

Manual for Morphological Annotation  
Revision for  
Prague Dependency Treebank – Consolidated  
2020 release

Marie Mikulová  
Jiří Hana  
Jan Hajič  
Hana Hanová  
Barbora Hladká  
Jaroslava Hlaváčová  
Emil Jeřábek  
Barbora Štěpánková  
Daniel Zeman

December 2020

# Contents

<b>1 Preface</b>	<b>4</b>
<b>2 Introduction</b>	<b>5</b>
<b>3 Morphological Dictionary</b>	<b>6</b>
3.1 Principle of unique analysis . . . . .	6
3.2 Principle of morphological differentiation . . . . .	7
3.3 Principle of unique paradigm . . . . .	9
<b>4 Lemma Structure</b>	<b>11</b>
4.1 Lemma proper . . . . .	12
4.1.1 Lemma number . . . . .	12
4.2 Additional information about the paradigm (AddInfo) . . . . .	13
4.2.1 Reference . . . . .	13
4.2.2 Name label . . . . .	13
4.2.3 Style label . . . . .	14
4.2.4 Variant info . . . . .	14
4.2.5 Derivation info . . . . .	15
4.2.6 Explanational comment . . . . .	15
<b>5 Tag Structure</b>	<b>16</b>
5.1 Part of speech (1st position) . . . . .	16
5.2 Detailed part of speech (2nd position) . . . . .	17
5.3 Gender (3rd position) . . . . .	22
5.4 Number (4th position) . . . . .	22
5.5 Case (5th position) . . . . .	23
5.6 Possessor's gender (6th position) . . . . .	24
5.7 Possessor's number (7th position) . . . . .	24
5.8 Person (8th position) . . . . .	25
5.9 Tense (9th position) . . . . .	25
5.10 Degree of comparison (10th position) . . . . .	25
5.11 Negation (11th position) . . . . .	26
5.12 Voice (12th position) . . . . .	26
5.13 Verbal aspect (13th position) . . . . .	26
5.14 Aggregate (14th position) . . . . .	27
5.15 Variant, style, abbreviation (15th position) . . . . .	27
<b>6 Stylistic Characteristics</b>	<b>29</b>
6.1 Style labeling . . . . .	29
<b>7 Derivative Relations</b>	<b>31</b>
7.1 Automatically derived lemmas . . . . .	32
7.2 Derivative relation types . . . . .	33
<b>8 Semantic Description</b>	<b>35</b>
<b>9 Orthographic and Stylistic Variants</b>	<b>36</b>
9.1 Full-paradigm variants . . . . .	36
9.1.1 Types of full-paradigm variants . . . . .	37
9.2 Wordform variants . . . . .	38
9.3 Boderline cases of full-paradigm and wordform variants . . . . .	38
9.3.1 Variation is not in the full paradigm . . . . .	38
9.3.2 Variation in the base form . . . . .	40

<b>10 Part of Speech Determination (problematic cases)</b>	<b>42</b>
10.1 Part of speech of inflexible words . . . . .	42
10.1.1 Frozen wordforms ( <i>krážem, bycha, domácku</i> ) . . . . .	42
10.2 Nouns from adjectives . . . . .	43
10.3 Part of speech of predicatives (words with suffix <i>-o</i> ) . . . . .	43
<b>11 Detailed Part of Speech</b>	<b>45</b>
11.1 Subtypes of pronouns . . . . .	45
11.2 Subtypes of numerals . . . . .	47
11.3 Subtypes of adverbs . . . . .	50
<b>12 Negation</b>	<b>51</b>
<b>13 Names and Terms</b>	<b>52</b>
13.1 Personal names . . . . .	52
13.2 Geographical names . . . . .	54
13.2.1 Countries, cities, rivers, mountains . . . . .	54
13.2.2 Streets, squares, stations . . . . .	55
13.2.3 Buildings . . . . .	55
13.3 Scientific terminology . . . . .	55
13.4 Other proper names . . . . .	56
<b>14 Abbreviations</b>	<b>57</b>
14.1 Fixed abbreviations of a single word . . . . .	57
14.2 Other abbreviations . . . . .	58
14.2.1 Well-known abbreviations composed of uppercase letters . . . . .	58
14.2.2 Less familiar abbreviations and abbreviations with many meanings . . . . .	58
14.2.3 Author's signature . . . . .	59
<b>15 Isolated Letters</b>	<b>60</b>
<b>16 Segments</b>	<b>61</b>
<b>17 Foreign Words</b>	<b>63</b>
17.1 Citation use . . . . .	63
17.2 Single word use . . . . .	64
17.3 Domesticated words of foreign origin . . . . .	65
<b>18 Aggregates</b>	<b>66</b>
<b>19 Hyphenated Composites</b>	<b>67</b>
<b>20 Typo, Distortion, Misspelling</b>	<b>68</b>
<b>21 Note on Tokenization</b>	<b>69</b>
<b>22 Appendix</b>	<b>70</b>
22.1 Detailed Part of Speech (SUBPOS): Quick Reference . . . . .	70
22.2 Categories Relevant for POS and SUBPOS Combinations . . . . .	73

# 1 Preface

Although the title of this report inherits the word "Manual" from the previous versions, it is no more intended to guide the annotators. Rather it attempts to describe the current state of the morphological annotation in the Prague Dependency Treebank – Consolidated 1.0 (PDT-C 1.0).<sup>1</sup> We believe that the guidelines can be of use to the users of the PDT-C 1.0 data, as well as for possible preparation of new data.

PDT-C 1.0 consists of four different datasets coming from PDT-corpora of Czech published earlier:

- dataset of written texts (the core PDT corpus in version 3.5),<sup>2</sup>
- dataset of translated texts (Prague Czech-English Dependency Treebank),<sup>3</sup>
- dataset of spoken texts (Prague Dependency Treebank of Spoken Czech),<sup>4</sup>
- datasets of user-generated texts (unpublished small treebank PDT-Faust).<sup>5</sup>

In the PDT-C project, we aim to provide all these treebanks with full manual annotation at the lower layers and unify and correct annotation at all layers. Specifically, the data in PDT-C 1.0 is (mainly) enhanced with a manual annotation at the morphological layer, consistently across all the four original treebanks. Altogether, the consolidated treebank contains almost 3,900,000 tokens with manual morphological annotation. The Czech morphological dictionary Morfflex,<sup>6</sup> which is now an integral part of the PDT-C 1.0 release, consists of more than 1 million lemmas/paradigms.

**Acknowledgment.** The research and language resource work reported in the paper has been supported by the LINDAT/CLARIAH-CZ projects funded by Ministry of Education, Youth and Sports of the Czech Republic (project LM2018101).

---

<sup>1</sup><https://ufal.mff.cuni.cz/pdt-c>

<sup>2</sup><http://ufal.mff.cuni.cz/pdt3.5>

<sup>3</sup><https://ufal.mff.cuni.cz/pcedt2.0>

<sup>4</sup><https://ufal.mff.cuni.cz/pdtsc2.0>

<sup>5</sup><https://ufal.mff.cuni.cz/grants/faust>

<sup>6</sup><https://ufal.mff.cuni.cz/morfflex>

## 2 Introduction

We do not want to substitute a grammarbook of Czech. So we will not systematically define word classes, morphological categories, paradigms, etc. All the annotators and users of the data and dictionary should understand the fundamentals of the Czech morphology, as most native Czech speakers do (the stuff is being taught in elementary schools). What we will describe the main principles of morphological annotation and focus on difficult and unusual phenomena.

The morphological annotation is based on a manual disambiguation of an automatic, dictionary-based morphological analysis of the annotated texts. For such automatic preprocessing, we use the MorphoDiTa tool.<sup>7</sup> In the annotation, a lemma (see Sect. 4) and a tag (see Sect. 5) is assigned to each wordform. The Lemma and the tag together uniquely identify the wordform (see Sect. 3.1). The annotation contains no syntactic structure, no attempt is even made to put together e.g. analytical verb forms or other types of multiword expressions (see Sect. 21).

Key element to annotation consistency is the Czech morphological dictionary **MorfFlex**, which is now an integral part of the PDT-C 1.0 release. The MorfFlex dictionary (see Sect. 3) is a flat list of lemma-tag-wordform triples.<sup>8</sup> It is in fact only an (automatic) derivative of the original, so-called “source format”, in which the dictionary is still being maintained. The source format is based on paradigm pattern system and a substantial part of the dictionary (65% lemmas) is mapped onto so-called derivational patterns. If a word belongs to a derivational patterns, several other words can be automatically derived from it. All automatically derived lemmas have the derivational information stored as a technical suffix of the lemma (Sect. 4.2.5). The suffix is really technical, primarily, it carries information about the automatic creation of a lemma; the manifested word-formation relation may not be correct or complete.<sup>9</sup> This is important for annotation. Because the creating lemmas from derivational patterns is an automatic process, there is no possibility to manually annotate stylistic and other categories of the derived lemmas (see more in Sect. 7).

The dictionary itself has undergone a long development process. It has been developed gradually since 1988. During this long time period, some phenomena originally included in the dictionary (e.g. word-formation relations, detailed annotation of terms and names) has been delegated to separate projects. In the source format, this information is preserved, but it is inconsistent and/or incomplete and it did not make it into the currently released version of the dictionary. We do not describe the source format here,<sup>10</sup> occasional reference to the source format is made only if it is necessary to explain a phenomenon in the currently described version of the MorfFlex 2020. The process of transformation the source format into the resulting dictionary, the description of which is, however, rather a technical matter, is not covered also in this document. We mention only some procedural aspects of this process (particularly the way of handling homonymy) that affect the way of representation of some phenomena in the dictionary, and therefore in the data.

The MorfFlex dictionary serves as a basis for annotation consistency. The goal of the annotation is full consistency between all the data and the dictionary. An inconsistency between the data and the dictionary indicates an annotation problem or error in the dictionary. All inconsistencies are corrected and there are only full matches now, except for a small amount of wordform occurrences in the data that are not in the dictionary (but have manual analysis in the data); this applies mostly to foreign wordforms and non-standard, sparse forms of Czech. However, if a wordform is in the dictionary, the dictionary contains all its morphological analyses (all possible lemma-tag pairs) that have been found in annotated data (and in other sources). A paradigm included in the dictionary (identified by the lemma) contains not only Standard Czech wordforms, it also provides non-standard variants and contains all forms found in the data, even defective forms, misspelling, typos, which are properly marked (see Sect. 6 and 20). The morphological annotation of a wordform that is in the data but not in the dictionary follows the same principles as applied to the dictionary.

<sup>7</sup><https://ufal.mff.cuni.cz/morphodita>

<sup>8</sup>The latest version of the dictionary contains 125,348,899 lemma-tag-wordform triples.

<sup>9</sup>The word-formation relations in Czech has been delegated to derivational data sources, such as Derinet: <https://ufal.mff.cuni.cz/derinet>

<sup>10</sup>The source format will be described in a separate document, which will be available in the first half of 2021.

### 3 Morphological Dictionary

The MorfFlex dictionary covers words (tokens) that occur in real Czech texts, i.e. Czech words, loan words and foreign words, proper nouns, abbreviations, isolated letters, part of words, numbers, and also punctuation and other non-alphanumeric characters (see also note on tokenization in Sect. 21). It captures both the singular and the plural set of wordforms of all inflected words, even of proper nouns. It is not focused only on standard Czech, the paradigms provide also non-standard variants and capture style characteristics of wordforms (see Sect. 6).

MorfFlex is a flat list of lemma-tag-wordform triples (see examples in Table 1). For each *wordform*, full inflectional information is coded in a *positional tag* (see Sect. 5). Wordforms are grouped into *paradigms* according to their formal morphological behavior. The paradigm is a set of all wordforms of the word. It is represented by a unique *lemma* (see Sect. 4). Apart from traditional morphological categories, the description also contains some semantic (see Sect. 8), stylistic (see Sect. 6) and derivational (see Sect. 7) information.

Wordform	Lemma	Tag
<i>podle</i>	<i>podle-1</i>	Dg-----1A----
<i>nepodle</i>	<i>podle-1</i>	Dg-----1N----
<i>podleji</i>	<i>podle-1</i>	Dg-----2A----
<i>nepodleji</i>	<i>podle-1</i>	Dg-----2N----
<i>podlejc</i>	<i>podle-1</i>	Dg-----2A---6
<i>nepodlejc</i>	<i>podle-1</i>	Dg-----2N---6
<i>nejpodleji</i>	<i>podle-1</i>	Dg-----3A----
<i>nejnepodleji</i>	<i>podle-1</i>	Dg-----3N----
<i>nejpodlejc</i>	<i>podle-1</i>	Dg-----3A---6
<i>nejnepodlejc</i>	<i>podle-1</i>	Dg-----3N---6
<i>podle</i>	<i>podle-2</i>	RR--2-----

Table 1: Example of the wordform-lemma-tag triples in the dictionary

#### 3.1 Principle of unique analysis

The principle of unique analysis, so called “golden rule of morphology”, is applied to the dictionary. The rule says that there must not exist more than one wordform with the same lemma and tag. Lemma and tag together uniquely identify the wordform. Two different wordforms always differ either in lemma or in tag. There are two means to ensure the principle is valid:

- **lemma numbering** (see Sect. 4.1.1),
- **tag numbering** at the 15<sup>th</sup> position (see Sect. 5.15),

Wordform	Lemma	Tag
<i>po lesich</i>	<i>les</i>	NNIP6-----A----
<i>po lesech</i>	<i>les</i>	NNIP6-----A---1

Table 2: Example of tag numbering

We use these means mainly for capturing homonymy of lemmas and different types of wordform variants. Each of these problematic issues is addressed differently. The former one is solved by adding a numerical index to homonymous lemmas (see Tab. 3), the latter one by adding a numerical index to 15<sup>th</sup> position of tag (see Tab. 2).<sup>11</sup>

<sup>11</sup>In the tables with examples, we present the analysis (lemma and tag) for the wordform given in the left column.

Wordform	Lemma	Tag
<u>náš stát</u>	stát-1_^(státní_útvart)	NNIS1-----A---
<u>stanu se vojákem</u>	stát-2_^(stanu_staneš)	VB-S---1P-AAP--
<u>stojím tu už dlouho</u>	stát-3_^(stojím_stojíš)	VB-S---1P-AAI--
<u>sníh staje</u>	stát-5_^(sníh)	VB-S---3P-AAI--

Table 3: Example of lemma numbering

### 3.2 Principle of morphological differentiation

MorfFlex captures primarily such phenomena that are of formal morphological nature. It does distinguish words with the same spelling (by adding numbers to lemmas; see Sect. 4.1.1), but different formal morphological characteristics.

Within one paradigm, all wordforms are allowed to appear with morphological tags of:

- **the same POS** (Sect. 5.1). The POS value is the same for the whole paradigm (identified by one lemma). Thus, there are two indexed lemmas for *drát*: **drát-1** for noun (meaning ‘wire’) and **drát-2** for verb with meaning ‘to pluck’ (see Tab. 4); or there are the lemmas **podle-1** and **podle-2** assigned in the dictionary to the string *podle* which can be either an adverb (meaning ‘meanly’) or a preposition (‘along’; see Tab. 1).

Similarly, homonymous forms of inflexible words (as well as inflexible loanwords) have as many lemmas in the dictionary as they express (inflexible) POS (see also Sect. 10.1). For example, *přece* which is in accordance to its function in a sentence interpreted as a conjunction (meaning ‘despite’) or as a particle (‘after all’), has two lemmas with different indexes and with different POS values at tags (see Tab. 4).

- **the same SUBPOS** (Sect. 5.2). The SUBPOS value (detailed part of speech) is the same for the whole paradigm except for a few exceptions:

- Paradigm of a verb is traditionally formed by several sets of forms (present/future forms, past participles, passive participles, transgressives, etc.) which are distinguished in the second SUBPOS position. However, all set of verbal wordforms are identified by one lemma.
- Short (nominal) forms of adjectives (e.g. *mlád* ‘young’) have a special value at the second SUBPOS position but together with long forms (e.g. *mladý* ‘young’) they are part of one paradigm. See examples in Table 4.
- In the paradigms of personal pronouns, a different value of SUBPOS indicates enclitic forms (cf. Sect. 11.1).

The SUBPOS value serves as an indicator which tag positions are to be filled and which not (i.e. the categories of GENDER, NUMBER, CASE, TENSE, etc.; see combination in Sect. 22.2). We are using unique values for SUBPOS category so that the value of the major speech category (the POS value) can be determined unambiguously from the value of the SUBPOS category. The only exceptions are abbreviations (POS = B) and segments (POS = S); potentially each SUBPOS value are possible for these POS (see more in Sect. 14 and Sect. 16).

- **the same GENDER** (Sect. 5.3) in case of nouns. The GENDER value of nouns is the same for the whole paradigm. For example, there are the lemma **rys-1** for masculine animate noun with meaning of an animal (‘lynx’), and lemma **rys-2** for masculine inanimate noun with meaning of ‘feature in face’ or ‘drawing’). Or there are two distinct paradigms of word

---

If there is context in the left column, then the wordform for which the analysis is given is underlined.

*kredenc*: *kredenc-1* as masculine and *kredenc-2* as feminine, even if they have the same meaning ('cupboard').<sup>12</sup>

The rule of different gender applies only to nouns (NN). For other part of speech with agreement gender (adjectives, etc.), a different value of GENDER does not necessarily mean a different paradigm.

- **the same ASPECT** (Sect. 5.13) in case of verbs. The value of ASPECT is the same for the whole paradigm. If a verb appears in two aspects, there should be two lemmas distinguished by the number. An example is the verb *pootevírat*. There must be two verbs with different aspects: *pootevírat-1* with imperfective aspect (meaning 'open slightly') and *pootevírat-2* with perfective aspect ('gradually open').

See examples in Table 4.

Wordform	Lemma	Tag
<i>ostrý drát</i>	drát-1	NNIS1-----A----
<i>drát peří</i>	drát-2	Vf-----A-I--
<i>ne peníze, ale přece lásku</i>	přece-1	J^-----
<i>přece jen nelhal</i>	přece-2	TT-----
<i>je mlád</i>	mladý	ACYS-----A----
<i>je mladý</i>	mladý	AAMS1----1A----
<i>dej mu</i>	on-1	P5ZS3---3-----
<i>jemu to nevadí</i>	on-1	PEZS3---3-----
<i>rys ostrovid</i>	rys-1	NNMS1-----A----
<i>rys v obličeji</i>	rys-2	NNIS1-----A----
<i>oškliví pavouci</i>	pavouk-1	NNMP1-----A----
<i>pavouky dvouher v tenise</i>	pavouk-2	NNIP1-----A----
<i>pod kredencem</i>	kredenc-1	NNIS7-----A----
<i>pod kredencí</i>	kredenc-2	NNFS7-----A----
<i>moje dítě</i>	dítě-1	NNNS1-----A----
<i>moje děti</i>	dítě-2	NNFP1-----A----
<i>oko na polévce i lidské</i>	oko-1	NNNS1-----A----
<i>oka na polévce</i>	oko-1	NNNP1-----A----
<i>oči lidské</i>	oko-2	NNFP1-----A----
<i>pootevírat trochu dveře</i>	pootevírat-1_^(otevírat_trochu)	Vf-----A-I--
<i>postupně pootevírat dveře</i>	pootevírat-2_^(postupně_otevírat)	Vf-----A-P--

Table 4: Examples: Principle of morphological differentiation

<sup>12</sup>The rule of same GENDER in noun paradigms is quite problematic in some cases. One group of problematic cases consists of words whose GENDER fluctuates even if they have the same meaning (cf. *ta kredenc* (fem.) - *ten kredenc* (masc.) 'cupboard'; *brambora* (fem.) - *brambor* (masc.) 'potato'; *ty bacily* (masc. inam.) - *ti bacilové* (masc. anim) 'bacillus'). Each GENDER variant is captured by a separate (indexed) paradigm/lemma. Particularly problematic are the cases of fluctuation between the masculine animate and non-animate gender. Some naturally inanimate words may have endings of the masculine animate in the genitive and accusative singular (e.g. *mám nového forda* 'I have a new Ford', *drží toho hřiba* 'she holds the mushroom'). The genitive/accusative wordforms with an animate ending are captured as a wordform variant within the non-animate paradigm (with I value on GENDER position). However, this solution is not adequate, especially due to gender agreement of the dependent adjective or predicate.

Another problematic group are nouns that have different grammatical gender in the singular and plural set of wordforms (cf. *dítě* (neut.) - *děti* (fem.) 'child - children', including derivatives of the *biodítě* 'bio-child', etc.). A similar problem is with words that have two sets of plural wordforms of two different gender (e.g. *ta oka* (neut.) - *ty oči* (fem.) 'loops - eyes'). To follow the rule of the same GENDER, we have separated the respective singular and/or the plural sets of wordforms into different paradigms. See examples in Tab. 4.

### 3.3 Principle of unique paradigm

MorfFlex does not capture any differences in meanings of homonymous or polysemous words. It means that there are no two identical paradigms (set of lemma-tag pairs) in the dictionary. It tries to avoid lexical distinctions that are not morphologically based: lexical homonymy and polysemy are not marked for words which do behave morphologically in the same way. Due to a large number of complicated cases, we do not take into account derivational, stylistic nor semantic differences.

