: Vector space models of semantics,” *J. Artificial Intelligence Res.*, vol. 37, pp. 141–188, Feb. 2010.
- [5] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, “Distributed representations of words and phrases and their compositionality,” in *Advances in Neural Information Processing Systems 26*. La Jolla, CA: NIPS, 2013, pp. 3111–3119.
- [6] S. Arora, Y. Li, Y. Liang, T. Ma, and A. Risteski, “A latent variable model approach to PMI-based word embeddings,” *Transactions of the Association for Computational Linguistics*, vol. 4, pp. 385–399, 2016.
- [7] Y. Zhou and X. Zhuang, “Robust reconstruction of the rate constant distribution using the phase function method,” *Biophys. J.*, vol. 91, no. 11, pp. 4045–4053, 2006.
- [8] N. Haydn, Y. Lacroix, and S. Vaienti, “Hitting and return times in ergodic dynamical systems,” *Ann. Probab.*, vol. 33, pp. 2043–2050, 2005.
- [9] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. Hoboken, NJ: Wiley Interscience, 2006.
- [10] M. Pollicott and M. Yuri, *Dynamical Systems and Ergodic Theory*, ser. London Mathematical Society Student Texts. Cambridge, UK: Cambridge University Press, 1998, vol. 40.
- [11] M. Ružička, “Anwendung mathematisch-statistischer Methoden in der Geobotanik (Synthetische Bearbeitung von Aufnahmen),” *Biológia (Bratislava)*, vol. 13, pp. 647–661, 1958.
- [12] E. Deza and M.-M. Deza, *Dictionary of distances*. Amsterdam, The Netherlands: Elsevier, 2006.
- [13] G. Gilbert, “Distance between sets,” *Nature*, vol. 239, no. 5368, pp. 174–174, 1972.
- [14] F. McSherry and M. Najork, “Computing information retrieval performance measures efficiently in the presence of tied scores,” in *Advances in Information Retrieval. ECIR 2008*, C. Macdonald, I. Ounis, V. Plachouras, I. Ruthven, and R. W. White, Eds. Berlin, Germany: Springer, 2008, pp. 414–421.
- [15] S. Brin and L. Page, “The anatomy of a large-scale hypertextual web search engine,” *Comput. Networks ISDN*, vol. 30, no. 1-7, pp. 107–117, 1998.
- [16] M. Bressan and E. Peserico, “Choose the damping, choose the ranking?” *Journal of Discrete Algorithms*, vol. 8, no. 2, pp. 199–213, 2010.
- [17] D. Lardiere, “Words and their parts,” in *An Introduction to Language and Linguistics*, R. Fasold and J. Connor-Linton, Eds. Cambridge, UK: Cambridge University Press, 2006, ch. 2, pp. 55–96.
- [18] A. D. Friederici, “The neurobiology of language comprehension,” in *Language Comprehension: A Biological Perspective*, A. D. Friederici, Ed. Berlin, Germany: Springer, 1999, ch. 9, pp. 265–304.
- [19] R. Jakobson, “The identification of phonemic entities,” in *Recherches structurales 1949. Interventions dans le débat glossématique*, ser. Travaux du Cercle Linguistique de Copenhague. Copenhagen, Denmark: Nordisk Sprog- og Kulturforslag, 1949, vol. V, pp. 205–213.
- [20] C. Everett, D. E. Blasí, and S. G. Roberts, “Climate, vocal folds, and tonal languages: Connecting the physiological and geographic dots,” *Proc. Natl. Acad. Sci. USA*, vol. 112, pp. 1322–1327, 2015.
- [21] C. R. Ember, “Commentary: Considering language as expressive culture,” *Journal of Language Evolution*, vol. 1, no. 1, pp. 60–61, 2016.
- [22] C. Everett, “Languages in drier climates use fewer vowels,” *Frontiers in Psychology*, vol. 8, p. Article 1285, 2017.

- [23] J. G. Frazer, *The Golden Bough: a study of magic and religion*, abridged ed. New York, NY: MacMillan, 1922.
- [24] W. Labov, *Principles of Linguistic Change: Social Factors*. Malden, MA: Blackwell Publishers, 2001.
- [25] ——, *Principles of Linguistic Change: Cognitive and Cultural Factors*. Malden, MA: Wiley-Blackwell, 2010.
- [26] K. Allan and K. Burridge, *Forbidden Words: Taboo and the Censoring of Language*. Cambridge, UK: Cambridge University Press, 2006.