Wordform	Lemma	Tag
přiletěl <u>jeřáb</u>	jeřáb-1_^(pták)	NNMS1-----A----
u silnice roste <u>jeřáb</u>	jeřáb-2_^(strom;stroj)	NNIS1-----A----
<u>jeřáb</u> na staveništi	jeřáb-2_^(strom;stroj)	NNIS1-----A----
vodovodní <u>kohoutek</u>	kohoutek-2_^(květina;uzávěr)	NNIS1-----A----
luční <u>kohoutek</u>	kohoutek-2_^(květina;uzávěr)	NNIS1-----A----
slepička a <u>kohoutek</u>	kohoutek-1_^(pták)	NNMS1-----A----
kovová <u>matka</u>	matka	NNFS1-----A----
<u>matka</u> spí	matka	NNFS1-----A----
<u>palička</u> a <u>palič</u>	palička_^(*)2	NNFS1-----A----
<u>palička</u> na maso	palička_^(*)2	NNFS1-----A----
pracuje jako <u>ekonomka</u>	ekonomka_^(*)2	NNFS1-----A----
studuje <u>ekonomku</u>	ekonomka_^(*)2	NNFS4-----A----
<u>Barča</u> a <u>Helča</u>	Barča_;;G_-;Y	NNFS1-----A----
sezdeme se na <u>Barče</u>	Barča_;;G_-;Y	NNFS6-----A----

Table 5: Examples: Principle of unique paradigm

Thus, we do not distinguish lemmas that have the same paradigm and that have:

- **different meaning.** If two (or more) words share all the morphological properties, there is only one lemma/paradigm in the dictionary. It means that we do not distinguish between the words *kohoutek* as ‘a flower’ and *kohoutek* as ‘a tap’ because both have the same inflectional model, namely for masculine inanimate noun. On the other hand, there is a different word *kohoutek* (‘a small cock’), that has different inflectional model for masculine animate noun. Similarly, there is one paradigm/lemma for word *palička* which has two meaning. First meaning is the meat knocking tool ‘tenderizer’ and with this meaning, the word *palička* is derived from the word *palice* ‘mallet’. Another meaning is ‘arsonist-female’ and with this meaning, the word *palička* is derived from the word *palič* ‘arsonist-male’. These two meanings have the same inflectional paradigm, so there is only one paradigm in the dictionary. The lemma in the dictionary, *palička\_^(\*)2*, contains the derivational comment (see Sect. 4.2.5) indicating that the paradigm was automatically derived from the word *palič* ‘arsonist-male’ (see more in Sect. 7). This does not mean that the paradigm with lemma *palička\_^(\*)2* cannot be used to analyze the wordforms of the word *palička* derived from *palice* ‘mallet’.
- **different derivational model.** If two (or more) words share all the morphological properties, there is only one lemma/paradigm, even if they differ in word-formation relations. E.g., the word *jeřábník* (‘man who works with a crane-device’) is derived from *jeřáb* as a ‘device’. It is not possible to derive *jeřábník* from *jeřáb* as a ‘tree’, but we do not distinguish the meanings ‘tree’ and ‘device’, because they do not differ from the inflectional point of view. Another example is word *matka*, that has two different meanings, namely a ‘nut’ and ‘mother’. These two meanings have the same inflectional paradigm, but their derivational behavior differs. It is possible to derive a possessive adjective only from word *matka* as a ‘mother’. But these derivational differences are not captured in the dictionary.
- **different style value.** If two (or more) words share all the morphological properties, there is only one lemma/paradigm, even if there is a difference in style value. E.g., there is a

standard word *ekonomka* as ‘female economist’ (automatically derived from word *ekonom* ‘male economist’; see derivational comment in the lemma) and a non-standard one with meaning ‘school of economics’. These two meanings have the same inflectional paradigm, they only differ in terms of the stylistic characteristics, but we do not capture this difference in our description. There is only one lemma/paradigm without any style label. See more about style labeling in Sect. 6.

- **different semantic label.** If two (or more) words share all the morphological properties, there is only one lemma/paradigm, even if there is a difference in so-called name label (Sect. 4.2.2). E.g., there is a female name *Barča* (captured by name label Y) and also common name of Barricade House of Culture *Barča* (captured by label G). The names have the same inflectional paradigm, so there is only one lemma/paradigm with all semantic flags relevant.

See examples in Tab. 5.

## 4 Lemma Structure

Lemma represents the whole paradigm. It has two parts. The first part, the lemma proper (Sect. 4.1), has to be a unique identifier of the lexical item/paradigm. The second part, so-called AddInfo (Sect. 4.2), is optional, it is not part of the identifier and contains additional information about the paradigm, e.g. semantic or derivational information, style label. They are related to the whole paradigm (set of wordforms) belonging to the given lemma.

A lemma is the same for all wordforms in the paradigm. Two lemmas with different AddInfo must differ in lemma proper. To achieve that, numbers are used to distinguish the lemmas.

The formal description of the lemma structure is in Tab. 6. Spaces were inserted between nonterminals to improve readability. In fact, there are no spaces in lemmas. Capitalized multi-character symbols are nonterminals. All other symbols are terminals.<sup>13</sup>

**TODO upravit: Lemma number nikdy není 0! TODO Comment - odlišila bych jako typ (VariantInfo) ^DD\*\*Word DS GC ?? TODO Term přejmenovat na Name. Term může být v lemmatu víckrát, Style jen jednou**

```

Lemma ::= LemmaProper | LemmaProper AddInfo
LemmaProper ::= Word | Word - Number | Number0 | SpecialChar
Word ::= Letter | Letter Word
Letter ::= A | a | Á | á | Ä | ä | ... | Z | z | Ž | ž | ,
Number ::= NonZero | NonZero Number0
Number0 ::= Digit | Digit Number0
NonZero ::= 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
Digit ::= 0 | NonZero
SpecialChar ::= ! | " | # | $ | % | & | ' | ( | ) | * | + | , | -
| . | / | : | ; | < | = | > | ? | @ | [ | \ | ] |
| ‘ | { | | | } | § | °
AddInfo ::= Reference Term Style Comment
Reference ::= <empty> | ` LemmaProper
Term ::= <empty> | _; Term1 Term
Style ::= <empty> | _, Style1
Comment ::= <empty> | _^ Comment1
Term1 ::= Y | E | G | U | m | o
Style1 ::= n | a | s | h | e | l | v | i
Comment1 ::= ( Explanation ) | ( Derivation ) | ( Explanation )_( Derivation )
Explanation ::= CommentChar | CommentChar Explanation
Derivation ::= * Number Word | * Word
CommentChar ::= Letter | Digit |
| ! | " | # | $ | % | & | ' | + | , | . | ~ | -
| / | : | ; | < | = | > | ? | @ | [ | \ | ] |
| ‘ | { | | | } | ^ | § | °

```

Table 6: Formal structure of the lemma

---

<sup>13</sup>Any character that is Letter in the Unicode standard can appear in place of the Letter nonterminal. In the non-ASCII area this most frequently applies to the Czech accented characters: ÁáČčĎďÉéĚěÍíŇňÓóŘřŠšŤťÚúŮůÝýŽž. However, other characters occur in names (e.g. German ÄäÖöÜü, Serbo-Croatian Ćć) and in foreign words (e.g. Slovak ĽĽÍÓóŔŕ). Standard HTML entities (such as &amp; for & or &agrave; for à) are also allowed. PDT 1.0 was encoded in the ISO Latin 2 codepage, so representing any West European characters required using entities. PDT 2.0 and the current version of PDT-C 1.0 shall be encoded in UTF8, so few entities are needed.

The single quote (') in some transcriptions of non-Latin alphabets (e.g. in Chinese *Mao C'-tung*, Hebrew *Be'er Sheva'*) and if it marks deleted parts of words (e.g. English *don't*, French *d'Artagnan*), is considered a SpecialChar and it splits the string into three tokens (*d ' Artagnan*). See also in Sect. 21.

## 4.1 Lemma proper

Lemma proper has to be a unique identifier of the paradigm. Usually it is the base form of the word (e.i. nominative singular for nouns, the same plus masculine positive for adjectives, similarly for pronouns and numerals; verbs are represented by their infinitive forms),<sup>14</sup> possibly followed by a number distinguishing different lemmas with the same spelling but different formal morphological behavior.

**Lowercase and uppercase letters in lemma.** Words that have to be always capitalized in writing, have their lemma capitalized as well (ex. *Král* as a ‘surname’, *Písek* as a ‘surname’ and as a ‘town’; cf. Table 7 and Table 8). On the other hand, a lower-case lemma is automatically assumed to have also a wordform with uppercase letters; cf. Tab. 7. Capitalized lemmas are usually, but not always assigned a name label (see Sect. 4.2.2).

Base form in lemma is case-sensitive. The case of the initial letter (any different character in a base form) is a sufficient distinction to consider the two lemmas to be different. As a consequence, *král* ('king') and *Král* ('surname') need not be distinguished by numbers; cf. Tab. 7.<sup>15</sup>

Wordform	Lemma	Tag
<i>král a královna</i>	král	NNMS1-----A----
<i>pan Král</i>	Král_ ; Y	NNMS1-----A----
<i>Král a královna odešli.</i>	král	NNMS1-----A----
<i>KRÁL a královna</i>	král	NNMS1-----A----
<i>pan KRÁL</i>	Král_ ; Y	NNMS1-----A----

Table 7: Examples: Lowercase and uppercase letters in lemma

### 4.1.1 Lemma number

Lemma number distinguishes different lemmas with the same spelling but different formal morphological behavior (cf. Sect. 3.2). Thus, we have the lemma *pět-1'5* for number ('five') and lemma *pět-2\_^(zpívat)* for verb ('to sing'); cf. Tab. 8.

However, there are unnumbered and numbered lemmas with the same base form in the dictionary and data. The used numbers also often do not form a continuous sequence. However, all different paradigms with lemmas of same base form are distinguished. Cf. lemmas for word *Písek* in Tab. 8.

Wordform	Lemma	Tag
<i>pět dětí</i>	pět-1'5	Cn-S1-----
<i>nemohl pět</i>	pět-2_^(zpívat)	Vf-----A-I--
<i>ta sršeň</i>	sršeň-1	NNFS1-----A----
<i>ten sršeň</i>	sršeň-2	NNMS1-----A----
<i>pan Písek</i>	Písek_ ; Y	NNMS1-----A----
<i>do Písku</i>	Písek-2_ ; G	NNIS2-----A----
<i>přivezl písek</i>	písek	NNIS4-----A----

Table 8: Examples: Lemma number

**Special lemma index number.** There are special lemma numbers used for a predefined word types (see ex. in Tab. 9 and respective sections):

<sup>14</sup>The cases where more than one "base form" is available are described in Sect. 9.3.2.

<sup>15</sup>The capitalization inside words (e.g. *AIDSu*, *McDonald*) are not sufficiently addressed in the dictionary and annotation.

- 33** The lemma number -33 is for an isolated letter (see more in Sect. 15),
- 77** The lemma number -77 is for a foreign word (see more in Sect. 17),
- 88** The lemma number -88 is for a special type of abbreviation (see more in Sect. 14),
- 99** The lemma number -99 is for an authors' signature (see more in Sect. 14.2.3).

Conventions of this kind exist solely for the convenience of a human reader.

Wordform	Lemma	Tag
písmeno <i>A</i>	A-33	Q3-----
cizí slovo <i>black</i>	black-77	F%-----
zkratka <i>OP</i>	OP-88	BNXXX----A---
šifra <i>mat</i>	mat-99_-;Y	BNXXX----A---

Table 9: Examples: Special lemma index number

## 4.2 Additional information about the paradigm (AddInfo)

The second part of lemma is optional, it is not part of the identifier and contains additional information about the paradigm (AddInfo in sequel), namely semantic and derivational information, style label. The information attached to lemma is related to the whole paradigm belonging to the given lemma (e.g. a style label (see Sect. 4.2.3) attached to lemma is valid for all wordforms of the paradigm).

### 4.2.1 Reference

The reference is one of the means of explaining the meaning of the lemma. It is systematically used only with spelled-out numbers and with Roman numbers: in the additional lemma field, there is a corresponding number in Latin numerals. The numeric item is separated from the proper lemma by an apostrophe (cf. examples in Tab. 10).

The numerical reference is used also for the derivatives of the basic definite numerals (i.e. for words such as *šest* ‘six’, *šestý* ‘sixth’, but also *šestina* ‘sixth’, *šestkrát* ‘six times’). The numerical reference only indicates which number the word refers to. Replacement of the word with the number in the reference is not assumed in all cases (e.g. *šestina* ‘sixth’ is not replaceable with 6).

Wordform	Lemma	Tag
v <i>šest hodin</i>	šest'3	C1-S4-----
<i>šestina</i>	šestina'6_^(*)3	NNFS1-----A---
rímské číslo <i>M</i>	M-3'1000	C}-----

Table 10: Examples: Reference

### 4.2.2 Name label

Words that have to be always capitalized in writing (names and other terms), have their lemma capitalized as well (e.g. *Novák* as a surname, *Praha* as a name of a town). Capitalized lemmas are usually (but not always) assigned a so-called name label.

The capitalized word type is indicated by ”;” followed by a letter (see Tab. 11). More than one name label may apply to one lemma. The name labels are sorted alphabetically except for *m* which is always the last. Examples can be found in the respective sections, mostly in Sect. 13.

Label	Description	Example
Y	personal name	Petr, Novák, Aristoteles
E	member of nation, inhabitant of territory	Slovák, Newyorčan
G	geographical name	Praha, Tatry, Václavák
U	scientific terminology	Australopithecus, Rh, Hydrosulfit
m	other proper name	Madeta, Opel, Sázavafest
o	color indication	červený, modrobílý

Table 11: Name labels

**Color indication.** Adjective color names are marked by a special label `_;o` (e.g. lemma `zelený_;``o` for adjective `zelený`). This is mainly due to the specific morphological behavior of these words. Color names can be put together almost indefinitely into compound names (`červenomodrobílý` ‘lit. red-blue-white’ etc.). Compared to adjectives, they have different derivative behavior: adjective color names, besides having the usual adverb derivation (`zeleně` ‘greenly’), have a very regular adverb derivation with the suffix `-o` and the prefix `na-`; (ex. `tam je zeleno` lit. ‘there is greenly’, `natřít nazeleno` lit. ‘paint greenly’).

#### 4.2.3 Style label

Lemmas can be stylistically classified. The style flag is indicated by ”`_`,” followed by a letter (see the values in Tab. 12). Standard lemmas have no stylistic flag but any lemma intended for special usage (bookish, dialect, slang, etc.) should be marked as such. At most one style flag is attached to any lemma. For automatically derived lemmas (lemmas with Derivation Info; see Sect. 4.2.5), the style label is not inherited (see more in Sect. 7). More about style labeling is in Sect. 6.

Label	Description	Example
s	standard word, bookish	asyl, kurs
a	archaic word	šlojíř, stout, these
h	non-standard word	Karlák, šutrák,
i	distortion, typo	součastník, intensívni
n	dialect	šufánek, čučkař
e	expressive word	kládička, sajtnička
l	slang, argot	genštáb, vertajmka
v	vulgar word	vlezdoprdelizmus

Table 12: Style labels

#### 4.2.4 Variant info

Orthographic and stylistic variants of a word (e.g. an archaic variant `these`, a standard variant `teze`, and a non-standard variant `téze` ‘thesis’) are captured in separate paradigms with different variant lemmas, but they are ”interconnected” using the Variant Info field in the AddInfo part of lemma. One of the variants is selected as “basic” (the standard one, i.e. `teze`) and other variants (non-standard `téze` and archaic `these`) refer to it: in comment brackets (following the caret sign `^`), the type of variant (see Tab. 13) is given after the other caret sign; after the two asterisks, there is the lemma of the basic variant. For example, lemma `these_,a_^(^DD**teze)` ”interconnect” archaic variant `these` with the basic one and lemma `téze_,h_^(^GC**teze)` ”interconnect” non-standard variant `téze`. More about capturing variants is in Sect. 9.1.

Type	Description
DD	standard variant
GC	non-standard variant
DS	distortion, typo

Table 13: Variant info

#### 4.2.5 Derivation info

For regular derivations, the lemma contains information about base lemma it is derived from. For example, lemmas of possessive adjectives (e.g. `otcův_^(*3ec)`) contain information about the noun they are derived from (`otec` ‘father’s’ ← `otec` ‘father’). The originating lemma is (for space saving reasons only) written in the form of a rule, which has two parts:

1. number of characters to remove from the end of the current lemma identification, expressed as a decimal number
2. characters to add to form the original lemma.

Only the proper lemmas are both input and output of this process (including the lemma number, if present). Each such rule must start with an asterisk to distinguish it from the explanation Comment (see next Sect. 4.2.6). Thus, for example, derivation info in lemma `otcův_^(*3ec)` means remove 3 characters, add `ec` to get `otec` or derivation info in lemma `koníčkův-2_^(*5ek-2)` means remove 5 characters, add `ek-2` to get `koníček-2`.

In the current version of the dictionary, only lemmas automatically derived from derivational patterns have filled-in derivation info. In the source format (from which the dictionary is made), there are also manually created comments on word-formation relation. However, due to their incompleteness and inconsistency, they were not transferred into the dictionary. More about derivations is in Sect. 7.

#### 4.2.6 Explanational comment

Any string in parentheses can be used as explanation of the lemma/paradigm usage (meaning). The string cannot contain spaces or parentheses. The underscore character is used to replace space, square brackets are used instead of parentheses. The explanation is in Czech. Example of usage, synonym etc. can also be used or both a verbal description and an example can be mixed.

An explanatory comment is optional in all cases. There is no rule when the comment can be used and when cannot. The lemmas automatically derived from derivational patterns (with Derivation info; see Sect. 4.2.5) do not have the explanatory comment (see Sect. 7).

**Note.** The inconsistency in the comments is historically caused by how careful or active the dictionary administrator was. The comments were never intended to replace the interpretation of meaning, they are primarily perceived as “hints” for annotators and as such they can be very useful. In the context of the principle of a unique paradigm (Sect. 3.3), it must be emphasized that the paradigm/lemma is applicable to all cases with same morphological behavior, regardless of the explanation given in the comment.

## 5 Tag Structure

Every tag is represented as a string of 15 symbols. Each position in the string corresponds to one morphological category according to a more or less traditional system of formal morphology.

A value in each category is represented as a single symbol, mostly an uppercase letter of the English alphabet (for example, P for plural), sometimes also another symbol (f for an infinitive, ^ for conjunction). In some cases no distinction among “traditional” values is being made where the possibility of correctly distinguishing them is low based on local context. For example, possessive pronouns in third person plural are not distinguished in gender and number, nor in case; passive participles (both active and passive) in masculine are not distinguished in animateness, etc. Typically, a letter X is used where all possible values might be considered in a more detailed tagset, or a special letter is used with a more restricted choice (e.g. Y is used for masculine animate/masculine inanimate “non-distinction”).

Non-applicable values are denoted by a single hyphen – (e.g. tense for nouns). Which categories are applicable/non-applicable is determined by the 2<sup>nd</sup> position of the tag (SUBPOS; Sect. 5.2). The tables of applicability/non-applicability of the tag categories related to a SUBPOS value are in Sect. 22.2.

An overview of the 15 tag positions is in Table 14. The categories and their values are described in the following sections.

#	Category Name	Description in English	Description in Czech
1	POS	Part of Speech	Slovní druh
2	SUBPOS	Detailed Part of Speech	Slovní poddruh
3	GENDER	Gender	Rod
4	NUMBER	Number	Číslo
5	CASE	Case	Pád
6	POSSGENDER	Possessor’s Gender	Rod vlastníka
7	POSSNUMBER	Possessor’s Number	Číslo vlastníka
8	PERSON	Person	Osoba
9	TENSE	Tense	Čas
10	GRADE	Degree of Comparison	Stupeň
11	NEGATION	Negation (by prefix)	Negace
12	VOICE	Voice	Slovesný rod
13	ASPECT	Aspect	Vid
14	AGGREGATE	Aggregate	Agregát
15	VAR	Variant	Varianta

Table 14: Categories in positional tag

### 5.1 Part of speech (1st position)

The POS category denotes the main part of speech, according to the traditional Czech scheme known from both comprehensive as well as high-school grammars. However, the assignment of the POS values is driven mainly by the requirements of consistency in further processing, therefore it is not always in line with traditional grammars. In addition to the ten traditional parts of speech, we distinguish also the abbreviations (Sect. 14), foreign words (Sect. 17), segments (Sect. 16), isolated letters (Sect. 15) and punctuation (Sect. 21). An overview of the 1<sup>st</sup> position values is in Table 15.<sup>16</sup>

<sup>16</sup>In morphological analysis using automatic tool, there can appear one more value of POS, namely X. In that case, the SUBPOS value is always @. It means that the wordform is not in the dictionary. In the PDT-C annotation, these values are not used; we analyze all wordforms according to the rules described here, even those that are not in the dictionary.

POS	Description	POS	Description
A	Adjective	P	Pronoun
B	Abbreviation	Q	Letter
C	Numeral	R	Preposition
D	Adverb	S	Segment
F	Foreign word	T	Particle
I	Interjection	V	Verb
J	Conjunction	Z	Punctuation
N	Noun		

Table 15: Values of POS category (1<sup>st</sup> position)

## 5.2 Detailed part of speech (2nd position)

This category is the most detailed one; it contains values for fine-grained distinction of the major part of speech category. Its primary technical purpose, however, is to serve as an indicator of applicability/non-applicability of other categories (i.e. the categories GENDER, NUMBER, CASE, etc. up to the last category, VAR). We are using unique values so that the value of the major speech category can be determined unambiguously from the value of the SUBPOS category. There are two exceptions: abbreviations (B) and segments (S), where potentially each SUBPOS value are possible; see Sect. 14 and 16).

An alphabetical list of all 66 SUBPOS values is in the appendix (Sect. 22.1). There are also tables of applicability/non-applicability of other categories of tag related to SUBPOS value (Sect. 22.2). Here, for each of the POS values, a list of detailed part of speech (SUBPOS) values and examples is given (the tables are arranged alphabetically according to the POS value).

### A Adjective

POS	Detailed part-of-speech used (SUBPOS)
A	A C G M O U

Examples:

POS & SUBPOS	wordform & translation	description
AA	<i>technický</i> ‘technical’	adj. general
AC	<i>mlád</i> ‘(be) young’	adj. nominal
AG	<i>kouřící</i> ‘smoking’	adj. derived from present transgressive
AM	<i>zvítězivší</i> ‘having-won’	adj. derived from past transgressive
AO	<i>svůj</i> ‘to be himself’	<i>svůj, nesvůj</i> in specific usage, <i>tentam</i>
AU	<i>Martinin</i> ‘Martina’s’	adj. possessive

## B Abbreviation

For POS with a value of B, any SUBPOS value is allowed (see Sect. 14). In the tables, the SUBPOS values that occurred in the PDT-C data are listed.

POS	Detailed part-of-speech used (SUBPOS)
B	$\wedge$ A b N

Examples:

POS & SUBPOS	wordform & translation	description
B $\wedge$	<i>tzn</i> ‘i.e.’	abbr. of conjunction
BA	<i>aj</i> ‘and other’	abbr. of adjective phrase
Bb	<i>atd</i> ‘and so on’	abbr. of adverb phrase
BN	<i>USA</i>	abbr. of noun phrase

## C Numeral

POS	Detailed part-of-speech used (SUBPOS)
C	= a d h j k l n o r v w y z }

Examples:

POS & SUBPOS	wordform & translation	description
C=	1.23 ‘1.23’	num. written using digits
C}	XIV ‘XIV’	Roman numeral
Ca	<i>několik</i> ‘several’	num. cardinal indef., non-adjectival declension
Cd	<i>čtverý</i> ‘four-kinds’	num. generic, adjectival declension
Ch	<i>několikerý</i> ‘several-kind’	num. generic indef., adjectival declension
Cj	<i>čtvero</i> ‘four-kinds-of’	num. generic, noun usage
Ck	<i>několikero</i> ‘several-kinds’	num. generic indef., noun usage
C1	<i>čtyři</i> ‘four’, <i>pět</i> ‘five’	num. cardinal, no gender
Cn	<i>jeden</i> ‘one’, <i>dva</i> ‘two’	num. cardinal, agreement gender
Co	<i>několikrát</i> ‘several-times’	num. multiplicative indef.
Cr	<i>druhý</i> ‘second’	num. ordinal
Cv	<i>sedmkrát</i> ‘seven-times’	num. multiplicative
Cw	<i>několikátý</i> ‘several’	num. ordinal indef.
Cy	<i>nejeden</i> ‘not-one’	num. cardinal indef., agreement gender
Cz	<i>sto</i> ‘hundred’	num. cardinal, noun usage

More details about detailed part of speech of numerals are in Sect. 11.2.

## D Adverb

POS	Detailed part-of-speech used (SUBPOS)
D	b g

Examples:

POS & SUBPOS	wordform & translation	description
Db	<i>nahoru</i> ‘up’	no degrees of comparison, no negation
Dg	<i>rychle</i> ‘quickly’	negation, degrees of c. possible

More details about detailed part of speech of adverbs are in Sect. 11.3.

F Foreign word

POS	Detailed part-of-speech used (SUBPOS)
F	%

Examples:

POS & SUBPOS	wordform & translation	description
F%	<i>The</i>	foreign word

I Interjection

POS	Detailed part-of-speech used (SUBPOS)
I	I

Example:

POS & SUBPOS	wordform & translation	description
II	<i>ach</i> 'oh!'	interjection

J Conjunction

POS	Detailed part-of-speech used (SUBPOS)
J	* , ^

Examples:

POS & SUBPOS	wordform & translation	description
J*	<i>krát</i> 'times'	binary math. operations
J,	<i>že</i> 'that'	conj. subordinate
J^	<i>a</i> 'and'	conj. coordinating

N Noun

POS	Detailed part-of-speech used (SUBPOS)
N	N

Example:

POS & SUBPOS	wordform & translation	description
NN	<i>robot</i> 'robot'	any noun incl. proper

P Pronoun

POS	Detailed part-of-speech used (SUBPOS)
P	1 4 5 6 7 8 9 D E H K L P Q S W Y Z

Examples:

POS & SUBPOS	wordform & translation	description
P1	<i>jehož, jejíž</i> ‘whose’	p. relative possessive
P4	<i>jaký</i> ‘what’	p. interrogative/relative with adj. declension
P5	<i>mu</i> ‘him’	p. personal for 3 <sup>rd</sup> person, clitic
P6	<i>sebe</i> ‘himself’	p. personal reflexive in long forms
P7	<i>se, si</i>	p. personal reflexive clitic
P8	<i>svůj</i> ‘his’	p. personal reflexive possessive
P9	<i>jeho</i> ‘his’, <i>její</i> ‘her’	p. personal possessive for 3 <sup>rd</sup> person
PD	<i>tento</i> ‘this’	p. demonstrative
PE	<i>on</i> ‘he’, <i>ona</i> ‘she’	p. personal for 3 <sup>rd</sup> person
PH	<i>mě</i> ‘me’, <i>ti</i> ‘you’	p. personal, no gender, clitic
PK	<i>někdo</i> ‘somebody’	p. indefinite, no gender
PL	<i>všechn</i> ‘all’, <i>sám</i> ‘alone’	p. delimiting
PP	<i>já</i> ‘I’, <i>ty</i> ‘you’	p. personal, no gender
PQ	<i>kdo</i> ‘who’	p. interrogative/relative, no gender
PS	<i>můj</i> ‘my’	p. personal possessive
PW	<i>nicí</i> ‘nobody’s’	p. negative, adjectival declension
PY	<i>nic</i> ‘nothing’	p. negative, no gender
PZ	<i>nějaký</i> ‘some’	p. indefinite with adj. declension

More details about detailed part of speech of pronouns are in Sect. 11.1.

Q Letter

POS	Detailed part-of-speech used (SUBPOS)
Q	3

Examples:

POS & SUBPOS	wordform & translation	description
Q3	A	isolated letter

R Preposition

POS	Detailed part-of-speech used (SUBPOS)
R	F R V

Examples:

POS & SUBPOS	wordform & translation	description
RF	<i>nehledě</i> ‘regardless’	part of preposition
RR	<i>v</i> ‘in’	prep. general (no vocalization)
RV	<i>ve</i> ‘in’	prep. with vocalization

## S Segment

For POS with a value of S, any SUBPOS value is allowed (see Sect. 16). In the table, the SUBPOS values that occurred in the PDT-C data are listed.

POS	Detailed part-of-speech used (SUBPOS)
S	2 A b N n

Examples:

POS & SUBPOS	wordform & translation	description
S2	<i>sci</i>	prefixal segm.
SA	<i>upový</i>	postfixal segm. of adjective
Sb	<i>line</i>	postfixal segm. of adverb
SN	<i>upista</i>	postfixal segm. of noun
S1	<i>ti</i>	postfixal segm. of numeral

## T Particle

POS	Detailed part-of-speech used (SUBPOS)
T	T

Examples:

POS & SUBPOS	wordform & translation	description
TT	<i>jen</i> ‘only’	particle

## V Verb

POS	Detailed part-of-speech used (SUBPOS)
V	B c e f i m p q s t

Examples:

POS & SUBPOS	wordform & translation	description
VB	<i>dělám</i> ‘(I) do’	present/future form
Vc	<i>býchom</i> ‘(we) would’	conditional of <i>být</i> ‘to be’
Ve	<i>dělajíce</i> ‘(they-)doing’	transgressive present
Vf	<i>dělat</i> ‘(to) do’	infinitive
Vi	<i>dělejme</i> ‘(let’s) do’	imperative
Vm	<i>udělav</i> ‘having-done’	transgressive past
Vp	<i>dělali</i> ‘(they) did’	past participle
Vq	<i>dělalť</i> ‘(he) did’	archaic past participle with -ť
Vs	<i>děláno</i> ‘(it) was-being-done’	passive participle
Vt	<i>dělámť</i> ‘(I) do’	archaic present/future form with -ť

## Z Punctuation

POS	Detailed part-of-speech used (SUBPOS)
Z	:

Examples:

POS & SUBPOS	wordform & translation	description
Z:	, %	punctuation, non-alphanumeric character

### 5.3 Gender (3rd position)

Grammatical gender is being described at the 3<sup>rd</sup> position – both the lexical gender of nouns, as well as the agreement gender of verbs, adjectives, pronouns and numerals.

Czech grammatical gender is considered to have four different values: masculine animate, masculine inanimate, feminine, and neuter (M, I, F and N, respectively). The tags denoting ambiguous combinations of the four basic gender tags are: H (feminine or neuter), Q (feminine singular or neuter plural), T (masculine inanimate or feminine in plural), Y (for masculine, regardless of animateness), and Z (for not feminine forms). The tag X is used in its usual sense (“any gender”). Except for X, the other ambiguous tags are never used for nouns or long adjectival forms. An overview of the 3<sup>rd</sup> position values is in Tab. 16.

GENDER	Description	Examples
F	Feminine	píseň, malá,
H	{F, N} Feminine or Neuter	udělajíc, dvě, moje
I	Masculine inanimate	dům, malý
M	Masculine animate	učitel, mladí, oni
N	Neuter	město, malé, běhalo
Q	Feminine (singular only) or Neuter (plural only) – only with participles and nominal forms of adjectives	schopna, běhala, ráda
T	Masculine inanimate or Feminine (plural only) – only with participles and nominal forms of adjectives	schopny, běhaly, rády
X	Any	otcovic, jejich, mymi
Y	{M, I} Masculine (either animate or inanimate)	schopen, běhal, jeden, on
Z	{M, I, N} Not feminine – only for (some) pronoun forms and certain numerals	měho, jednoho (genitive)

Table 16: Values of GENDER category (3<sup>rd</sup> position)

NUMBER	Description	Examples
D	Dual	dvěma malýma nohama
P	Plural	dvě malé nohy
S	Singular	malá noha
W	Singular for feminine gender, plural with neuter – only in participle, nominal adjective with Q gender	schopna, běhala
X	Any	finále, jejich

Table 17: Values of NUMBER category (4th position)

### 5.4 Number (4th position)

The number category mostly takes on only one of the two standard values: S for singular or P for plural. Nevertheless, Czech still uses so-called dual number for several nouns denoting symmetrical body parts: *oči* ‘eyes’, *ramena* ‘arms’, *nohy* ‘legs’, *uši* ‘ears’ (but not, for example, *kolená* ‘knees’). Due to strong agreement rules, this distinction applies also to adjectives, pronouns and for numeral *dvě* ‘two’. However, as the dual manifests itself only in the instrumental case, it is not necessary to make this distinction for verb number (the nominal agreement only applies to the nominative case). Also, as in all the other cases (nominative through locative) this distinction does not appear on the surface, we also do not make this distinction there. Moreover, all the nouns displaying the dual forms are in feminine (the plural of *prsa* ‘breasts’ (of neuter gender) is considered an ordinary

plural), thus the D value is actually used only in tags having also the feminine value of the gender category.

The special value W is used in connection with the GENDER value Q, in order to distinguish the combination of (feminine, singular) and (neuter, plural) from the simple cases (see ex. in Tab. 18).

The value X is used with undeclinable nouns and adjectives, and with pronouns with systematically ambiguous number. It is also used for abbreviations, where the grammatical number distinction is systematically hidden. An overview of the 4<sup>th</sup> position values is in Tab. 17.

Wordform	Lemma	Tag
<i>matka ja churava</i>	churavý	ACQW-----A---
<i>děvčata jsou churava</i>	churavý	ACQW-----A---
<i>matka přišla</i>	přijít	VpQW----R-AAP--
<i>děvčata přišla</i>	přijít	VpQW----R-AAP--

Table 18: Examples: Special value W of NUMBER category

## 5.5 Case (5th position)

Czech traditionally distinguishes among seven cases: nominative, genitive, dative, accusative, vocative, locative and instrumental. Traditionally, the cases are numbered, see Tab. 19).

CASE	Description	Example
1	Nominative	<i>žena, ženy</i>
2	Genitive	<i>ženy, žen</i>
3	Dative	<i>ženě, ženám</i>
4	Accusative	<i>ženu, ženy</i>
5	Vocative	<i>ženo, ženy</i>
6	Locative	<i>ženě, ženách</i>
7	Instrumental	<i>ženou, ženami</i>
X	Any	<i>finále</i>

Table 19: Values of CASE category (5th position)

Once a case is relevant for some detailed part of speech category (i.e. for some SUBPOS), it can take any of the seven values (with negligible – and questionable – exceptions for vocative for personal pronouns). Virtually all nouns, adjectives, pronouns, and numerals express case, and for agreement purposes, we use the CASE category also for all prepositions. Having decided that verb valency is not part of morphological processing, there is no case agreement information linked to verbs. However, there is one case where CASE is present at a verb form: in passive participle (Vs). However, we distinguish only accusative of feminine forms which is unique. An accusative ending might be attached if the participle is used as a nominal adjective form (typically in a verbal attribute syntactic function). See ex. in Tab. 20. (A similar situation occurs with nominal adjectives with the tag starting with AC.)

Wordform	Lemma	Tag
<i>už mám polévku uvařenu</i>	uvařit	VsFS4---X-APP--
<i>za Marie Terezie</i>	za	RR--2-----
<i>za týden</i>	za	RR--4-----

Table 20: Examples: Special cases of the CASE category application

## 5.6 Possessor's gender (6th position)

At the 6<sup>th</sup> position, the possessor's gender of possessive adjectives (with AU at POS and SUBPOS positions) and pronouns possessive for 3<sup>rd</sup> person (with P9 and P1) is captured. The values are in Table 21.

POSSGENDER	Description	Example
F	Feminine	<i>matčin, její</i>
M	Masculine animate	<i>otcův</i>
Z	{M, I, N} Not feminine	<i>jeho</i>
X	Any	<i>jejich</i>

Table 21: Values of POSSGENDER category (6th position)

Possessive adjectives and possessive pronouns (personal and relative) for the 3<sup>rd</sup> person refer to a possessor, and a possessive agreement rule applies for number and gender of the possessor (cf. Tab. 22). Though it could be argued (similarly as we did for the categories POS and SUBPOS), that this category is more lexically than inflectionally based, it has been decided to treat it as all the other morphological categories.

For possessive adjectives, given the nature of possible possessors, it can separately take on only the values of masculine animate and feminine. For possessive pronouns, both masculine genders and the neuter gender are always homonymous. Thus (cf. also the method used for the GENDER category) we do not use any of the three basic values (M, I, N) – instead, the value of Z is used.

If a possessive pronoun usage relates to the subject of the sentence (i.e. wordforms of *svůj*), then there is no gender distinction at all, and therefore the category POSSGENDER and also POSSNUMBER is not used; cf. Tab. 22.

Wordform	Lemma	Tag
<i>Petr opravil Janovi jeho auto.</i>	jeho	P9XXXZS3-----
<i>Petr opravil Janě její auto.</i>	jeho	P9NS4FS3-----
<i>Petr opravil jejich auto.</i>	jeho	P9XXXP3-----
<i>Petr opravil svoje auto.</i>	svůj-1	P8NS4-----

Table 22: Examples: Application of the POSSGENDER and POSSNUMBER categories

## 5.7 Possessor's number (7th position)

This category is closely related to the POSSGENDER category (cf. Sect. above). Similar reasons and rules apply here. This category is used only for possessive pronouns (personal and relative; with the values P9 and P1 at POS and SUBPOS positions). For possessive adjectives, it does not make sense to use it – there is no plural at least – the possessor, lexically present in the form, is always considered to be in singular. The values are in Table 23.

POSSNUMBER	Description	Example
P	Plural	<i>náš, váš, jejich</i>
S	Singular	<i>můj, tvůj, jeho</i>

Table 23: Values of POSSNUMBER category (7th position)

## 5.8 Person (8th position)

The PERSON category expresses the person of verb forms (if applicable), and person of personal pronouns. Personal pronouns, however, have a separate lemma for each person (and also for number; i.e. *já* for the 1<sup>st</sup> person singular, *my* for the 1<sup>st</sup> person plural, *ty* for the 2<sup>nd</sup> person singular, *vy* for the 2<sup>nd</sup> person plural and *on-1* for the 3<sup>rd</sup> person; see more in Sect. 11.1). This redundancy has been introduced for better human readability, since these lemmas are traditionally considered to be separate words). The values are in Table 24.

There is no attempt to assign person to all participles, transgressives, etc. For capturing compound wordfoms with the auxiliary verb *být* (“to be”), e.g. *přišels*, *bych*, *abys*, see Sect. 18.

PERSON	Description	Example
1	1 <sup>st</sup> person	<i>píšu</i> , <i>píšeme</i> , <i>my</i>
2	2 <sup>nd</sup> person	<i>píšeš</i> , <i>písete</i> , <i>ty</i>
3	3 <sup>rd</sup> person	<i>píše</i> , <i>píšou</i> , <i>ony</i>

Table 24: Values of PERSON category (8th position)

## 5.9 Tense (9th position)

The TENSE category belongs to verb forms only. The values are in Table 25.

Contrary to the traditional tense category assignment, tense is meant here in the purely morphological sense, without regard to the actual tense of an analytical verb form used within a particular sentence. Thus e.g. the verb form *pracoval* ‘(he) worked’ is assigned only the past tense value R, even though it can appear as a part of present conditional (*pracoval by* ‘(he) would work’ as well as the true past tense (or “perfective” tense, as Czech does not really distinguish these two): *pracoval jsem* ‘(I) (have) worked’.

In Czech, future tense is traditionally assigned also to present perfective verb forms. However, morphologically, these forms are simply present tense. On the contrary, there are real future tense forms (but only for a handful of verbs), such as *pojedu* ‘(I) will go’, which is a future tense form of the verb *jet* ‘to go’ (vs. *jedu* ‘(I) go’), and of course *budu* ‘(I) will’, *budeš* ‘(you) will’ ... – future forms of the verb *být* ‘to be’ used both in the existential sense as well as an auxiliary verb form for analytically expressed future tense of imperfective verbs.

TENSE	Description	Example
F	Future	<i>pojedeme</i> , <i>budu</i>
P	Present	<i>napíšu</i> , <i>píšu</i> , <i>jsem</i>
R	Past	<i>psal</i> , <i>byl</i>

Table 25: Values of TENSE category (9th position)

## 5.10 Degree of comparison (10th position)

The category of degrees of comparison is used for adjectives and adverbs. Traditionally, numbers are being used for the degrees of comparison; see Table 26.

The GRADE category is considered to be an inflectional category (as opposed to a derivation). This means that the wordforms of the same stem in positive, comparative and superlative (e.g. *velká* ‘big’, *větší* ‘bigger’, *největší* ‘biggest’) are understood as the wordforms of one paradigm. If the graded wordforms are not of the same stem (e.g. *dobrý* ‘good’ and *lepší* ‘better’), there are the separate paradigms represented by different lemmas. See examples in Tab. 27.

GRADE	Description	Example
1	Positive	<i>velká, pěkně</i>
2	Comparative	<i>větší, pekněji</i>
3	Superlative	<i>největší, nejpěkněji</i>

Table 26: Values of GRADE category (10th position)

For certain types of adjectives, such as possessive adjectives, this category is considered irrelevant (the value of Not applicable (–) is used). On the other hand, the value positive (1) is used for adjectives which do not form comparatives or superlatives only on semantic grounds (such as the often discussed adjective *optimální* ‘optimal’). In most cases, however, the morphological analysis would not reject even the comparative or superlative forms of such adjectives, marking their degree of comparison correctly and consistently according to their prefix and suffix.

Wordform	Lemma	Tag
<i>velký přítel</i>	<i>velký</i>	AAMS1----1A----
<i>nevětší přítel</i>	<i>velký</i>	AAMS1----3A----
<i>dobré podmínky</i>	<i>dobrý</i>	AAFP1----1A----
<i>lepší podmínky</i>	<i>lepší</i>	AAFP1----2A----
<i>nejlepší podmínky</i>	<i>lepší</i>	AAFP1----3A----

Table 27: Examples: Degree of comparison

## 5.11 Negation (11th position)

Similarly to the degree of comparison category, the NEGATION category is treated as a fully inflectional category, as negation is expressed in Czech by a prefix. Negation can be attached to verbs, adverbs, adjectives, and in principle, also to nouns (although only verbal nouns are typically used with negation regularly). The values are in Table 28. (See more in Sect. 12.)

NEGATION	Description	Example
A	Affirmative (not negated)	<i>velká, pěkně, přišel, ochota</i>
N	Negated	<i>nevelká, nepěkně, nepřišel, neochota</i>

Table 28: Values of NEGATION category (11th position)

## 5.12 Voice (12th position)

The VOICE category is used for verb forms only (mainly for verb participles). It is not used for verbal adjectives, even if they are derived from passive participle forms. The values are in Table 29.

## 5.13 Verbal aspect (13th position)

At the 13<sup>th</sup> position, the verbal ASPECT is coded. This category is used for verb forms only. It is not used for verbal adjectives and nouns. The values are in Table 30.

In fact, verbal aspect is rather lexical than morphological property but it is practical to keep it in the tags. Anyway, no paradigm/lemma is allowed to occur with two different verbal aspects in the accompanying tags.

VOICE	Description	Example
A	Active	<i>píšu, pojedu, psala</i>
P	Passive	<i>napsán</i>

Table 29: Values of VOICE category (12th position)

ASPECT	Description	Example
P	Perfect verb	<i>napsat</i>
I	Imperfect verb	<i>psát</i>
B	Biaspectual verb	<i>absolvovat</i>

Table 30: Values of ASPECT category (13th position)

## 5.14 Aggregate (14th position)

At the 14<sup>th</sup> position, so-called aggregates are coded. An aggregate is a wordform created by combining two or more forms into one and cannot be simply assigned any POS category (e.g. wordform *proň* consists of a pronoun *on* ‘he’ and the preposition *pro* ‘for’). The tag describes the main component of the aggregate (i.e. the pronoun *on*) and the joined components are coded at the 14<sup>th</sup> position of the tag (for joined preposition *pro*, there is value p). The values for capturing aggregates are in Table 31. For more information about aggregates, see Sect. 18.

AGGREGATE	Joined component	Description	Example
d	<i>do-</i>	preposition <i>do</i>	<i>doň</i>
n	<i>na-</i>	preposition <i>na</i>	<i>nač, načpak, naň</i>
o	<i>o-</i>	preposition <i>o</i>	<i>oč, očpak, oň</i>
p	<i>pro-</i>	preposition <i>pro</i>	<i>proň</i>
v	<i>ve-</i>	preposition <i>ve</i>	<i>več</i>
z	<i>za-</i>	preposition <i>za</i>	<i>zač, začpak, zaň</i>
D	<i>do- + -s</i>	prep. <i>do- + 2<sup>nd</sup></i> pers. of <i>být</i>	<i>doňs</i>
N	<i>na- + -s</i>	prep. <i>na- + 2<sup>nd</sup></i> pers. of <i>být</i>	<i>načs, načpaks, naňs</i>
O	<i>o- + -s</i>	prep. <i>o- + 2<sup>nd</sup></i> pers. of <i>být</i>	<i>očs, očpaks, oňs</i>
P	<i>pro- + -s</i>	prep. <i>pro- + 2<sup>nd</sup></i> pers. of <i>být</i>	<i>proňs</i>
Z	<i>za- + -s</i>	prep. <i>za- + 2<sup>nd</sup></i> pers. of <i>být</i>	<i>začs, začpaks, zaňs</i>
c	<i>-ch/sem</i>	1 <sup>st</sup> pers. sg. of cond. verb <i>být</i>	<i>bych, abych, kdybysem</i>
s	<i>-s</i>	2 <sup>nd</sup> pers. sg. of verb <i>být</i>	<i>příšels, kdyžs, kdybys</i>
m	<i>-chom/sme</i>	1 <sup>st</sup> pers. pl. of cond. verb <i>být</i>	<i>abychom, kdybysme</i>
e	<i>-ste</i>	2 <sup>nd</sup> pers. pl. of cond. verb <i>být</i>	<i>byste, abyste, kdybyste</i>

Table 31: Values of AGGREGATE category (14th position)

## 5.15 Variant, style, abbreviation (15th position)

Variant information (cf. values in Tab. 32) is used whenever all morphological categories together with a given lemma can be expressed by more than one wordform.