- [27] A. D. Friederici, J. Bahlmann, S. Heim, R. I. Schubotz, and A. Anwander, “The brain differentiates human and non-human grammars: Functional localization and structural connectivity,” *Proc. Natl. Acad. Sci. USA*, vol. 103, no. 7, pp. 2458–2463, 2006.
- [28] S. M. Wilson, S. Galantucci, M. C. Tartaglia, K. Rising, D. K. Patterson, M. L. Henry, J. M. Ogar, J. DeLeon, B. L. Miller, and M. L. Gorno-Tempini, “Syntactic processing depends on dorsal language tracts,” *Neuron*, vol. 72, no. 2, pp. 397–403, 2011.
- [29] P. F. Brown, P. V. deSouza, R. L. Mercer, V. J. Della Pietra, and J. C. Lai, “Class-based  $n$ -gram models of natural language,” *Computat. Linguit.*, vol. 18, no. 4, pp. 467–479, 1992.
- [30] S. Pinker, *The Stuff of Thought: Language as a Window into Human Nature*. New York: Viking, 2007.
- [31] ——, *Learnability and Cognition: the Acquisition of Argument Structure*. Cambridge, MA: The MIT Press, 2013.
- [32] N. Chomsky, *The Minimalist Program*. Cambridge, MA: MIT Press, 2015.
- [33] M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul, “An introduction to variational methods for graphical models,” in *Learning in graphical models*, M. I. Jordan, Ed. Cambridge, MA: MIT Press, 1999, pp. 105–161.
- [34] D. Heckerman, “A tutorial on learning with Bayesian networks,” in *Learning in graphical models*, M. I. Jordan, Ed. Cambridge, MA: MIT Press, 1999, pp. 301–354.
- [35] G. Carlsson, “Topology and data,” *Bull. Amer. Math. Soc. (N.S.)*, vol. 46, no. 2, pp. 255–308, 2009.
- [36] W. E, J. Lu, and Y. Yao, “The landscape of complex networks—critical nodes and a hierarchical decomposition,” *Methods Appl. Anal.*, vol. 20, no. 4, pp. 383–404, 2013.
- [37] S. B. Needleman and C. D. Wunsch, “A general method applicable to the search for similarities in the amino acid sequence of two proteins,” *Journal of Molecular Biology*, vol. 48, no. 3, pp. 443–453, 1970.
- [38] T. F. Smith and M. S. Waterman, “Identification of common molecular subsequences,” *Journal of Molecular Biology*, vol. 147, no. 1, pp. 195–197, 1981.
- [39] M. F. Porter, “An algorithm for suffix stripping,” *Program*, vol. 14, no. 3, pp. 130–137, 1980.
- [40] ——, “Snowball — a small string processing language designed for creating stemming algorithms for use in Information Retrieval,” <http://snowball.tartarus.org/>, 2014.
- [41] S. Pinker, “Words and rules in the human brain,” *Nature*, vol. 387, no. 6633, pp. 547–548, 1997.
- [42] W. D. Marslen-Wilson and L. K. Tyler, “Dissociating types of mental computation,” *Nature*, vol. 387, no. 6633, pp. 592–594, 1997.
- [43] E. Lieberman, J.-B. Michel, J. Jackson, T. Tang, and M. A. Nowak, “Quantifying the evolutionary dynamics of language,” *Nature*, vol. 449, no. 7163, pp. 713–716, 2007.
- [44] Y. Yang, W.-t. Yih, and C. Meek, “WikiQA: A challenge dataset for open-domain question answering,” in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Lisbon, Portugal: Association for Computational Linguistics, 2015.
- [45] J. Yeon and L. Brown, *Korean: a comprehensive grammar*. London, UK: Routledge, 2011.
- [46] T. Lundskær-Nielsen and P. Holmes, *Danish: An Essential Grammar*, 2nd ed. London, UK: Routledge, 2000.
- [47] B. C. Donaldson, *Dutch: A Comprehensive Grammar*, 2nd ed. London, UK: Routledge, 2008.
- [48] C. Rounds, *Hungarian: An Essential Grammar*. London, UK: Routledge, 2001.

- [49] S. Starostin, A. Dybo, and O. Mudrak, *Etymological dictionary of the Altaic languages*. Leiden, The Netherlands: Brill, 2003.
- [50] J. I. Hualde and J. Ortiz de Urbina, Eds., *A grammar of Basque*. Berlin, Germany: Mouton de Gruyter, 2003.
- [51] G. Eryiğit and E. Adalı, “An affix stripping morphological analyzer for Turkish,” in *Proceedings of the IASTED International Conference ARTIFICIAL INTELLIGENCE AND APPLICATIONS*, Innsbruck, Austria, 2004.