For example, Czech masculine nouns can typically have one or two variants in the singular genitive, singular locative, plural nominative, and plural locative cases (the presence or absence and the number of variants differs in different paradigms, and often is just lexically based (e.g. both *orli* and *orlové* (‘eagles’) are the wordforms of the noun *orel* (‘eagle’) and express plural masculine nominative)).

VAR	Description	Example
-	Not applicable – basic variant, standard contemporary style	<i>orlové, přijdeme, myslit</i>
1	Standard variant	<i>orli, myset</i>
2	Standard variant	<i>mysliti</i>
3	Standard variant	<i>mysleti</i>
4	Standard variant	<i>pomažemť</i>
5	Non-standard variant	<i>přídeme</i>
6	Non-standard variant	<i>přijdem</i>
7	Non-standard variant	<i>přídem</i>
8	Non-standard variant	<i>přijdeme</i>
9	Non-standard variant, misspelling	<i>příjdem</i>
b	Abbreviated form	<i>s (=sekunda)</i>
a	Other abbreviated form	<i>sec (=sekunda)</i>
c	Other abbreviated form	<i>sek (=sekunda)</i>

Table 32: Values of VAR category (15th position)

Another important distinction is a style of the form. Many standard noun, adjectival, pronominal, numeral as well as some verbal endings have their non-standard counterpart, and the morphology must be able to handle them properly, i.e. to distinguish them from their contemporary counterparts. Numbers 1 to 4 mark standard variants, while numbers 5 to 9 relate to non-standard ones. The 9 value is also for misspellings, typos or another distortions (see Sect. 20). More information about variants can be found in Sect. 9.2.

We have also added new values – letter b, c and a – to the 15<sup>th</sup> position of the tag for marking abbreviation of a (single) word which is captured as a special wordform of the paradigm of that word. For more information, see Sect. 14.

## 6 Stylistic Characteristics

Stylistic characteristics are captured to a limited extent. We used different means for capturing style of a paradigm and style of a particular wordform:

- Style label at AddInfo part of a lemma
- Variant Info at AddInfo part of a lemma
- Numbering of wordform variants at the 15<sup>th</sup> VAR position of tag

For instance, *šlojíř* is an archaic word meaning ‘veil’. Its lemma *šlojíř\_*,**a** bears the archaic style label **a** at AddInfo part (Sect. 4.2.3). For the stylistic and orthographic variants of a whole word (for the full-paradigm variants; see Sect. 9.1), we also state their basic word variant at AddInfo part (Sect. 4.2.4). For example, for the non-standard word *mlejn* ‘mill’, there is the standard variant *mlýn* ‘mill’. That fact is indicated at AddInfo part (*mlejn\_*,**h**<sup>~</sup>(^GC\*\*mlýn)). On the other hand, *lvové* ‘lions’ is an archaic wordform of a standard lemma *lev* ‘lion’. In this case, the variant of wordform is captured in the tag describing the wordform (there is the number 1 at the 15<sup>th</sup> position). For wordforms, standard and non-standard variants are only distinguished (not archaic, dialect, etc. see Sect. 5.15).

Rules for style labeling are described below and there is a separate Sect. 9 dedicated to orthographic and stylistic variants of the paradigms and wordforms.

Wordform	Lemma	Tag
<i>středověký šlojíř</i>	<i>šlojíř_</i> , <b>a</b>	NNIS1-----A----
<i>peknej mlejn</i>	<i>mlejn_</i> , <b>h</b> <sup>~</sup> (^GC**mlýn)	NNIS1-----A----
<i>lvové a psové</i>	<i>lev</i>	NNMP1-----A---1

Table 33: Examples: Capturing style

### 6.1 Style labeling

A word can be stylistically classified with style label attached to the lemma. Only one style label (or none, i.e. no more) is assigned to a lemma, the most general one applicable to most of contexts. Possible style labels attached to a lemma are:

**Ø – no style label.** No style flag is applicable for:

- standard, non-marked words, i.e. usable in neutral texts with neutral stylistic feature, they do not appear inappropriate or expressive.
- literary words of higher, bookish style, the expressions that may appear slightly peculiar in some texts but not completely improperly or incorrectly (e.g. *trýzen*, *blankyt*),
- words not yet incorporated into the language norm, occasionalisms, i.e. expressions well formed, meaning comprehensible, but authored for a particular occasion, unless they clearly belong to other group below (e.g. *blábolizmus*),
- special terms that cannot be considered non-standard, but they are marked because of their low frequency of use (e.g. *ikonostas*, *sakristie*),
- neologisms unless they clearly belong to other groups below (e.g. *europoslanec*).

**s – standard** style label. The style label with the value **s** is applied for:

- second (and possibly another) standard variant of the word (e.g. *kafeterie* - *kafetérie*, *kurs* - *kurz*). This style label is used only in connection with Variant Info (Sect. 4.2.4) when capturing orthographic and stylistic variants of word (see more in Sect. 9.1).

**a – archaic style label.** The style label with the value **a** is applied for:

- obsolete, archaic expressions (archaisms), which are not considered to be non-standard, but their use in neutral, standard texts is marked; they are almost with no evidence in a contemporary corpora, or they are only documented in historical texts; sometimes, they have been replaced by other (newer) words that are considered neutral (e.g. *uštipáček*, *servít*, *šlojíř*, *slout*),
- words written according to the original or older spelling standard; there have usually variants written according to the newer spelling standard (e.g. *thema*, *these*),
- words (so-called historisms) that denote objects that no longer exist; their use in a neutral text is marked and their meaning is unclear to the speaker without detailed knowledge of the area (e.g. words from the area of military terminology of the Austro-Hungarian Empire, although they are passively known from fiction: *maršrúta*, *vachman*).

**h – non-standard style label.** The general non-standard label with the value **h** is used if a non-standard word may be annotated with more stylistic features or if it cannot be classified with any of the specific labels below; typically:

- expressions from common Czech which are not considered a dialect (e.g. *vocelárna*, *mlejn*),
- univerbations, i.e. words formed from collocations (e.g. *pedák*, *hotelovka*, *tirák*, *stadiónovka*),

**n – dialect style label.** The style label with the value **n** is applied for:

- dialectal expressions; they might have marked phonic composition and/or there might be less awareness of the meaning among speakers, e.g. *šufánek*, *čučkař*.

**l – slang, argot style label.** the style label with the value **l** is applied for:

- non-standard expressions that are used by a very narrow group of people with the same professional or other special interests; e.g. *vertajmka*, *genštáb*.

**e – expressive style label.** The style label with the value **e** is applied for:

- expressions with positive and negative emotional feature evident from their formation; expressive words may not be considered non-standard, but their use in neutral texts is marked; they can be diminutives from expressions that have a different style label (e.g. word *sajtnička* with the style label **e** is derived from the word *sajtna* with the style label **1**).

**v – vulgar style label.** From expressive expressions we only separate vulgar words, the use of which may have a more serious impact. The style label with the value **v** is applied for:

- non-standard expressions that have a negative emotional feature, very often with obscene and/or vulgar content, they usually have an expressive phonic composition, semantically often refer to taboo areas (e.g. *vlezdoprdelizmus*, *čurák*).

**i – special style label.** The style label with the value **i** is applied for:

- distorted words, misspelling which are frequently used (e.g. *součastník*, *vánoce*, *intensívni*).

**Note.** The lemmas automatically derived from derivational patterns (with Derivation info; see Sect. 4.2.5) do not have any style label (see Sect. 7).

## 7 Derivative Relations

The dictionary primary handles inflection, not word-formation (for modeling word-formation relations in the lexicon of Czech, there is another project - Derinet<sup>17</sup>). However, as mentioned in Introduction (Sect. 2), a substantial part of the dictionary is mapped onto so-called derivational patterns in the source format. If a lemma belongs to a derivational patterns, several other lemmas/paradigms can be derived from it. All lemmas created by the derivational pattern have the derivational information stored in AddInfo part of the lemma (Sect. 4.2.5). In the current version of the dictionary, information on derivation is attached only to the lemmas created by the automatic procedure.<sup>18</sup>

The Derivation Info is primarily of technical or procedural nature: it carries information about the automatic creation of a lemma; the manifested word-formation relation may not be correct or complete.

Originating Lemma	→	Derivative Lemma
vychovat	→	vychovaný_^(*2t)
vychovat	→	vychovatelný_^(*4)
vychovat	→	vychování_^(*3at)
vychovat	→	vychovací_^(*2t)
vychovat	→	vychovavší_^(*3t)
vychovat	→	vychovávat_^(*4at)
vychovaný_^(*2t)	→	vychovanost_^(*3ý)
vychovaný_^(*2t)	→	vychovaně_^(*1ý)
vychovatelný_^(*4)	→	vychovatelnost_^(*3ý)
vychovatelný_^(*4)	→	vychovatelně_^(*1ý)
vychovávat_^(*4at)	→	vychovávání_^(*3at)
vychovávat_^(*4at)	→	vychovávací_^(*2t)
vychovávat_^(*4at)	→	vychovávající_^(*4t)
vychovávat_^(*4at)	→	vychovávaný_^(*2t)
vychovávat_^(*4at)	→	vychovávatelný_^(*4)
vychovávaný_^(*2t)	→	vychovávanost_^(*3ý)
vychovávaný_^(*2t)	→	vychovávaně_^(*1ý)
vychovatelný_^(*4)	→	vychovatelnost_^(*3ý)
vychovatelný_^(*4)	→	vychovatelně_^(*1ý)

Table 34: Example of automatic derivation

The derivational subsystem is a very useful, effective and economical tool that saves space in the source format. Using derivation patterns, 65 percent of lemmas/paradigms are created in the dictionary. For example, 19 paradigms are automatically derived from the verb *vychovat* ‘to bring up’ with lemma *vychovat*. See examples in Table 34. In the source format, these paradigms are written using only one line: *vychovat* *atd* = *vychovat*, where *atd* is a derivative pattern, the derivation rules of which is shown in Table 35.

On the other hand, the usage of derivation patterns has two important consequences for the structure of the dictionary:

- **Over-derivation.** Derivational patterns sometimes lead to only hypothetical derivatives, which are correctly formed but they are attested neither in the corpora nor in the dictionaries; e.g. lemma *obchodovávatelnost\_^(\*3ý)* with meaning possibility of repeatedly trade (‘retradability’), or hypothetical iterative adjective *převolovávatelný-2\_^(\*6-2)* derived from verb *převolovat-2-2\_^(přemíra.volū)* with meaning ‘to cause an excess number of oxen’.

<sup>17</sup><https://ufal.mff.cuni.cz/derinet>

<sup>18</sup>In the source format, there are also manually created comments on word-formation relation which did not make it into the currently described version for various reasons (inconsistency, incompleteness).

- **No manual labeling.** It is not possible to manually assign an additional information about the paradigm (i.e. AddInfo; e.g. value of term label or style labels) to automatically derived lemmas. During the automatic derivation process, derived lemmas can only automatically inherit values from originating lemma. See rules below.

```

atd 0,atd,r0,at,0,0,*
atd an,y,r0,aný,r0,at,-
atd ateln,y,r0,atelný,r0,at,-
atd ání,stn,r0,ání,r0,at,-
atd ac,i,r0,ací,r0,at,-
atd avš,iavv,r0,avší,r0,at,-
atd áv,atn,r0,ávat,r0,at,-
atd anost,kt1n,r0,anost,r0,aný,-
atd an,jev,r0,aně,r0,aný,-
atd atelnost,kt1n,r0,atelnost,r0,atelný,-
atd ateln,jev,r0,atelně,r0,atelný,-
atd ávání,stn,r0,ávání,r0,ávat,-
atd ávac,i,r0,ávací,r0,ávat,-
atd ávajíc,iavg,r0,ávající,r0,ávat,-
atd ávan,y,r0,ávaný,r0,ávat,-
atd ávateln,y,r0,ávatelný,r0,ávat,-
atd ávanost,kt1n,r0,ávanost,r0,ávaný,-
atd ávan,jev,r0,ávaně,r0,ávaný,-
atd ávatelnost,kt1n,r0,ávatelnost,r0,ávatelný,-
atd ávateln,jev,r0,ávatelně,r0,ávatelný,-

```

Table 35: Example of derivative pattern

## 7.1 Automatically derived lemmas

Automatically derived lemmas are identified by their derivation comment, which is created automatically. As mentioned above, automatically derived lemmas cannot be manually annotated. Values of AddInfo part of the lemma can be only automatically inherited from the original lemma. The rules are as follows. Automatically derived lemmas:

- **inherit reference** (Sect. 4.2.1). For example, from the base lemma *šest* ‘six’, the following derivatives are automatically derived: *šestkrát* ‘six times’, *šestina* ‘sixth’, *šestery* ‘sise’. The reference to the number is valid for each derived word.
- **inherit name labels** (Sect. 4.2.2). Any name label added to the original lemma always applies to the derived word. For example, from male surname (e.g. *Hromádka* ‘(Mr.) Hromádka’), female surname (e.g. *Hromádková* ‘(Mrs.) Hromádka’) and possessive adjective (e.g. *Hromádkův* ‘Hromádka’s’) are derived. The name label Y belongs to the original lemma and also to all its derivatives.
- **do not inherit style label** (Sect. 4.2.3). Derived words often do not have the same stylistic value as the original lemma, therefore derived lemmas do not inherit stylistic labels. For example, the word *výtržnost* ‘riot’ is automatically derived from the word *výtržný* ‘riotous’. While the word *výtržný* is an archaic word in the current Czech, the word *výtržnost* is not.
- **inherit variant info** (Sect. 4.2.4). An orthographic / phonetic change in the word stem also occurs in words derived from that word. Therefore, automatically derived words inherit Variant Info from the original lemma. The lemma of the basic variant in Variant Info is

automatically rewritten to the lemma of the corresponding derived variant. E.g. from the verb *votevřít* ‘to open-informal’ (which is a non-standard variant of verb *otevřít* ‘to open’), the adjective *votevřený* ‘open-informal’ and the noun *votevřenost* ‘openness-informal’ (non-standard variants of adjective *otevřený* ‘open’ and noun *otevřenost* ‘openness’, respectively) are derived.

- **do not inherit explanatory comment** (Sect. 4.2.6). Since there is no guarantee that the explanatory comment added to the original lemma is also valid for the lemma derived, during the automatic derivation derivative lemma does not inherit the explanatory comment.

See examples in Tab. 36.

Wordform	Lemma	Tag
<i>šest</i>	<i>šest'6</i>	Cn-S1-----
<i>šestina</i>	<i>šestina'6_^(*)3</i>	NNFS1----A---
<i>šestero</i>	<i>šestero'6_^(*)3</i>	CjNS1-----
<i>šestkrát</i>	<i>šestkrát'6_^(*)4</i>	Cv-----
<i>pan Hromádka</i>	<i>Hromádka-1;Y</i>	NNMS1----A---
<i>paní Hromádková</i>	<i>Hromádková-1;Y_^(*)5a-1</i>	NNFS1----A---
<i>Hromádkův pes</i>	<i>Hromádkův-1;Y_^(*)4a-1</i>	AUIS1M-----
<i>musíme votevřít</i>	<i>votevřít_;h_^(^GC**otevřít)</i>	Vf-----A-P--
<i>votevřený dveře</i>	<i>votevřený_^(^GC**otevřený)_(*3ít)</i>	AAFP1---1A---6
<i>votevřeně</i>	<i>votevřeně_^(^GC**otevřeně)_(*1ý)</i>	Dg-----1A----
<i>jeho votevřenost</i>	<i>votevřenost_^(^GC**otevřenost)_(*3ý)</i>	NNFS1----A----
<i>výtržný fanoušek</i>	<i>výtržný_,a</i>	AAMS1----1A----
<i>výtržnosti na zápase</i>	<i>výtržnost_^(*)3ý</i>	NNFP1----A----

Table 36: Examples: Automatic derivation

## 7.2 Derivative relation types

Here, we present several typical derivation patterns that are applied regularly automatically to a great number of original lemmas.

- to verbs:
  - iterative verbs with suffix *-[áié .. etc.]vat*
  - adjectives with suffix *-ný, -cí, -jící, -vaný, -vší, -telný, -vatelný*
  - nouns with suffix *-ní, -vání, -nost, -telnost, -vatelnost*
  - adverbs with suffix *-ně, -telně, -vatelně*
- to nouns:
  - possessive adjectives with suffix *-ův, -in*
  - female counterparts with suffix *-kyně, -ka* are derived from the masculine animate nouns with suffixes *-∅, -ec*
  - female surname are derived from some male surnames
- to adjectives:
  - nouns with suffix *-skost, -ckost* are derived from adjectives with suffix *-ský, -cký*

- adverbs with suffix *-ě* and nouns with suffix *-ost* are derived from adjectives with hard declension
- adverbs with suffix *-ě* are derived from adjectives with soft declension
- derivatives with suffix *-ina*, *-ery* and *-krát* are derived from basic numerals.

Examples of derivatives from verb are in Table 34 above. Examples of derivatives from nouns, adjectives and numerals are in Table 37.

Originating Lemma	→	Derivative Lemma
Milan_ ; Y	→	Milanův_ ; Y_-^(*2)
Eva_ ; Y	→	Evin_ ; Y_-^(*2a)
chodec	→	chodkyně_-^(*4ec)
doktor	→	doktorka_-^(*2)
Nový_ ; Y	→	Nová_ ; Y_-^(*1ý)
Konvička_ ; Y	→	Konvičková_ ; Y_-^(*3a)
světský	→	světskost_-^(*3ý)
vědecký	→	vědeckost_-^(*3ý)
barevný	→	barevně_-^(*1ý)
veliký	→	velikost_-^(*3ý)
letní	→	letně_-^(*1í)
sedm'7	→	sedmina'7_-^(*3)
sedm'7	→	sedmkrát'7_-^(*4)
sedm'7	→	sedmery'7_-^(*3)

Table 37: Examples: Derivative relation types

**Derivative collisions.** Sometimes, when deriving from several different originating lemmas, the morphologically same result is obtained. E.g. the same possessive adjective *Janův* ‘John’s’ is derived from the names *Jan*, *Jano* and *Janus*. Derivative collisions can also occur between a derived lemma and non-derived one (e.g., the female surname *Černá* is automatically derived from the male surname *Černý* with Y name label. There is also an identical paradigm for the name of the village *Černá* with the name label G).

In this case, there is only one lemma/paradigm in the dictionary (according to the principle of a unique paradigm; see Sect. 3.3). When generating the dictionary from the source format, these identical lemmas are merged into one. All derivative comments are stored in the AddInfo part of the lemma. The “path” to all originating lemmas is thus preserved. The merged lemma (like all derivatives) inherits name labels and variant infos and does not inherit the style labels and explanatory comments. Thus, all name labels and variant infos (from all originating lemmas) are stored in the AddInfo part of the merged lemma. See examples in Tab. 38.

Wordform	Lemma	Tag
<i>Jan a jeho, Janův sen</i>	Janův_ ; Y_-^(*2)_(*2o)_(*2us)	AUIS1M-----
<i>Jano a jeho, Janův sen</i>	Janův_ ; Y_-^(*2)_(*2o)_(*2us)	AUIS1M-----
<i>Janus a jeho, Janův sen</i>	Janův_ ; Y_-^(*2)_(*2o)_(*2us)	AUIS1M-----
<i>zboží dovážené z Německa</i>	dovážený_-^(*2t)_(*3it)	AANS1----1A----
<i>zboží dovážené na váze</i>	dovážený_-^(*2t)_(*3it)	AANS1----1A----
<i>paní Černá</i>	Černá_ ; G_ ; Y_-^(*1ý)	NNFS1-----A----
<i>Černá v Pošumaví</i>	Černá_ ; G_ ; Y_-^(*1ý)	NNFS1-----A----

Table 38: Examples: Derivative collisions

## 8 Semantic Description

The description of word meanings and other semantic distinctions is the least pursued goal in building the morphological dictionary and in morphological annotation. Compare Sect. 3 which describes the main principles of the dictionary. We capture primarily the formal morphological behavior of words regardless of semantic differences. We build a dictionary of wordforms, not a dictionary of words or meanings.

However, partial semantic description is contained in the following attributes and values (their description is in the relevant sections):

- Explanational comment at AddInfo part of the lemma (Sect. 4.2.6),
- Numeric reference at AddInfo part of the lemma (Sect. 4.2.1),
- Name label at AddInfo part of the lemma (Sect. 4.2.2),
- SUBPOS position of the tag, namely of pronouns and numbers (Sect. 5.2).

## 9 Orthographic and Stylistic Variants

There are two types of orthographic and stylistic variants:

- **full-paradigm variant**: an orthographic, phonetic or stylistic change applies to the whole paradigm (e.g. non-standard *mlejn* ‘mill’ vs. the standard *mlýn* ‘mill’),
- **wordform variant**: an orthographic, phonetic or stylistic change is manifested only in single wordforms (e.g. nominative singular of masculine inanimate *zelenej* ‘green’ vs. the standard *zelený*).

The full-paradigm variants are captured in the variant info field of the lemma (see more in Sect. 9.1), the wordform variants are indicated by the VAR position of the morphological tag (see more in Sect. 9.2).

### 9.1 Full-paradigm variants

If an orthographic, phonetic or stylistic change applies to the full paradigm, i.e. to all wordforms (cf. wordforms of *okno* ‘window’ in common Czech: *vokno*, *vokna*, *voknem* vs. in standard Czech: *okno*, *okna*, *oknem*), each set of wordforms is captured in separate paradigm with different lemma: we select one of the variants as “basic” (the standard one) and other variants (second standard, non-standard, archaic) refer to it in an additional descriptive element (variant info), attached to the lemma (Sect. 4.2.4). Non-basic paradigms (except derivatives, see below) also have a style label attached to the lemma (Sect. 4.2.3).

The basic variant is the one that is the least marked. Common alternations, that appear repeatedly, are solved the same way. For instance variants with the prosthetic consonant *v-* at the beginning of words have the basic variant without the prosthetic *v-* (e.g. *okno* - *vokno* ‘window’). Cf. examples in Tab. 39. Types of the most common full-paradigm variants are listed below in Sect. 9.1.1.

**Multiple variants.** In the case of multiple full-paradigm variants, one basic variant is chosen and other variants refer to it. Cf. example of variants of word *Afghánistán* in Tab. 39.

**Variants and derivation.** Paradigms derived automatically inherit variant info from the original lemma. This does not apply to style labels. The automatic derivatives have never any style label. Cf. examples of words *mlíkař* ‘milkman’ and *mlíkařka* ‘milkwoman’ in Tab. 39. See more information about derivatives in Sect. 7.

Wordform	Lemma	Tag
<i>okno</i>	<i>okno</i>	NNNS1----A----
<i>vokno</i>	<i>vokno_</i> , h <sup>^</sup> (^GC**okno)	NNNS1----A----
<i>mlíkař</i>	<i>mlíkař_</i> , h <sup>^</sup> (^GC**mlékař)	NNMS1----A----
<i>mlíkařka</i>	<i>mlíkařka</i> <sup>^</sup> (^GC**mlékařka)_(*2)	NNFS1----A----
<i>Afghánistán</i>	<i>Afghánistán_</i> ; G	NNIS1----A----
<i>Afganistan</i>	<i>Afganistan_</i> ; G, s <sup>^</sup> (^DD**Afghánistán)	NNIS1----A----
<i>Afgánistán</i>	<i>Afgánistán_</i> ; G, s <sup>^</sup> (^DD**Afghánistán)	NNIS1----A----
<i>Afghánistán</i>	<i>Afghánistán_</i> ; G, s <sup>^</sup> (^DD**Afghánistán)	NNIS1----A----
<i>Afganistán</i>	<i>Afganistán_</i> ; G, s <sup>^</sup> (^DD**Afghánistán)	NNIS1----A----
<i>Afghanistán</i>	<i>Afghanistán_</i> ; G, s <sup>^</sup> (^DD**Afghánistán)	NNIS1----A----
<i>Afghanistan</i>	<i>Afghanistan_</i> ; G, s <sup>^</sup> (^DD**Afghánistán)	NNIS1----A----
<i>Afghánistan</i>	<i>Afghánistan_</i> ; G, s <sup>^</sup> (^DD**Afghánistán)	NNIS1----A----
<i>Afgánistán</i>	<i>Afgánistán_</i> ; G, s <sup>^</sup> (^DD**Afghánistán)	NNIS1----A----

Table 39: Examples: Full-paradigm variants

### 9.1.1 Types of full-paradigm variants

There are a lot of types of variations in the Czech language. Alomorphs – both in the roots (*mlýn* – *mlejn*) and in prefixes and suffixes (*výlet* – *vejlet*) – are seen as cases of variation. Here, we list the most common ones. Examples of all listed types are in Table 40.

Wordform	Lemma	Tag
<i>kurs</i>	<i>kurs</i> _s_ <sup>^</sup> (^DD**kurz)	NNIS1----A----
<i>optimismus</i>	<i>optimismus</i> _s_ <sup>^</sup> (^DD**optimizmus)	NNIS1----A----
<i>citron</i>	<i>citron</i> _s_ <sup>^</sup> (^DD**citrón)	NNIS1----A----
<i>Abel</i>	<i>Abel</i> _s_ <sup>^</sup> (^DD**Ábel)	NNMS1----A----
<i>parfumérie</i>	<i>parfumérie</i> _s_ <sup>^</sup> (^DD**parfumerie)	NNFS1----A----
<i>archívní</i>	<i>archívni</i> _s_ <sup>^</sup> (^DD**archivní)	AAIS1----1A----
<i>príma</i>	<i>príma</i> -2_h_ <sup>^</sup> (^GC**prima-2)	Db-----
<i>mlejn</i>	<i>mlejn</i> _h_ <sup>^</sup> (^GC**plýtvat)	Vf-----A-I--
<i>plejtvat</i>	<i>plejtvat</i> _h_ <sup>^</sup> (^GC**mlýn)	NNIS1----A----
<i>vokolo</i>	<i>vokolo</i> -1_h_ <sup>^</sup> (^GC**okolo-1)	RR--2-----
<i>theorie</i>	<i>theorie</i> _a_ <sup>^</sup> (^DD**teorie)	NNFS1----A----
<i>dyž</i>	<i>dyž</i> _h_ <sup>^</sup> (^GC**když)	J,-----
<i>součastník</i>	<i>součastník</i> _i_ <sup>^</sup> (^DS**současník)	NNMS1----A----
<i>vánoce</i>	<i>vánoce</i> _i_ <sup>^</sup> (^DS**Vánoce)	NNMS1----A----
<i>management</i>	<i>management</i> _s_ <sup>^</sup> (^DD**manažment)	NNIS1----A----

Table 40: Examples: Full-paradigm variants

**Standard variants** of type DD in variant info. Style label of variants is most usually **s**, alternatively **a**.

**kurz** – **kurs**. Variants with *-z/s-* spelling. The variant with *-z-* is captured as a basic one. A frequent subtype of this alternation are variants with *-zmus/smus* spelling (e.g. *optimizmus* – *optimismus*).

**teorie** – **theorie**. Consonant variation *-t/th-*. The variant with *-t-* is captured as a basic one.

**citrón** – **citron**. Variation in vowel length. The variant with long vowel is captured as a basic one. This solution applies to most vowel-length variation cases but not in general (see the types below).

**parfumerie** – **parfumérie**. Vowel-length variation in the suffixes *-erie/érie*, *-iv(ní)/ív(ní)* (e.g. *archiv(ní)* – *archív(ní)*), *-ivum/ívum* (e.g. *pasivum* – *pasívum*), *-emie/émie* (e.g. *leukemie* – *leukémie*), *ped/péd* (e.g. *logoped* – *logopéd*). The short-vowel spelling is captured as a basic variant.

**manažment** – **management**. Variation between Czechized spelling and original foreign language spelling. More domesticated variant is captured as a basic standard one.

Foreign-language names are an area where there are usually a lot of different spellings in Czech texts. Determining the basic variant can be difficult. In complicated cases, we always just select one spelling as the basic; cf. variants of word *Afghánistán* in Tab. 39.

**Non-standard variants** of type GC in variant info. Style label of non-standard variants is most usually **h**, alternatively **n**, **e**, **l** or **v**. Standard variant is captured as a basic one.

**príma** – **príma**. Non-standard variation in vowel length.

**mlýn** – **mlejn**. Change from *-ý/í-* to non-standard *-ej-*.

**okolo – vokolo.** Addition of non-standard prosthetic consonant *v-* to the beginning of a word.

**když – dyž.** Consonant group reduction in non-standard expressions.

**Distortion variants** of type DS in variant info. Style label of non-basic variants is i.

**Vánoce – vánoce** Distorted words, words with outdated spelling or misspelling which are frequently used, intentional typos, phonetic transcription of words, etc. are captured as a variant of standard spelling variant.

## 9.2 Wordform variants

For capturing orthographic and stylistic variants of a wordform manifested usually in the ending (e.g. both *orli* and *orlové* ‘eagles’ are the wordforms of the noun *orel* ‘eagle’ and express plural masculine nominative), we use the the 15<sup>th</sup> position of the tag. No number at the the 15<sup>th</sup> position indicates the basic wordform. Numbers 1 to 4 mark standard variants (e.g. *orli* – *orlové* ‘eagles’), while numbers 5 to 9 relate to non-standard ones (e.g. *malý* – *malej* ‘little’; cf. examples in Tab. 41). The number 9 is also used for distortions, typos, misspellings (e.g. *o mě* ‘about me’). An overview of the values for 15<sup>th</sup> tag position see in Sect. 4.2.3.

Wordform	Lemma	Tag
<i>kde hnízdí orlové</i>	<i>orel</i>	NNMP1-----A----
<i>kde hnízdí orli</i>	<i>orel</i>	NNMP1-----A---1
<i>malý dům</i>	<i>malý</i>	AAIS1----1A----
<i>malej dům</i>	<i>malý</i>	AAIS1----1A---6
<i>přijdeme</i>	<i>přijít</i>	VB-P---1P-AAP--
<i>přídeme</i>	<i>přijít</i>	VB-P---1P-AAP-5
<i>příjdem</i>	<i>přijít</i>	VB-P---1P-AAP-6
<i>přýjdeme</i>	<i>přijít</i>	VB-P---1P-AAP-7
<i>přídem</i>	<i>přijít</i>	VB-P---1P-AAP-8
<i>přýdem</i>	<i>přijít</i>	VB-P---1P-AAP-9
<i>o mě</i>	<i>já</i>	PP-S6--1-----9

Table 41: Examples: Wordform variants

The main function of wordform numbering at the 15<sup>th</sup> tag position is to distinguish the forms according to the principle of unique analysis (see Sect. 3.1). Marking the stylistic value is secondary (if a form should have more standard variants and the values 1, 2, 3, 4 would not be enough, the values primarily designated for the non-standard variants could be used and vice versa. However, so far there was no case for which the values 1-9 were not enough; cf. example of *přijdeme* ‘we come’ in Tab. 41.).

## 9.3 Boderline cases of full-paradigm and wordform variants

In this section, boundary cases – where it is difficult to decide whether a particular case is a full-paradigm variant of or wordform variant – are described.

### 9.3.1 Variation is not in the full paradigm

Some variations are not applied throughout the full paradigm, but only in some wordforms (e.g. within the verb *mýt* ‘to wash’, the change *-ý* to *-ej* occurs only in the present tense wordforms and in the infinitive). We consider such cases to be wordform variants. It means that all variant wordforms are captured within one paradigm and they are distinguished at the 15<sup>th</sup> tag position.<sup>19</sup>

<sup>19</sup>The solution is based on the strong tradition, that lemma is represented by nominative for nouns and by infinitive for verbs, etc. In cases with no base form available, the variant wordforms are not perceived as wordforms

Wordform	Lemma	Tag
mýt	mýt	Vf-----A-I--
mejt	mýt	Vf-----A-I-6
myju	mýt	VB-S---1P-AAI--
myji	mýt	VB-S---1P-AAI-1
meju	mýt	VB-S---1P-AAI-3
být	být	Vf-----A-I--
nebejt	být	Vf-----N-I-6
jsem	být	VB-S---1P-AAI--
sem	být	VB-S---1P-AAI-6
myslit	myslit	Vf-----A-I--
myslet	myslit	Vf-----A-I-1
myslil	myslit	VpYS---R-AAI--
myslel	myslit	VpYS---R-AAI-1
začít	začít-1	Vf-----A-P--
začnout	začít-1	Vf-----A-P-1
přijal	přijmout	VpYS---R-AAP--
přijmul	přijmout	VpYS---R-AAP-1
svatější	svatý-1	AAIS1---2A----
světější	svatý-1	AAIS1---2A---1
světější	svatý-1	AAIS1---2A---1
rozepnuli	rozepnout	VpMP---R-AAP--
rozepjali	rozepnout	VpMP---R-AAP-1
rozpjali	rozepnout	VpMP---R-AAP-2
rozepli	rozepnout	VpMP---R-AAP-3

Table 42: Examples: Wordform variants

The main types are as follows (examples of all listed types are in Tab. 42)

**mýt – mejt.** Within the verb *mýt* ‘to wash’, the change *-ý* to *-ej-* (and similarly *-í* to *-ej-* within the verb *sít*) occurs only in the present tense wordforms and in the infinitive. Within the verb *být* ‘to be’ the alternation *-ý* to *-ej-* is only in the infinitive. There is only one paradigm with one lemma (the basic one: *mýt*, *sít*, *být* etc.) and all variant wordforms are marked at the 15<sup>th</sup> tag position.

Similarly: *omýt – omejt, rýt – rej; lít – lejt, nalít – nalejt*, etc.

**myslit – myslit.** Within the verb *myslit* ‘to think’, the suffix alternation *-e/i-* occurs only in the past participle wordforms and in the infinitive. There is only one paradigm with one lemma (with the basic one according to the Dictionary of Standard Czech by Czech Academy of Sciences, i.e *myslit*) and all variant wordforms are distinguished at the 15<sup>th</sup> tag position.

Similarly: *bydlet – bydlit, bulet – bulit; bystřet – bystřit; bělet – bělit*, etc.

**začít – začnout.** Within the verb *začít* ‘to begin’, there is a suffixal alternation *-∅/nu-* and it does not occur with all wordforms, so there is only one paradigm with one lemma and variant wordforms are distinguished at the 15<sup>th</sup> position of the tag.

Similarly: *přijal – přijmul; tnout – tít; nájedl – nájmul*, etc.

---

of two different paradigms. Variants *přijal – přijmul* (for which there is only one infinitive wordform *přijmout*) are of the same kind as variants *začal – začnul* (for which there are two variant infinitive wordforms *začít – začnout*). Both *přijal – přijmul* and *začal – začnul* are captured in the same way. Dividing the variant forms to different paradigms (for example with numbered lemmas *přijmout-1* for the complete paradigm with forms *přijal*, etc. and *přijmout-2..h* for the incomplete paradigm only including the non-standard wordforms; *přijmul*) would be very difficult, maybe impossible due to large amount of variant forms; cf. wordforms of *rozepnout* in Table 42.

**zmrazen – zmražen.** Within the verb *zmrazit* ‘to freeze’, there is an alternation *-z/z-* and it does not occur with all wordforms (it occurs in the passive participle), so there is only one paradigm with one lemma and the wordforms are distinguished at the 15<sup>th</sup> tag position.

**svatější – světější.** The root alternation *svat/svět* is only present in comparative and superlative. There is only one paradigm with one lemma (the basic one: *svatý*) and the variant wordforms are marked at the 15<sup>th</sup> position of the tag.

Similarly: *bělejší – bílejší*.

### 9.3.2 Variation in the base form

If the wordform that is used as the lemma (i.e. mainly nominative, infinitive) has more variants in the respective paradigm, one of the variant wordforms, the one that is least marked, is chosen as the lemma and other forms are captured as variants at the 15<sup>th</sup> tag position.

Wordform	Lemma	Tag
<i>pracovat</i>	<i>pracovat</i>	Vf-----A-I--
<i>pracovati</i>	<i>pracovat</i>	Vf-----A-I-2
<i>péci</i>	<i>péci</i>	Vf-----A-I--
<i>péct</i>	<i>péci</i>	Vf-----A-I-1
<i>nit</i>	<i>nit</i>	NNFS1----A---
<i>nitč</i>	<i>nit</i>	NNFS1----A---6
<i>konsenzus</i>	<i>konsenzus</i>	NNIS1----A---
<i>konsenz</i>	<i>konsenzus</i>	NNIS1----A---1
<i>virem</i>	<i>vir-1</i>	NNIS7----A---
<i>virusem</i>	<i>vir-1</i>	NNIS7----A---1
<i>Patricie</i>	<i>Patricie_</i> ;Y	NNFS1----A---
<i>Patricia</i>	<i>Patricie_</i> ;Y	NNFS1----A---1
<i>Avia</i>	<i>Avia_</i> ;m_^(vozidlo)	NNFS1----A---
<i>Avie</i>	<i>Avia_</i> ;m_^(vozidlo)	NNFS1----A---1
<i>Polanský</i>	<i>Polanský_</i> ;Y	NNMS1----A---
<i>Polanski</i>	<i>Polanský_</i> ;Y	NNMS1----A---1
<i>Polanskij</i>	<i>Polanský_</i> ;Y	NNMS1----A---2
<i>Rucký</i>	<i>Rucký_</i> ;Y	NNMS1----A---
<i>Ruckij</i>	<i>Rucký_</i> ;Y	NNMS1----A---1
<i>Ruckoj</i>	<i>Rucký_</i> ;Y	NNMS1----A---2
<i>Ruckej</i>	<i>Rucký_</i> ;Y	NNMS1----A---6
<i>nejvíce</i>	<i>více</i>	Dg-----3A---
<i>nejvíc</i>	<i>více</i>	Dg-----3A---1
<i>zase</i>	<i>zase-1</i>	Db-----
<i>zas</i>	<i>zas-1_</i> ;s^(^DD**zase-1)	Db-----

Table 43: Examples: Wordform variants

The main types of this case are as follows (examples of all listed types are in Tab. 43):

**pracovat – pracovati.** There is *-t/ti* alternation in most infinitive forms. The *-t* variant is the basic one, it is captured as a lemma. Infinitive ending with *-ti* (which is less common, slowly disappearing, bookish, archaic) usually has value 2 at the 15<sup>th</sup> tag position.

**péci – péct.** Another alternation in infinitive forms is *-ci/t* type. Here, infinitive ending with *-ci* is captured as a base form.

**niť – nit.** A noun can have multiple variants of the singular nominative (e.g. *nit* vs. *nitč*). We mark variants at the 15<sup>th</sup> position of the tag and one of them is used for the lemma.

**konsenzus – konsenz.** The masculine inanimate nouns with -∅/us alternation in ending of singular nominative and singular accusative (and sometimes also in all cases, e.g. *vir* – *virus*) are captured as one paradigm. The masculine animate nouns with -∅/us alternation (e.g. *Josephus* – *Joseph*) are captured as the two separate paradigms/lemmas.

**Patricie – Patricia.** There is an alternation of endings -ia/ie in the loan female personal names (e.g. *Patricie* vs. *Patricia*) and other names (e.g. *Istrie* vs. *Istria*). Nouns with variant nominative ending -ia/ie are represented in one paradigm with lemma ending -ie. The -ia form is only in nominative (with value 1 at the 15<sup>th</sup> position of the tag). In some cases, however, the basic variant is the one with -ia ending (e.g. *Avia* vs. *Avie*).

**Polanský – Polanski – Polanskij.** Great variation in the endings also occurs in the transcription of Polish and Russian surnames into Czech, e.g. *Polanský* vs. *Polanski* vs. *Polanskij* or *Rucký* vs. *Ruckij* vs. *Ruckoj* vs. *Ruckej*. A variant with a Czech ending -ý is captured as the lemma.

**více – víc.** When it comes to inflexible words, the distinction between morphological, spelling and word-forming variants is not very clear. Close, similar words with a flexible POS (e.g. variants of multiple numbers *dvakrát* vs. *dvakráte*) and similar words that compare (e.g. variants of comparative adverb *nejvíce* vs. *nejvíc*) are captured as wordform variants (they are distinguished at 15<sup>th</sup> tag position). In other cases, for non-comparative adverbs (e.g. *zase* vs. *zas*), variants are captured as full-paradigm (using a variant reference in AddInfo part of lemma). Formally or semantically distant words (e.g. *jen* vs. *jenom*, *aspoň* vs. *alespoň*) are not interlinked at all.

## 10 Part of Speech Determination (problematic cases)

### 10.1 Part of speech of inflexible words

The determination of traditional inflexible POS does not depend on morphological properties. Distinguishing some adverbs, particles, prepositions, conjunctions, and interjections is based on the syntactic function of a given word in a sentence.

Some inflexible words can perform various functions in a sentence; e.g. the word *tak* ‘so’ can be a conjunction (*pršelo, tak nešel* ‘it was raining, so he didn’t go’), a particle (*tak už jdeme* ‘so here we go’) or an adverb (*uděláme to tak* ‘we’ll do it that way’). In some cases, POS of inflexible words (like *tak* ‘so’, *totiž* ‘that-is’, *však* ‘but’, *vůbec* ‘not-at-all’, *jen* ‘only’, etc.) is difficult to determine. The decision on POS is left to the annotator in problematic cases. See examples in Tab. 44.

Wordform	Lemma	Tag
<i>tak už jdeme</i>	tak-1	TT-----
<i>pršelo, tak nešel</i>	tak-2	J^-----
<i>uděláme to tak</i>	tak-3	Db-----
<i>legendy rocku</i>	rock-1	NNIS2----A---
<i>rock festival</i>	rock-2	AAXXX----1A---
<i>rock’n’roll</i>	rock-3	S2-----A---
<i>skladba We will rock you</i>	rock-77	F%-----
<i>cyklo-oblečení</i>	cyklo-1	S2-----A---
<i>cyklo výlet</i>	cyklo-2	AAXXX----1A---
<i>pro příznivce cyklo</i>	cyklo-3	NNXXX----A---
<i>v Sazka Aréně</i>	Sazka_;	NNFS1----A---
<i>m</i>		

Table 44: Examples: Part of speech of inflexible words

Similarly, according to a syntactic function of a word in a sentence, we also determine the part of speech of domesticated loanwords such as *online*, *cyklo* ‘cyclo’, *rock* (see also Sect. 17.3), and of inflexible words in general. E.g. the loanword *rock* has four lemmas in the dictionary that capture the usage of this word in various syntactic function (as we see in the Tab. 44, the loanword *rock* can be inflected if it is a noun).

**Note.** Commonly inflected names in an attributive position (e.g. *Gambrinus* in *Gambrinus liga* ‘Gambrinus league’, *Sazka* in *Sazka Aréna* ‘Sazka Arena’) have a noun (not an adjective) tag. Cf. Tab. 44.

#### 10.1.1 Frozen wordforms (*krážem*, *bycha*, *domácku*)

In the case of frozen wordforms of words that occur in only one type of collocation (such as *krážem* ‘cross’ in idiom *křížem krážem* ‘cross by cross’; *bycha* in idiom *pozdě bycha honit*; *domácku* ‘home’ in *po domácku* ‘home-made’), the lemma of paradigms equals the frozen form, and there is only a single form in the paradigm. See examples in Tab. 45.

Wordform	Lemma	Tag
<i>křížem krážem</i>	krážem	Db-----
<i>pozdě bycha honit</i>	bycha	NNIS4----A---
<i>po domácku</i>	domácku	NNNS6----A---

Table 45: Examples: Frozen wordforms

## 10.2 Nouns from adjectives

In accordance with the tradition, we capture the so-called substantiated adjectives as nouns (the value of the SUB/POS tag position is NN), although their inflection is adjectival and they fulfill the role of the noun mainly syntactically. Due to the large number of unclear cases, only a limited number of cases are captured in this way (i.e. clear historical examples *vrátný* ‘porter’, *průvodčí* ‘conductor’ which in contemporary language already appear only as nouns; further *cestující* ‘traveller’, *zemřelý* ‘dead’, *milá* ‘girlfriend’, *nevidomí* ‘the-blind’, but for example not *zaměstnaný* ‘employee’, *milovaný* ‘beloved’, *dojízdějící* ‘commuting’, *místní* ‘local’, *nastávající* ‘groom-to-be’.)

Cf. examples in Tab. 46.

Wordform	Lemma	Tag
<i>náš vrátný</i>	<i>vrátný-1_</i> <sup>^</sup> ( <i>osoba</i> )	NNMS1-----A----
<i>to je moje milá</i>	<i>milá-2_</i> <sup>^</sup> (*3ý-2)	NNFS1-----A----
<i>moje milá žena</i>	<i>milý-1_</i> <sup>^</sup> ( <i>příjemný</i> )	AAFS1-----1A----
<i>nevidomí maséři</i>	<i>nevidomý-1</i>	AAMP1-----1A----
<i>telefon pro nevidomé</i>	<i>nevidomý-2</i>	NNMP4-----A----
<i>lidé dojízdějící do zaměstnání</i>	<i>dójízdějící_</i> <sup>^</sup> (*4t)	AGMP1-----A----
<i>zpráva pro dojízdějící</i>	<i>dójízdějící_</i> <sup>^</sup> (*4t)	AGMP4-----A----

Table 46: Examples: Nouns from adjectives

## 10.3 Part of speech of predicatives (words with suffix *-o*)

Deadjectival and deverbal words with suffix/ending *-o* (e.g. *teplo* ‘warm’, *zataženo* ‘cloudy’) which serve as a predicative in a sentence are captured as adverbs. However, these words can be homonymous with nouns and also with nominal forms of adjectives / passive participles and prefixal segments. Cf. Tab. 47.

In some contexts, it can be difficult to decide on morphological analysis. Here are guidelines for determining the part of speech of these derivatives.

Wordform	Lemma	Tag
<i>období tepla a sucha</i>	<i>teplo-1</i>	NNNS2-----A----
<i>je velmi teplo</i>	<i>teplo-2_</i> <sup>^</sup> ( <i>být_někomu_teplo</i> )	Dg-----1A----
<i>teplo-vodní topení</i>	<i>teplo-3</i>	S2-----A----
<i>je tu docela hořko</i>	<i>hořko-1</i>	Db-----
<i>nálada přesla v hořko</i>	<i>hořko-2</i>	NNNS4-----A----
<i>hořko-sladké vzpomínky</i>	<i>hořko-3</i>	S2-----A----
<i>je zataženo</i>	<i>zataženo-2_</i> <sup>^</sup> ( <i>být_zataženo</i> )	Dg-----1A----
<i>fotím i v zataženu</i>	<i>zataženo-1</i>	NNNS6-----A----
<i>obloha je zatažena</i>	<i>zatáhnout</i>	VsQW----X-APP--
<i>bylo tam prázdro</i>	<i>prázdro-1</i>	Dg-----1A----
<i>má v hlavě prázdro</i>	<i>prázdro-2</i>	NNNS4-----A----
<i>síň byly pusty a prázdn</i>	<i>prázdný</i>	ACTP-----A----

Table 47: Examples: Words with suffix *-o*

**Adverb.** As adverb, we capture the derivatives with *-o* in predicative and adverbial positions, unless they follow a preposition. Cf. examples:

Dg *Je velmi teplo.* ‘It’s very warm.’; *Bylo mi teplo.* ‘I was warm.’; *Máte tu lacino.* ‘You’re cheap here.’; *Je zataženo a bude zataženěji.* ‘It is cloudy and it will be more cloudy.’

**Db** *Země medu, v níž je docela hořko.* ‘A land of honey in which it is quite bitter.’; *známe se krátko* ‘We know each other for a short time.’; *draho prodat* ‘sell dearly’.

Deciding whether the adverb is of the Dg type (i.e. it forms a comparative and superlative) or whether it does not form degrees of comparision (Db type; cf. Sect. 11.3) is relatively difficult. In addition, in many cases, there is competition between adverbs ending in *-o* and adverbs ending in *-e*. We evaluate wordforms of comparative and superlative preferably as adverbs ending in *-e*. Cf. examples in Tab. 48.

Wordform	Lemma	Tag
<i>bylo mu lehko po těle</i>	lehko-1	Dg-----1A----
<i>bylo mu lehce u srdce</i>	lehce	Dg-----1A----
<i>to se lehko řekne</i>	lehko-1	Dg-----1A----
<i>to se lehce řekne</i>	lehce	Dg-----1A----
<i>to se lehčejí řekne</i>	lehce	Dg-----2A----
<i>cestuje na lehko</i>	lehko-2	NNNS4----A----
<i>lacino získat</i>	lacino-1	Dg-----1A----
<i>lacině získat</i>	lacině_^(*)1ý)	Dg-----1A----
<i>lacinějí získat</i>	lacině_^(*)1ý)	Dg-----2A----
<i>prodává dost nelacino</i>	lacino-1	Dg-----1N----
<i>draho prodat</i>	draho	Db-----
<i>draze prodat</i>	draze	Dg-----1A----
<i>prodat nejdráž</i>	draze	Dg-----3A----

Table 48: Examples: Adverbs ending in *-o* vs. adverbs ending in *-e*

**Noun.** As nouns with neutrum gender, we capture the words ending *-o* if they are in subject or object position, if there is an adjectival or pronominal modifier or if they follow a preposition (and are inflected). Cf. examples:

NNN *Je velké teplo.* ‘It’s very hot.’; *To jsou velká tepla.* ‘These are the heat.’; *Pojďme do tepla.* ‘Let’s get warm.’; *Dobrá nálada přešla v hořko.* ‘The good mood turned bitter.’; *Cítí to úzko, jež ho obcházelo.* ‘He felt the tightness around him.’; *Jsem právě naladěna na něžno a milo.* ‘I’m just in the mood for being gentle and kind.’; *Prodat za lacino.* ‘Sell for cheap.’; *Fotím i v zataženu.* ‘I take photos even in the cloudy.’

**Nominal adjective or passive participle.** If a derivative has an adjectival agreement with a governing noun (usually in a subject position), it is a short form of the adjective or a passive participle. Cf. examples:

AC *Pusty, prázdný byly síně ostatní.* ‘The other halls were desolate and empty’; *Je daleka toho, aby mu šla opět zachraňovat život.* ‘She is far from saving his life again.’; *Žena ležela na nosítkách blízka kómatu.* ‘The woman lay on a stretcher near a coma.’

Vs *Obloha je zatažena a zamračena.* ‘The sky is cloudy and overcast.’

In this case, the word with *-o* is a wordform of respective adjective or verb. Therefore, in the case of adjectival nominal wordfom, the lemma is a respective long adjective. In the case of passive participle, the lemma is an infinitive of the respective verb.

**Prefixal segment.** Derivatives with suffix *-o* which are part of hyphenated composites (Sect. 19) are captured as a prefixal segment (Sect. 16).

S2 *černo-bílý svět* ‘black-and-white world’; *hořko-sladké vzpomínky* ‘bitter-sweet memories’

## 11 Detailed Part of Speech

In this section we describe criteria according to which the words within the given POS are classified into subtypes (described by a value on the 2<sup>nd</sup> tag position - SUBPOS). We describe subtypes for pronouns (Sect. 11.1), numerals (Sect. 11.2) and adverbs (Sect. 11.3).

### 11.1 Subtypes of pronouns

We identify several features that can serve as criteria for dividing pronouns into various subtypes:

- morphological behavior
- semantic function
- possession
- reflexivity
- clitichood

So called agreement gender and semantic function are chosen as the main criteria. The criterion of agreement gender affects which morphological categories (especially GENDER and NUMBER) are relevant to determine for a given pronoun. Semantic classification is based on the traditional division of pronouns into personal, indefinite, demonstrative, negative, etc. The subtypes of pronouns are described by a value on the 2<sup>nd</sup> tag position - SUBPOS, see overview in Sect. 5.2 and Tab. 49 here.

Type & Subtype	Gender	No Gender
<b>Personal</b>		
- Reflexive	PE <i>on, něj</i> ; clitic P5 <i>mu</i>	PP <i>já, ty</i> ; clitic PH <i>mi</i>
- Possessive	PS <i>můj, nás</i>	P6 <i>sebe</i> ; clitic P7 <i>se, si</i>
- Possessive, 3 <sup>rd</sup> pers.	P9 <i>jeho, její, jejich</i>	-
- Reflexive possessive	P8 <i>svůj</i>	-
<b>Relative</b>	P4 <i>který, jaký, čí, jenž</i>	PQ <i>kdo, co, kdož, copak</i>
- Possessive	P1 <i>jehož, jejíž, jejichž</i>	-
<b>Indefinite</b>	PZ <i>nějaký, čísi, sotvačterý</i>	PK <i>někdo, cosi, nevímco</i>
<b>Negative</b>	PW <i>nijaký, ničí, žádný</i>	PY <i>nikdo, nic</i>
<b>Demonstrative</b>	PD <i>ten, tentýž, takový</i>	-
<b>Delimiting</b>	PL <i>všechn, sám, veškerý</i>	-

Table 49: Subtypes of Pronouns

**Morphology.** The morphological criterion divides pronouns into two groups: gender pronouns and no gender pronouns.

**Gender pronouns** express variable values of the gender (and also number) depending on the gender (and number) of the governing noun (cf. *nějaký dům* ‘some house’ (masc. sg), *nějaká žena* ‘some woman’ (fem. sg.), *nějaké dítě* ‘some child’ (neut. sg.), *nějaké domy* ‘some houses’ (masc. pl.)) or according to the sense (*on* ‘he’, *ona* ‘she’, *ono* ‘it’, *oni* ‘they’). All forms are represented by one lemma (nom. sg. masc. anim.); similarly to adjectives. The GENDER and NUMBER tag positions are filled. Cf. examples in Tab. 50.

**No gender pronouns** are pronouns with no gender and number variation (e.g. *někdo, kdeco*; cf. Tab. 50). The GENDER and NUMBER tag positions are not filled (although we are aware that most of them could be classified as masculine (e.g. *kdo* and other various personal pronouns) or neutrum (*co* and other various non-personal pronouns)). No gender pronouns behave as syntactic nouns in sentences.

Wordform	Lemma	Tag
<u>nějaká</u> žena	nějaký	PZFS1-----
<u>někdo</u> tu je	někdo	PK--1-----
<u>mně</u> to dal	já	PP-S3--1-----
dal <u>mi</u> to	já	PH-S3--1-----
<u>ona</u> přišla	on-1	PEFS1--3-----
<u>jeho</u> si vážím	on-1	PEYS4--3-----
vážím si <u>ho</u>	on-1	P5ZS4--3-----
bojím se o <u>něho</u>	on-1	PEZS4--3----1
stará se o <u>sebe</u>	se_~(zvr._zájmeno/částice)	P6--4-----
nestarej <u>se</u>	se_~(zvr._zájmeno/částice)	P7--4-----
<u>můj</u> dům	můj	PSYS1-S1-----
<u>váš</u> dům	váš	PSYS1-P2-----
<u>svůj</u> dům	svůj-1	P8YS1-----
do <u>jeho</u> domu	jeho	P9XXXZS3-----
do <u>jejího</u> domu	jeho	P9ZS2FS3-----
do <u>její</u> chalupy	jeho	P9FXXFS3-----
<u>jejich</u> dům	jeho	P9XXXXP3-----
<u>žena</u> , <u>jejíž</u> dům vidím	jehož	P1IS4FS3-----

Table 50: Examples: Subtypes of pronouns

**Semantic function.** We classify six main semantic groups of pronouns, largely following the Czech grammar tradition. In each group, we use the unique SUBPOS value to distinguish between gender and no gender pronouns.

- **personal:**

PE gender (*on* ‘he’, *ona* ‘she’, *ono* ‘it’, *jím* ‘them’)  
 PP no gender (*já* ‘I’, *ty* ‘you’, *my* ‘we’, *vám* ‘to you’)

- **relative/interrogative:**

P4 gender (*jaký* ‘what’, *který* ‘which’, *čí* ‘whose’, *jenž* ‘who’, *kterýžto* ‘which’)  
 PQ no gender (*kdo* ‘who’, *co* ‘what’, *cožpak* ‘isn’t-it-true-that’, *kdož* ‘who’)

- **indefinite:**

PZ gender (*nějaký* ‘some’, *čísi* ‘somebody’s’, *číkolí* ‘anybody’s’, *sotvakterý* ‘hardly-some’)  
 PK no gender (*někdo* ‘somebody’, *bůhvíkdo* ‘whoever’, *cosi* ‘something’)

- **negative:**

PW gender (*nijaký* ‘no/none’, *ničí* ‘nobody’s’, *žádný* ‘no/none’)  
 PY no gender (*nic* ‘nothing’, *nikdo* ‘nobody’)

- **demonstrative:**

PD gender (*ten* ‘this’, *tamten* ‘that’, *onen* ‘that-over-there’, *tentýž* ‘same’, *takový* ‘such’)

- **delimiting:**

PL gender (*všechnen* ‘all’, *sám* ‘alone’, *veškerý* ‘whole’)

Unlike Czech grammars, we do not distinguish interrogative as a separate subtype because of its unclear distinction from the relative type (i.e. its enormous homonymy).

Other values of the **SUBPOS** category distinguish possessive, reflexive and clitic pronouns (see below).

**Clitichood.** Several personal pronouns have a clitic (short) form. They have a special value of the **SUBPOS**:

P5 personal gender clitic (e.g. *mu, ho* ‘him’);

PH personal no gender clitic (e.g. *mě, mi* ‘me’, *ti* ‘you’);

P7 personal (no gender) reflexive clitic (e.g. *se, si*).

Distinguishing clitic forms at **SUBPOS** tag position violates the principle that the 2<sup>nd</sup> position is the same for the whole paradigm (cf. in Sect. 3.2).

The pronouns in the forms requested after any preposition (with prefix *n-*: *něj, něho* ‘him’) are not distinguished in the **SUBPOS** position; these wordfoms are distinguished at the 15<sup>th</sup> position of the tag (cf. in Tab. 50).

**Reflexivity and Possession.** Futher subtypes of personal (and also relative) pronouns are introduced based on the feature of reflexivity and possession:

P8 personal possessive pronouns (e.g. *můj* ‘my’, *tvůj* ‘your’, *náš* ‘our’ and *váš*);

P6 personal reflexive pronoun (long forms; i.e. *sebe, sobě*);

P8 personal possessive reflexive pronoun (i.e. *svůj* ‘my’ / ‘your’ / ‘her’ / ‘his’).

Furthermore, two groups are specially divided:

P9 personal possessive pronouns for the 3<sup>rd</sup> person (e.g. *jeho* ‘his’, *její* ‘her’, *jejich* ‘their’);

P1 relative possessive pronouns (*jehož, jejíž* ‘whose’).

Except the gender and number of an object, these pronouns (P9 and P1) express also the gender and number of a subject, i.e. possessor (e.g. *jeho dům* ‘his house’ (masc. sg., possessor: masc. sg.), *do jejího domu* ‘to her house’ (masc. sg., possessor: fem. sg.), *do její chalupy* ‘to her cottage’ (fem. sg. possesor: fem. sg.), *jejich dům* ‘their house’ (masc. sg., possessor: pl.). That’s why the **POSGENDER** and **POSSNUMBER** tag positions are also filled here. Cf. examples in Tab. 50.

## 11.2 Subtypes of numerals

Within the category of numerals, words are associated mainly on the basis of numerical meaning. From a formal point of view, there are words that behave as nouns (*nula* ‘zero’, *milion* ‘million’), pronouns (*dvě* ‘two’, *jedny* ‘one-kind’), adjectives (*první* ‘first’, *paterý* ‘five-kinds’, *několikatý* ‘umpteenth’) and adverbs (*třikrát* ‘three-times’). We respect the traditional semantic concept of the POS of numerals, however, the formal behavior is the basic for the subdivision of numerals into individual subtypes.

Numerals (except the numbers written by arabic or roman numeral symbols) are distinguisthed into several subtypes (captured at the 2<sup>nd</sup> tag position – **SUBPOS**, see overview in Sect. 5.2 and Tab. 51 here) according various combination of the two basic features:

- morphological behavior
- semantics including definiteness

Type & Subtype		<b>Adjectival</b>	<b>Non-adjectival</b>
<b>Cardinal</b>	def.	Cn <i>jeden, dva, oba</i>	C1 <i>tři, čtyři, pět, půl</i> Cz <i>sto, tisíc, miliarda, nula</i> Ca <i>mnoho, málo, t/kolik, pár</i>
	indef.	Cy <i>nejeden</i>	-
<b>Ordinal</b>	def.	Cr <i>první, druhý, šestsetdruhý</i>	-
	indef.	Cw <i>k/tolikátý, bůhvíkolikátý</i>	-
<b>Multiplicative</b>	def.	-	Cv <i>třikrát/e</i>
	indef.	-	Co <i>k/tolikrát/e, několikrát</i>
<b>Generic</b>	def.	Cd <i>jedny, dvoje/í, patery/y, obé</i>	Cj <i>patero, dvé, tré</i>
	indef.	Ch <i>nejedny, několikerý/y</i>	Ck <i>několikero</i>
<b>Arabic numerals</b>		C= 1.24	
<b>Roman numerals</b>		C} MXV	

Table 51: Subtypes of Numerals

**Morphology.** The morphological criterion is primary. We distinguish between numerals with adjectival declension and numerals with another declension.

The **adjectival numerals**, i.e. with agreement gender, express variable values of the gender (and also number) depending on the gender (and number) of the governing noun (e.g. *jeden člověk* ‘one man’ (masc. sg.), *jedna žena* ‘one woman’ (fem. sg.), *jedno dítě* ‘one child’ (neut. sg.) or *několikátý problém* ‘multiple problem’ (masc. sg.), *několikáte problémy* ‘multiple problems’ (masc. pl.)). All wordforms are represented by one lemma (nom. sg. masc. anim.); similarly to adjectives. However, there are not degrees of comparison in adjectival declension of numerals, and some paradigms have only singular (e.g. *jeden* ‘one’) or only plural wordforms (*dva* ‘two’). In the tag, the GENDER, CASE and NUMBER positions are filled. Cf. examples in Table 52.

There are several types of **non-adjectival numerals**. First, their morphological behavior is similar to that of nouns (e.g. *sto* ‘hundred’, *nula* ‘zero’, *patero* ‘five-kinds-of’). The position of GENDER is filled (one value for the whole paradigm).

In the other types, the morphological behavior is quite specific and expresses no gender; they can be inflected (e.g. *tři domy* ‘three houses’ (Nom.), *bez tří domů* ‘without three houses’ (Gen.), *mnoho domů* ‘many houses’ (Nom.), *bez mnoha domů* ‘without many houses’ (Gen.)), or they are inflexible (e.g. *několikrát* ‘several-times’ with only variant wordform *několikráte* ‘several-times’). The tag position of GENDER is not filled. Other tag positions (CASE, NUMBER) are filled in relation to other features. Only the subtype of cardinal non-adjectival definite numerals covers both the types, i.e. distinguishes “noun” and “no-gender” type by a specific value of SUBPOS in the tag (cf. value C1 and Cz in the Tab. 51).

In certain contexts, numerals do not display grammatical relations with other words in a sentence by means of declension, they are used in a rigid (indeclinable) form. Cf. *ke čtyřem stům dětem* ‘to four hundred children’ (inflectional form) vs. *ke čtyři sta dětem* ‘to four hundred children’ (indeclinable form); cf. also *sto* ‘hundred’ in *do sto lidí* ‘to a hundred people’ or *v pouhých sto výtiscích* ‘in just a hundred copies’,<sup>20</sup> *tisíc* ‘thousand’ in *až po stovky tisíc let* ‘for hundreds of thousands of years’, *jedna* ‘one’ in *sedmdesát jedna občanů* ‘seventy-one citizens’, *čtvrtě* ‘quarter’ in *před tří čtvrtě rokem* ‘three quarters of a year ago’, *pár* ‘few’ in *o pár stech tisících* ‘about a few hundred thousand’, etc. So for most (cardinal) numerals we introduce a subspecified analysis (with the value X for GENDER, NUMBER, CASE). Cf. examples in Tab. 52.

<sup>20</sup>Note that in *v pouhých sto výtiscích* ‘in just a hundred copies’ the adjective *pouhých* ‘just’ is in agreement not with the following numeral *v pouhých sto výtiscích* ‘hundred’, but with the noun in the genitive *výtiscích* ‘copies’ depending on the numeral.

Wordform	Lemma	Tag
<u>jeden</u> muž	jeden‘1	CnYS1-----
<u>jedna</u> žena	jeden‘1	CnFS1-----
<u>jedno</u> dítě	jeden‘1	CnNS1-----
<u>několikáčký</u> problém	několikáčký	CwYS1-----
<u>několikáčké</u> problémy	několikáčký	CwIP1-----
<u>sto</u> lidí	sto-1‘100	CzNS1-----
<u>nulou</u> dělit nelze	nula	CzFS7-----
<u>patero</u> přikázání	patero‘5	CjNS1-----
<u>tři</u> domy	tři‘3	C1-P1-----
<u>do</u> <u>mnoha</u> zemí	mnoho-1	Ca--2-----
<u>několikrát</u> zazvonil	několikrát	Co-----
<u>několikráte</u> zazvonil	několikrát	Co-----1
ke <u>čtyřem</u> stům détem	čtyři‘4	C1-P3-----
ke <u>čtyři</u> sta détem	čtyři‘4	C1-XX-----
ke <u>čtyři</u> <u>sta</u> détem	sto-1‘100	CzNXX-----1
do <u>sto</u> lidí	sto-1‘100	CzNXX-----
až po stovky <u>tisíc</u> let	tisíc‘1000	CzIXX-----
<u>sedmdesát</u> <u>jedna</u> občanů	jeden‘1	CnXXX-----
před tři <u>čtvrtě</u> rokem	čtvrt	CzFXX-----1

Table 52: Examples: Subtypes of numerals

**Semantics.** We classify the following subtypes of numerals:

- **Cardinal** – express quantity.

Cn adjectival definite: *jeden* ‘one’, *dva* ‘two’, *oba* ‘both’.

Cy adjectival indefinite: *nejeden* ‘not-only-one’.

C1 non-adjectival no gender definite: *tři* ‘three’, *čtyři* ‘four’, *pět* ‘five’, *půl* ‘half’.

Ca non-adjectival no gender indefinite: *kolik* ‘how much’, *mnoho* ‘much/many’, *málo* ‘little/few’, *tolik* ‘that much/many’, *několik* ‘some (number of)’, *kdovíkolik* ‘who-knows-how-much/many’, *pár* ‘some (number of)’.

Cz non-adjectival ”noun” definite: *sto* ‘hundred’, *milion* ‘million’, *nula* ‘zero’, *čtvrt* ‘quarter’.

- **Ordinal** – express position in a sequential order.

Cr adjectival definite: *třetí* ‘third’, *pátý* ‘fifth’.

Cw adjectival indefinite: *kolikáčký* ‘at-what-position-in-a-sequence’, *tolikáčký* ‘at-that-position-in-a-sequence’, *několikáčký* ‘umpteenth’.

- **Multiplicative** – express how many times/folds.

Cv definite: *pětkrát* ‘five-times’, *sedmkrát* ‘seven-times’.

Co indefinite: *kolikrát* ‘how-many-times’, *mnohofokrát* ‘many-times’, *tolikrát* ‘that-many-times’, *několikrát* ‘several-times’, *nejednou* ‘not-only-one-time’.

- **Generic** – express number of different kinds, types.

Cd adjectival definite: *jedny* ‘one-kind’, *dvojí* ‘two-kinds’, *desaterý* ‘ten-kinds’, *patery* ‘five-kinds’, *obé*.

C<sub>h</sub> adjectival indefinite: *kolikery* ‘how-many-kinds’, *nejedny* ‘not-only-one-kind’, *tolikery* ‘that-many-kinds’, *několikery* ‘several-kinds’, *několikery* ‘several-kinds’.

C<sub>j</sub> non-adjectival ”noun” definite: *čtvero* ‘four-kinds-of’, *desatero* ‘ten-kinds-of’, *dvé*, *tré*.

C<sub>k</sub> non-adjectival ”noun” indefinite: *kolikero* ‘how-many-kinds’, *několikero* ‘several-kinds-of’, *tolikero* ‘that-many-kinds’.

Unlike Czech grammars, we do not distinguish interrogative numerals as a separate type: the interrogative numerals are included in the corresponding types of indefinite numerals (e.g. the interrogative numeral *kolikátý* ‘at-what-position-in-a-sequence’ is included in the ordinal indefinite type, the interrogative numeral *kolikrát* ‘how-many-times’ is included in the multiplicative indefinite type).

Words composed from a numeral and a preposition, e.g. *podruhé* ‘the-second-time’, *poněkolikáté* ‘several-times’, *zadruhé* ‘the-second-time’, *naněkolikrát* ‘several-times’ are considered compound adverbs.

### 11.3 Subtypes of adverbs

We divide adverbs into two groups differentiated in the second position in the tag (i.e. SUBPOS):

D<sub>b</sub> Adverb without a possibility to form negation and degrees of comparison, i.e. positions of the NEGATION and GRADE are not filled. E.g. *pozadu* ‘behind’, *naplocho* ‘flatly’, *včera* ‘yesterday’, *kde* ‘where’.

D<sub>g</sub> Adverbs forming negation or/and degrees of comparison; positions of the GRADE and NEGATION are filled. E.g. *rychle* ‘fast’, *nerychle* ‘not-fast’ *zajímavější* ‘more interesting’. Some adverbs only form negative worfforms, but they do not form degrees of comparison; e.g. *mnoho* ‘many’.

See examples in Tab. 53.

Wordform	Lemma	Tag
<i>pozadu</i>	<i>pozadu</i>	D <sub>b</sub> -----
<i>zítra</i>	<i>zítra</i>	D <sub>b</sub> -----
<i>kde</i>	<i>kde</i>	D <sub>b</sub> -----
<i>rychle</i>	<i>rychle_</i> <sup>^</sup> (*1ý)	D <sub>g</sub> -----1A----
<i>nerychle</i>	<i>rychle_</i> <sup>^</sup> (*1ý)	D <sub>g</sub> -----1N----
<i>nejrychleji</i>	<i>rychle_</i> <sup>^</sup> (*1ý)	D <sub>g</sub> -----3A----
<i>pracuje nemnoho</i>	<i>mnoho</i> -2	D <sub>g</sub> -----1N----

Table 53: Examples: Subtypes of adverbs

## 12 Negation

Negated words, which are only the opposite of positive words (e.g. *neherec* ‘non-actor’, *nevelký* ‘non-big’, *nepracovat* ‘do not work’) are captured in the same paradigm (with a positive lemma) together with non-negated forms. Negated forms have the value N at the 11<sup>th</sup> NEGATION position in the tag (Sect. 5.11), affirmative forms have the value A. Cf. examples in table 54.

Negated forms of adjectives, adverbs and verbs are generated into the dictionary regularly within inflective patterns. Negated forms of nouns are added to the dictionary selectively.

A negated lemma (and thus a separate paradigm for negated forms) occurs when the negated forms of the word are not merely a negation of the positive form (there is a noticeable shift in meaning; e.g. *nemoc* ‘disease’, *nebývalý* ‘unprecedented’). Cf. examples in Tab. 54.

For numerals (with the value C of POS), pronouns (with the value P of POS) and possessive adjectives (with the value AU at the POS and SUBPOS positions), the NEGATION position is not filled. Negated forms of these words are rare (e.g. *nemnoho* ‘non-much’, *neúčastníkův* ‘non-participant’s’, *nemůj* ‘non-mine’). If they occur, these negated forms are captured as a separate paradigm with the lemma in the negated form. Cf. examples in Tab. 54.

Wordform	Lemma	Tag
<i>neherec</i>	<i>herec</i>	NNMS1-----N---
<i>nevelký dům</i>	<i>velký</i>	AAIS1-----1N---
<i>zase nepracuje</i>	<i>pracovat</i>	VB-S---3P-NAI--
<i>má moc</i>	<i>moc-1</i>	NNFS4-----A---
<i>má nemoc</i>	<i>nemoc</i>	NNFS4-----A---
<i>nebývalý zájem</i>	<i>nebývalý</i>	AAIS1-----1A---
<i>neúčastníkův zájem</i>	<i>neúčastníkův_</i> <sup>^</sup> (*2)	AUIS1M-----
<i>nemnoho lidí</i>	<i>nemnoho</i>	Ca---1-----

Table 54: Examples: Negation

## 13 Names and Terms

Proper names in Czech always start with a capital letter. From morphological point of view, they are "normal" words, but it appears to be useful to give them a simple semantic label. Thus, name labels are attached to their lemmas. Named entities are very often multiple-word expressions and as such cannot be optimally described at wordform level. Therefore, lemma of a proper name bears name label only if it is a proper name itself, i.e. not in combination with another word. We respect the original idea that the name label explains the meaning of the lemma, not the context it appears in. Thus for instance *Nový* 'New' is lemmatized as *nový* in *Nový Bydžov*, not *Nový\_-;G* and *Bydžov* is lemmatized *Bydžov\_-;G*. For details see below.

For overview of the labels for names and terms see Sect. 4.2.2.

### 13.1 Personal names

Personal names are assigned the label Y. Given names and surnames are not distinguished by the name label in their lemmas. The value Y is used for all personal names, also for nicknames, names of horses, pets, etc. Names of members of a particular nation, inhabitants of a particular territory are labeled with E. See examples in Tab. 55.

Wordform	Lemma	Tag
<i>Pythagoras</i>	<i>Pythagoras_-;Y</i>	NNMS1----A----
<i>Jiří Včelař Kotas</i>	<i>Včelař_-;Y</i>	NNMS1----A----
<i>Alík</i>	<i>Alík_-;Y</i>	NNMS1----A----
<i>Dán</i>	<i>Dán_-;E</i>	NNMS1----A----
<i>fotbalista Petr Čech je Čech</i>	<i>Čech_-;E_-;Y</i>	NNMS1----A----
<i>Pcheng je čínské jméno i příjmení</i>	<i>Pcheng_-;Y</i>	NNMS1----A----
<i>Italčino vítězství</i>	<i>Italčin_-;E_-^(*3ka)</i>	AUNS1F-----
<i>František Palacký</i>	<i>Palacký_-;Y</i>	NNMS1----A----
<i>Milena Jesenská</i>	<i>Jesenská_-;Y</i>	NNFS1----A----
<i>paní Nováková</i>	<i>Nováková_-;Y</i>	NNFS1----A----
<i>Angelina Jolieová</i>	<i>Jolieová_-;Y</i>	NNFS1----A----
<i>Angelina Jolie</i>	<i>Jolie_-;Y</i>	NNFXX----A----
<i>s Georgem Washingtonem</i>	<i>George_-;Y</i>	NNMS7----A----
<i>s George Washingtonem</i>	<i>George_-;Y</i>	NNMXX----A----

Table 55: Examples: Personal names

If a name can serve both as a personal name and as a member of a particular nation (e.g. name *Čech* as a surname and nation) and if there is a same declension for the both usages, then there is one lemma/paradigm covering all usages and all relevant name labels are attached to it (cf. ex. in Tab. 55). Possessive adjectives derived from proper names (*Novákův*, *Italčin*) inherit the name labels from the original lemma.

Personal names homonymous with a normal Czech word always have a lemma of their own. Thus *Zeman* (surname) is lemmatized as *Zeman-1\_-;Y*, not *zeman* 'squire' (cf. also ex. of *Písek* in Tab. 8 in Sect. 4.1.1).

Personal names are always tagged as nouns, even if they have an adjectival form (cf. ex. *Palacký*, *Jesenská*, *Nováková*, *Jolieová* in Tab. 55).

Foreign personal names are not marked as foreign words because in Czech texts, they are usually declined according to the Czech grammar; e.g. *Bill Clinton*, *bez Billa Clintonona*, *Billu Clintonovi*, *s Billem Clintonem....* Even if a name allows for a frozen (undeclined) form, there usually is a context in which it can be declined: *s George Washingtonem* vs. *s Georgem Washingtonem*. Some foreign names, such as *Steffi*, *Jolie* are never declined. Cf. Tab. 55.

Names of horses, pets, etc. have all kinds of names (e.g. *Vinná réva*, *He Shall Reign*, *La Paloma*, *Lučina*, *Areál*). Quite often one does not know whether it is male or female (sometimes even female-like names belong to a male horse). If any reasonable analysis is possible it should be used regardless the lemma is marked as a name or not. It will be marked as a name within a separate project on named entity recognition. However, if the name is a word that has no other meaning or if it has different gender (declension), a new lemma with the label Y is introduced.

**Prepositions and other function words in personal names.** Prepositions, conjunctions and other determiners form parts of personal names that indicate geographical roots of the family (*Jiří z Poděbrad*, *Kryštof Harant z Polžic a Bezdružic*). They are analyzed as normal words. It may not be always clear whether the part after the preposition shall be annotated as a surname or a geographical name. If the Czech preposition *z* is present, the following word is a geographical name (even if it is a foreign location as in *Blanka z Valois*). In the foreign personal names (e.g. *Ludwig van Beethoven*, *Miguel de Cervantes y Saavedra*, *Hans van den Broek*), the foreign prepositions (*von*, *van* and *de*) are analyzed as foreign words. The original geographical meaning is usually less obvious for a Czech reader and the following word is annotated as surname. See ex. in Table 56.

Wordform	Lemma	Tag
<i>Jiří z Poděbrad</i>	<i>z-1</i>	RR--2-----
<i>Jiří z Poděbrad</i>	<i>Poděbrady_</i> ;G	NNIP2----A----
<i>František Lobkovic</i>	<i>Lobkovic_</i> ;Y	NNMS1----A----
<i>František z Lobkovic</i>	<i>Lobkovice_</i> ;G	NNFP2----A----
<i>František z Lobkovic</i>	<i>z-1</i>	RR--2-----
<i>Blanka z Valois</i>	<i>Valois_</i> ;G	NNNX---A----
<i>Ludwig van Beethoven</i>	<i>van-77</i>	F%-----
<i>Ludwig van Beethoven</i>	<i>Beethoven_</i> ;Y	NNMS1----A----

Table 56: Examples: Function words in personal names

**Chinese and other Asian names** can consist of one syllable or of two (or more) syllables, often connected with a dash (however sometimes separated by a space). For multi-syllable (hyphenated) names (e.g. *Siao-Pching*, *Ir-Sen*), the tag usually specifies NUMBER and CASE position only for the last syllable. If the name is in the nominative according to the context (and therefore it is not possible to recognize whether it is used as a declinable word), the NUMBER and CASE positions are filled as a nominative (i.e. value X is not use). However, if a foreign name is never declined in Czech, the unspecified tag is used (with X-value at NUMBER and CASE positions). The last syllable in a hyphenated composite can also be captured as a segment, especially if it is written with lower case and expresses a case and a number (e.g. *Chuang-timu*). Cf. ex. in Tab. 57.

Wordform	Lemma	Tag
<i>Teng Siao-Pching odchází.</i>	<i>Teng_</i> ;Y	NNMS1----A----
<i>Teng Siao-Pching odchází.</i>	<i>Siao-3_</i> ;Y	NNMXX----A----
<i>Teng Siao-Pching odchází.</i>	<i>Pching_</i> ;Y	NNMS1----A----
<i>podobá se Čchin Š' Chuang-timu</i>	<i>Čchin-1_</i> ;Y	NNMXX----A----
<i>podobá se Čchin Š' Chuang-timu</i>	<i>Š-1_</i> ;Y	NNMXX----A----
<i>podobá se Čchin Š' Chuang-timu</i>	<i>Chunag-1_</i> ;Y	NNMXX----A----
<i>podobá se Čchin Š' Chuang-timu</i>	<i>ti-2_</i> ;Y	SNMS3----A----

Table 57: Examples: Chinese and other Asian names

## 13.2 Geographical names

Geographical names are assigned label G.

### 13.2.1 Countries, cities, rivers, mountains

The main word (head) in a multi-word geographical name is a noun; the same holds for a one-word city name. If it is homonymous with an adjective, a new noun lemma/paradigm is created for the name. Thus *Hluboká* in the name *Hluboká nad Vltavou* is captured as a noun with lemma *Hluboká\_*;G (not as an general adjective *hluboký* ‘deep’).

Nouns that are frequently used in names (such as *Újezd*, *Ústí* may have their own geographical lemmas even if they are homonymous with a normal word. For homonymous pairs where the non-geographical usage is much more common (such as *voda* ‘water’, *ves* ‘village’, *město* ‘city’) it is recommended to capture them with the non-geographical lemma even in geographical usages. Other words in multi-word names (adjectives, prepositions, conjunctions, etc.) are represented as ordinary words. Other nouns in names can be represented as geographical names only if they are themselves geographical names (names of rivers and mountains in city names, etc.). See Tab. 58.

Wordform	Lemma	Tag
<i>Hluboká nad Vltavou</i>	<i>Hluboká_</i> ;G	NNFS1-----A----
<i>Ústí nad Labem</i>	<i>Ústí_</i> ;G	NNNS1-----A----
<i>Ústí nad Labem</i>	<i>nad-1</i>	RR--7-----
<i>Ústí nad Labem</i>	<i>Labe_</i> ;G	NNNS7-----A----
<i>Ohrada u Hluboké</i>	<i>ohrada</i>	NNFS1-----A----
<i>Ohrada u Hluboké</i>	<i>Hluboká_</i> ;G	NNFS2-----A----
<i>Kostelec nad Černými lesy</i>	<i>Kostelec_</i> ;G	NNIS1-----A----
<i>Kostelec nad Černými lesy</i>	<i>černý_</i> ;o	AAIP7----1A----
<i>Kostelec nad Černými lesy</i>	<i>les</i>	NNIP7----1A----
<i>Karlovy Vary</i>	<i>Karlův_</i> ;Y_ ^(*3el)	AUIP1M-----
<i>Karlovy Vary</i>	<i>Vary_</i> ;G	NNIP1-----A----
<i>Orlické hory</i>	<i>orlický</i>	AAFP1----1A----
<i>Orlické hory</i>	<i>hora</i>	NNFP1-----A----
<i>Divoká Orlice</i>	<i>divoký</i>	AAFS1----1A----
<i>Divoká Orlice</i>	<i>Orlice_</i> ;G	NNFS1-----A----
<i>na Mont Blanka</i>	<i>Blanc-1_</i> ;G	NNIS6-----A----
<i>v Cincinnati</i>	<i>Cincinnati_</i> ;G	NNNXX-----A----
<i>v Los Angeles</i>	<i>Los-77</i>	F%-----
<i>v Los Angeles</i>	<i>Angeles-77</i>	F%-----

Table 58: Examples: Names of countries, cities, rivers, mountains

**Foreign geographical names** are mostly annotated as an undeclined noun (e.g. *Cincinnati*) or foreign words (mainly in the case of multi-word names, e.g. *Los Angeles*). If they are declined in Czech, they have analysis according to the Czech morphology (no matter how morphologically they behave in the original language). For instance, *blanc* is adjective in French *Mont Blanc* but it behaves as a noun in *na Mont Blanka*. *Mont* is annotated as a foreign word. See Tab. 58. See Sect. 17 for more information on foreign words.

**Adjectives derived from geographical names.** Adjectives derived from geographical names (e.g. *africký* ‘African’ from *Afrika* ‘Africa’) are not marked as geographical (no G label in lemma). These adjectives are not capitalized in Czech, while the original nouns are.

### 13.2.2 Streets, squares, stations

In street, station and square names, we suppose that a word such as *ulice* ‘street’, *náměstí* ‘square’, etc. is always present, even if elided on the surface. Therefore the tagging of the name of the street and square is not altered. Cf. Table 59.

Wordform	Lemma	Tag
<i>Bydlím v Horské ulici</i>	horský	AAFS6----1A----
<i>Bydlím v ulici Štěpánská</i>	štěpánský	AAFS1----1A----
<i>Bydlím v ulici Mezi Zahrádkami</i>	mezi-1	RR---7-----
<i>Bydlím v ulici Mezi Zahrádkami</i>	zahrádka	NNFP7----A----
<i>Palackého náměstí</i>	Palacký_-;Y	NNMS2----A----
<i>Sejdme se na Palackého</i>	Palacký_-;Y	NNMS2----A----
<i>Bydlím na náměstí Míru</i>	mír_-^(opak_války)	NNIS2----A----
<i>stanice Staroměstská</i>	staroměstský	AAFS1----1A----
<i>ve směru od Dejvické</i>	dejvický	AAFS2----1A----
<i>zastávka Na Knížecí</i>	Knížecí_-;G	NNNS6----A----
<i>stanice Anděl</i>	Anděl-2_-;G	NNIS1----A----

Table 59: Examples: Names of streets, squares, stations

### 13.2.3 Buildings

Building names can be annotated as a geographic name. On the other hand, many building names are composed of common words (apelatives) and we do not capture such names as geographical names. Cf. Tab. 60.

Wordform	Lemma	Tag
<i>Rudolfinum</i>	Rudolfinum_-;G	NNNS1----A----
<i>Pražský hrad</i>	pražský	AAIS1----1A----
<i>Pražský hrad</i>	hrad	NNIS1----A----
<i>Chrám sv. Barbory</i>	chrám	NNIS1----A----
<i>Chrám sv. Barbory</i>	svatý-1	AAXXX----1A---b
<i>Chrám sv. Barbory</i>	Barbora_-;Y	NNIS2----A----

Table 60: Examples: Names of buildings

## 13.3 Scientific terminology

Scientific terms from chemistry, medicine, natural science (written with uppercase letter) are assigned U-label. Cf. Tab. 61.

Wordform	Lemma	Tag
<i>Acylpirin</i>	Acylpirin_-;U	NNIS1----A----
<i>Australopithecus</i>	Australopithecus_-;U	NNMS1----A----
<i>Hydrosulfit</i>	Hydrosulfit_-;U	NNIS1----A----

Table 61: Examples: Scientific terminology

### 13.4 Other proper names

Other proper names, names of companies, foundations, shops, clubs, sport clubs, restaurants, unique product names, names of events, works of art etc. have lemmas flagged `m`.

However, “words”, the usage of which is not limited to the name, get their “normal” lemmas (and they are not annotated as a name). Only if a word cannot be explained another way or if its meaning has nothing to do with the company or in the case of well-known name (e.g. *Škoda*), the label and capitalized lemma is used.

The border between personal and other names is fuzzy: if it is clear that a surname is part of a company name (e.g. *Uzenářství Novák a syn*) it is lemmatized with personal name label. On the other hand, *Škoda* is annotated as a company no matter that it was also named after a person. This name is mostly known as a company name. The same applies to geographic and other names if they are part of a company name (cf. *Stavební společnost Šumava, s.r.o* and other examples in Tab. 62).

Wordform	Lemma	Tag
<i>Škoda-auto, a. s.</i>	Škoda-1_; <code>m</code>	NNFS1-----A----
<i>Škoda-auto, a. s.</i>	auto	NNNS1-----A----
<i>Uzenářství Novák a syn</i>	Novák_; <code>Y</code>	NNMS1-----A----
<i>restaurace U Medvídku</i>	u-1	RR--2-----
<i>restaurace U Medvídku</i>	medvídek	NNMP2-----A----
<i>TJ Sokol</i>	Sokol-2_; <code>m</code> _^(organizace)	NNIS1-----A----
<i>Stavební společnost Šumava</i>	Šumava_; <code>G</code>	NNFS1-----A----
<i>televize Nova</i>	Nova_; <code>m</code> _^(televize)	NNFS1-----A----
<i>Tatra, a. s.</i>	Tatra_; <code>m</code> _^(auto_;; mléko)	NNFS1-----A----
<i>nový Renault</i>	Renault_; <code>m</code>	NNIS1-----A----
<i>časopis Sluníčko</i>	sluníčko	NNNS1-----A----
<i>výstava Habitat 2019</i>	Habitat_; <code>m</code>	NNIS1-----A----
<i>Mistrovství světa v hokeji</i>	mistrovství	NNNS1-----A----

Table 62: Examples: Other proper names

## 14 Abbreviations

There are two ways to capture abbreviations:

- as a special (abbreviated) wordform of a non-abbreviated (single) word paradigm: there is a special value (b, a or c) at the 15<sup>th</sup> position of the tag,
- as a special POS Abbreviation with the code B at the 1<sup>st</sup> position of the tag.

In the first way, fixed abbreviations of a single word are captured (e.g. abbreviation *s.* ‘p.’ of word *strana* ‘page’; see more in Sect. 14.1). Other abbreviations (e.g. *USA*) are captured in the second way (see more in Sect. 14.2).

### 14.1 Fixed abbreviations of a single word

Fixed, well-known abbreviations of a single word usually followed by a period in Czech text (e.g. abbreviation *s.* ‘p.’ of word *strana* ‘page’) are captured as a special wordform of the paradigm of the non-abbreviated word. The fact that it is an abbreviated form is expressed at the 15<sup>th</sup> position by the letters b (a or c in the case of the other fixed abbreviation of the same word (e.g. *s.* and *str.* ‘p.’ are both used as the abbreviation for *strana* ‘page’)). Only those categories that are valid for each use of the abbreviation are coded in the tag. Usually, the NUMBER and CASE positions are not specified.

Wordform	Lemma	Tag
<i>délka 3 m</i>	<i>metr</i>	NNIXX-----A---b
<i>3m tyč</i>	<i>metrový</i>	AAXXX-----1A---b
<i>ve 2 h</i>	<i>hodina</i>	NNFXX-----A---b
<i>ve 2 hod.</i>	<i>hodina</i>	NNFXX-----A---a
<i>na s. 8</i>	<i>strana</i>	NNFXX-----A---b
<i>na str. 8</i>	<i>strana</i>	NNFXX-----A---a
<i>za 2 s</i>	<i>sekunda</i>	NNFXX-----A---b
<i>za 2 sec</i>	<i>sekunda</i>	NNFXX-----A---a
<i>za 2 sek.</i>	<i>sekunda</i>	NNFXX-----A---c
<i>kWh</i>	<i>kilowatthodina</i>	NNFXX-----A---b
<i>800 m n. m.</i>	<i>nad-1</i>	RR--7-----b
<i>800 m n. m.</i>	<i>moře</i>	NNNS7-----A---b
<i>č. 5</i>	<i>číslo</i>	NNNXX-----A---b
<i>300 n. l.</i>	<i>náš</i>	PSZS2-P1-----b
<i>300 n. l.</i>	<i>letopočet</i>	NNIS2-----A---b
<i>např.</i>	<i>například</i>	TT-----b
<i>It.</i>	<i>Itálie_</i> ;G	NNFXX-----A---b

Table 63: Examples: Fixed abbreviations of a single word

**Units of measurements.** In this way, the abbreviations of units of measurements are captured, too (e.g. *km* is abbreviated form of the word *kilometr*, *C* is abbreviated form of *Celsius*, etc.). This type of abbreviations usually appears without the period. Cf. examples in Tab. 63.

**Note.** It is impossible to capture the abbreviation of more than one word (e.g. abbreviation *atd.* ‘etc.’ of the words *a tak dále* ‘et cetera/and so on’) in the way described here. These abbreviations belong to the part of speech Abbreviations (with the B value at POS position). They are described in the following section 14.2. cf. Tab. 64.

Wordform	Lemma	Tag
<i>atd</i>	<i>atd_ ^ (a_tak_dále)</i>	Bb-----
<i>apod</i>	<i>apod_ ^ (a_podobně)</i>	Bb-----
<i>ap</i>	<i>ap_ ^ (a_podobně)</i>	Bb-----
<i>aj</i>	<i>aj-1_ ^ (a_jiný/á/é)</i>	BAXXX----1A----

Table 64: Examples: Fixed abbreviations of multiple words

## 14.2 Other abbreviations

Abbreviations which abbreviate at least two words (e.g. *atd* 'and so on'), or are composed of uppercase letters e.g. *USA*) are captured as a special POS with the value B at the 1<sup>st</sup> tag position. The lemma of such an abbreviation is the abbreviation itself. The belonging of an abbreviation to a traditional POS – although this may be difficult to determine in a number of cases (e.g. multi-word abbreviations; cf. abbreviation *atd.* in Tab. 64) – is reflected at the 2<sup>nd</sup> position of the tag. The value at the 2<sup>nd</sup> position also determines which other positions of the tag are filled.

Most abbreviations are nouns and can be used with more than one gender. Of course, abbreviations have no endings but the surrounding context can reveal their underlying gender whenever gender agreement is required by the Czech grammar (e.g. *staronový NKÚ* as masculine inanimate vs. *slavné NKÚ* as neuter). However, the GENDER, NUMBER and CASE positions are usually not specified in the tag of abbreviations (cf. examples in Tab. 65 and 66).

### 14.2.1 Well-known abbreviations composed of uppercase letters

The lemma of well-known abbreviations composed of uppercase letters (e.g. *USA*; everyone knows what this abbreviation means even without context) is accompanied with a name label (e.g. G in case of geographical name abbreviation). There is a semantic explanation attached to the lemma. Cf. examples in Tab. 65.

Similarly, the symbols of chemical elements and compounds (e.g. abbreviation *N* for nitrogen, *CO* for carbon monoxide) as well as academic titles (e.g. *JUDr*) are captured in this way.

Wordform	Lemma	Tag
<i>USA</i>	<i>USA_ ; G_ ^ (United_States_of_America)</i>	BNXXX-----A----
<i>NKÚ</i>	<i>NKÚ_ ; m_ ^ (Nár._kontrolní_úřad)</i>	BNXXX-----A----
<i>CD</i>	<i>CD-1_ ^ (Audio/Data,_Compact_Disc)</i>	BNXXX-----A----
<i>KB</i>	<i>KB_ ; m_ ^ (Komercní_banka)</i>	BNXXX-----A----
<i>N</i>	<i>N-1_ ; U_ ^ (zn._dusíku)</i>	BNXXX-----A----
<i>JUDr</i>	<i>JUDr_ ^ (doktor_práv)</i>	BNXXX-----A----

Table 65: Examples: Well-known abbreviations composed of uppercase letters

### 14.2.2 Less familiar abbreviations and abbreviations with many meanings

In the case of less familiar abbreviations and abbreviations with many meanings (cf. HK in Tab. 66) we do not distinguish partial meanings; there is only one analysis for all usages of the abbreviated form. These abbreviations have the lemma with the number -88. The tag is always BNXXX-----A----. Cf. Tab. 66.

**Note:** One-letter abbreviations with many meanings (e.g. *V. Havel* vs. *V. Mrštík*) are captured as a special POS for isolated letters (the tag starts with Q3), see Sect. 15.

Wordform	Lemma	Tag
<u>HK</u> je horní končetina	HK-88	BNXXX-----A----
<u>HK</u> je Hradec Králové	HK-88	BNXXX-----A----
<u>HK</u> je Hospodářská komora	HK-88	BNXXX-----A----
<u>OU</u> je Oxford University	OU-88	BNXXX-----A----
<u>OU</u> je odborné učiliště	OU-88	BNXXX-----A----
<u>pr.</u> volno	pr-88	BNXXX-----A----

Table 66: Examples: Less familiar abbreviations and abbreviations with many meanings

#### 14.2.3 Author's signature

The author's name abbreviations used in newspapers (e.g. *ber*, *mas*, etc. in "sentences" like *PRAHA (ČTK, ber)*) have the lemma equal to the wordform. They are numbered -99 as a human-readable indication and flagged Y. Cf. Tab. 67.

Wordform	Lemma	Tag
<i>Praha (haš)</i>	haš-99_-;Y	BNXXX-----A----
<i>Praha (ČTK, ber)</i>	ber-99_-;Y	BNXXX-----A----

Table 67: Examples: Author's signature

## 15 Isolated Letters

Isolated letters stand for many meanings. We do not distinguish between the usage as an abbreviation (e.g. *A. Franklin*), a label (e.g. *skupina A* ‘group A’), *A-konto* ‘A-account’), an item in a list (e.g. *a*), a separator in a text (e.g. *o o o o o o o o*), or even other meanings.

Any isolated letter is represented as a special POS with the value Q3 at the first and second position of the tag. The lemma has the number -33 as a human-readable indication. This analysis exists for all letters of the Czech alphabet (upper- and lower-case). Cf. examples in Tab. 68.

**Note.** Not every one-letter token in a text is annotated as an isolated letter. We distinguish letters as well-known abbreviations in the case of chemical symbols (e.g. *O* for *oxygen*; see Sect. 14.2.1), letters as one-syllable personal names (especially Chinese; e.g. *S* in *Wang S'tching*), and letters as foreign words (e.g. *I* in *I love you*; see Sect. 17) .

Wordform	Lemma	Tag
<i>kružnice A</i>	A-33	Q3-----
<i>jedeme do p.</i>	p-33	Q3-----
<i>V. Havel</i>	V-33	Q3-----
<i>V. Mrštík</i>	V-33	Q3-----
<i>Karlovy V.</i>	V-33	Q3-----
<i>Ch. Dickens</i>	Ch-33	Q3-----

Table 68: Examples: Isolated Letters

## 16 Segments

Segments are incomplete words. They are parts of words; in order to understand them, they must be joined with another string or word to create a complete word. They are usually joined with a separator, most often with a hyphen. Note that combinations of two complete words connected with a separator e.g. *Praha-Hradčany* ‘Prague-Hradcany’, *Anna-Marie*, *propan-butan* ‘propane-butane’) are not segments (see more in Sect. 19). The segments are treated as a special POS with the value **S**. According to their position in the complete word, we distinguish:

- prefixal segments,
- postfixal segments.

Wordform	Lemma	Tag
<u>česko-ruská kniha</u>	česko	S2-----A----
<u>nepoliticko-politická diskuse</u>	politicko	S2-----N----
<u>hudebně-zábavný pořad</u>	hudebně-2	S2-----A----
<u>pěti- až sedmiletý chlapec</u>	pěti'5	S2-----A----
<u>ultra-moderní</u>	ultra-1	S2-----A----
<u>ultra jemný prášek</u>	ultra-2	AAXXX----1A----
<u>mini-sukně</u>	mini-1	S2-----A----
<u>mini sukně</u>	mini-3	AAXXX----1A----
<u>ne/souhlasím</u>	ne-2	S2-----A----
<u>s manželem/kou</u>	ka	SNFS7-----A----
<u>n-tice</u>	tice_^(n-tice)	SNFS1-----A----
<u>řekl(a)</u>	a-2	SpQW----R-AA---
<u>uslyší to už po x-té</u>	tý-2_^(x-tý)	SAFS6----1A----
<u>přikládáme soubor/y</u>	y-2_^(soubor/y)	SNIP4-----A----
<u>12-tí hodinová směna</u>	ti_^(10-ti)	S1-XX-----
<u>na Tchaj-wanu</u>	wan	SNIS6-----A----

Table 69: Examples: Segments

**Prefixal segments** are strings that appear at the beginning of words. They are usually followed with a separator, most often with a hyphen (e.g. segment *česko* in example *česko-ruská kniha* ‘Czech-Russian book’), but also with a space or another separator (e.g. segment *pěti* in example *pěti až sedmiletý chlapec*).

A typical prefixal segment has a special form ending with *-o* (e.g. *česko-ruský* ‘Czech-Russian’) but there are other types of segments (cf. segments *kvazi* ‘quasi’, *hydroxy* in *kvazi-valenční* ‘quasi-valency’, *hydroxy-sloučeniny* ‘hydroxy-compounds’). The form of prefixal segment can be homonymous with an adverb or wordforms belonging to other POS (cf. *hudebně* ‘musically’ as a segment in *hudebně-zábavný* ‘musically-entertaining’ and *hudebně* as an adverb in *hudebně nadaný* ‘musically gifted’ or *mini* as a segment in *mini-sukně* ‘mini-skirt’ and as a noun in *nosit mini* ‘wear a mini’). Homonymous wordforms are common particularly in the case of inflexible loanwords (cf. *super* as a segment in *super-moderní* ‘super-modern’, as an adjective in *super zábava* ‘super fun’ and as an adverb in *mám se super* ‘I am super’; see other examples in Tab. 69).

If a possible segment is not expressed in a special ”segment” form like *česko*, *anglo*, *tří*, then we tag it as a segment only in those cases where the wordform is attached to the next word by a non-space separator (cf. *pop*, *mini* as segments in *pop-kultura* ‘pop-culture’, *mini-sukně* ‘mini-skirt’ and as adjectives in *pop kutura* ‘pop culture’, *mini sukňe* ‘mini skirt’).

Lemma of a prefixal segment is the string itself, unless it is in negative form. In that case, the positive form (without the negative prefix *ne-* ‘non’) is considered to be the lemma (cf. *nepoliticko* ‘nonpolitical’ in Tab. 69). The tag of all prefixal segments has the code 2 at the 2<sup>nd</sup> position. Moreover, we specify for them also the 11<sup>th</sup> position concerning negation (see examples in Tab. 69).

**Postfixal segments** are strings that appear at the end of a wordform. They are usually attached directly to the word they combine with (cf. segment *ka* in example *manžel/ka* ‘husband/wife’) or segment *kou* in example *s manželem/kou* ‘with husband/wife’). The separator is most often a hyphen, a parenthesis or a slash.

The postfixal segments express an affiliation to a specific POS. Thus, all the inflectional categories that describe the whole wordform, except for the first one (the code for POS, which is S), are filled in the tag (with the exception of the ASPECT for verbs). The lemma is the closest “basic wordform”. See examples in Tab. 69.

## 17 Foreign Words

Foreign words enter Czech texts in three different ways:

- in citations (i.e. whole phrase in a foreign language, multi-word foreign name),
- as individual words (i.e. single foreign word),
- as domesticated words of foreign origin.

### 17.1 Citation use

Whole phrases in a foreign language are sometimes inserted into Czech texts as citations. Multi-word foreign names, multi-word Latin nomenclature, names of songs and other works belong to this category, too (e.g. *European market research center*, *University of Colorado*, *Monomorium floricola*).

Wordform	Lemma	Tag
<i>práce v European market research center</i>	European-77	F%-----
<i>práce v European market research center</i>	market-77	F%-----
<i>práce v European market research center</i>	research-77	F%-----
<i>práce v European market research center</i>	center-77	F%-----
<i>na University of Colorado</i>	University-77	F%-----
<i>na University of Colorado</i>	of-77	F%-----
<i>na University of Colorado</i>	Colorado-77	F%-----
<i>v Colorado Springs</i>	Colorado-77	F%-----
<i>v Coloradu Springs</i>	Colorado-1_1;G	NNNS6----A---
<i>v Coloradu Springs</i>	Springs-77	F%-----
<i>chování mravence Monomorium minimum</i>	Monomorium-77	F%-----
<i>chování mravence Monomorium minimum</i>	minimum-77	F%-----
<i>Peugeot 306 Grand Prix</i>	Peugeot-77	F%-----
<i>Peugeot 306 Grand Prix</i>	306	C=-----
<i>Peugeot 306 Grand Prix</i>	Grand-77	F%-----
<i>Peugeot 306 Grand Prix</i>	Prix-77	F%-----
<i>The FIA Cup</i>	The-77	F%-----
<i>The FIA Cup</i>	FIA-88	BNXXX----A---
<i>The FIA Cup</i>	Cup-77	F%-----
<i>Třetí ročník Škoda Hockey Cupu</i>	Škoda-1_1;m	NNFS1----A---
<i>Třetí ročník Škoda Hockey Cupu</i>	Hockey	F%-----
<i>Třetí ročník Škoda Hockey Cupu</i>	cup-1_1^(pohár)	NNIS2----A---
<i>najlepšie slovenské piesne</i>	najlepšie-77	F%-----
<i>najlepšie slovenské piesne</i>	slovenský	AAFP1----1A---
<i>najlepšie slovenské piesne</i>	piesne-77	F%-----

Table 70: Examples: Foreign words in citations

All wordfoms in such a foreign piece of text are analyzed as the foreign word POS coded F% at the 1<sup>st</sup> and 2<sup>nd</sup> tag positions. Lemma is the same as the wordform, including uppercase and lowercase initial letter, it has an index of 77 and no other labels. The whole tag is always F%-----. There are two exceptions to this rule: Abbreviations of uppercase letters (cf. *FIA* in *The FIA Cup*) are analyzed as abbreviations (see Section 14) and numbers writtten in digits have always tag C=-----. See examples in Table 70.

If an apparently foreign word is used with a Czech suffix, we analyse it as a Czech word. In other words, its POS is not "foreign word". For instance the wordform *Cupu* in *Třetí ročník Škoda Hockey Cupu* has an analysis based on the Czech morphology. This rule applies mostly to multi-word foreign geographical names (e.g. *v Tel Avivu*, *v Coloradu Springs*, etc.). See examples in Table 70.

**Slavic languages and Czech dialects.** Slavic languages (most prominently Slovak, but also Czech dialects and old Czech) are related to contemporary standard Czech. Citations may contain words that are identical to their Czech counterparts. When a word has a foreign suffix it must be annotated as a foreign word even if its baseform belongs to Czech. If all words in a phrase are identical in their forms and meanings to Czech, the phrase should be annotated as Czech, even if we know that it is in fact Slovak or other language. For instance, if a Slovak song was named *Drahý otec*, there is no need to annotate it as foreign. However, if a single word does not fit the Czech grammar or vocabulary, the best would be to annotate whole citation as foreign words. See example *Vracaja sa dom* in Table 70.

## 17.2 Single word use

In word usages, Czech morphology takes precedence.

Wordform	Lemma	Tag
<i>s Georgem Washingtonem</i>	George_.;Y	NNMS7----A----
<i>s George Washingtonem</i>	George_.;Y	NNMXX----A----
<i>to je George Washington</i>	George_.;Y	NNMS1----A----
<i>pracuje v Coloradu</i>	Colorado-1_.;G	NNNS6----A----
<i>pracuje v Colorado</i>	Colorado-1_.;G	NNNXX----A----
<i>prapaguje nový business</i>	business	NNIS4----A----
<i>v novém businessu</i>	business	NNIS6----A----
<i>mluvil s Hillary</i>	Hillary-1_.;Y	NNFXX----A----
<i>následuje písnička Girls</i>	Girls-77	F%-----

Table 71: Examples: Single foreign words

The basic rule applies to separately used foreign words in a Czech sentence (often names, terms; e.g. *George*, *Colorado*): If a wordform takes a Czech suffix (e.g *s Georgem*, *v Coloradu*), it is not captured as a foreign word. It is analyzed according to the Czech morphology. See examples in Table 71. If a single foreign name or term does not take Czech suffixes (e.g *s George Washingtonem*, *pracuje v Colorado*), it is usually captured as an inflexible noun (with value X in NUMBER and CASE position). However, in the case of little known words (e.g *následuje písnička Girls*), the wordform is captured as the foreign word POS.

Wordform	Lemma	Tag
<i>online služby</i>	online-1	AAXXX----1A----
<i>pracuje online</i>	online-2	Db-----
<i>před úvodním buly</i>	buly	NNNXX----A----
<i>s pop artem</i>	art-1	NNIS7----A----
<i>art zóna</i>	art-2	AAXXX----1A----
<i>Museum of Art</i>	Art-77	F%-----
<i>láska, jídlo a faux pas</i>	faux-77	F%-----
<i>láska, jídlo a faux pas</i>	pas-77	F%-----

Table 72: Examples: Domesticated words of foreign origin

### 17.3 Domesticated words of foreign origin

Foreign words constantly enter Czech language, take Czech endings, settle with Czech declension paradigms and become normal Czech words. Words that entered Czech long ago are not felt as foreign any more (e.g. *kakao*). Nevertheless, even newer words (e.g. *hardware*) should not be treated as foreign if they fit into this category.

Domesticated loanwords are analyzed according to the Czech morphology. However, some loanwords are inflexible, they do not settle with any Czech declension paradigms, but they are part of the Czech vocabulary (e.g. *online*, *buly*, *wi-fi*). These domesticated loanwords are not captured with the foreign word POS, they are treated as inflexible nouns, adjectives, adverbs, etc. (cf. also Sect. 10.1). Multi-word inflexible loanword (e.g. *faux pas*, *de iure*, *de facto*, *play off*, *ad hoc*, *cash flow*), though domesticated, are captured as foreign word POS. The boundary between inflexible loanwords and foreign words is very blurred. See examples in Table 72.

## 18 Aggregates

An aggregate is a wordform that is created by joining two or more wordforms (components of the aggregate) into one and cannot be simply assigned any POS. Aggregates are common especially in agglutinative languages, but there are some aggregate types in Czech, too. The following types of aggregates are captured:

- **pronominal aggregate:** aggregate consisting of a preposition and the pronoun *on* ‘he’ or *co*, *copak* ‘what’ (e.g. *pro* + *on* → *pron*; *za* + *co* → *zac*).
- **verbal aggregate:**
  - aggregate containing the contracted *-s* which stands for the wordform *jsi* ‘you are’. It can be appended to the end of a wordform of almost any POS (e.g. *promluvil* + *jsi* → *promluvils*; *dobře* + *jsi* → *dobřes*).
  - aggregate consisting of the conditional verbal wordform *by* ‘would’ or conditional conjunction such as *aby*, *kdyby* ‘so that, if-would’ and contracted form of the auxiliary verb *být* ‘to be’: *-ch* for the 1<sup>st</sup> person singular, *-s* for the 2<sup>nd</sup> person singular, *-chom* for the 1<sup>st</sup> person plural, *-ste* for the 2<sup>nd</sup> person plural (e.g. *kdyby* + *byste* → *kdybste*; *aby* + *bychom* → *abychom*).

The lemma of a pronominal aggregate is the lemma of the pronoun. The lemma of a verbal aggregate (containing the contracted form of the present tense of the auxiliary verb *být* ‘to be’) is the lemma of its first component.

The fact that a wordform is an aggregate is coded at the 14<sup>th</sup> position of the tag. The code of pronominal aggregates corresponds to the initial letter of the preposition that forms their first component. Verbal aggregates are coded with the letter **c** for *-ch*, **s** for *-s*, **m** for *-chom* and **e** for *-ste*; see also Section 5.14. Verbal and pronominal aggregates can combine; such aggregates are marked with the initial letter of the preposition, but in an uppercase letter (see the example *začs* in Table 73). The lemma of such (combined) aggregates is the pronoun involved.

Wordform	Lemma	Tag
<i>zač</i>	<i>co</i>	PQ--4-----z-
<i>začs</i>	<i>co</i>	PQ--4-----Z-
<i>doň</i>	<i>on-1</i>	P5ZS2--3-----d-
<i>načpaks</i>	<i>copak</i>	PQ--4-----N-
<i>dobřes</i>	<i>dobře</i>	Dg-----1A--s-
<i>promluvils</i>	<i>promluvit</i>	VpYS----R-AAPs-
<i>bych</i>	<i>být</i>	Vc-----Ic-
<i>bysme</i>	<i>být</i>	Vc-----Im6
<i>kdybychom</i>	<i>kdyby</i>	J, -----m-
<i>dybychom</i>	<i>dyby_,h_^(^GC**kdyby)</i>	J, -----m-
<i>dybysme</i>	<i>dyby_,h_^(^GC**kdyby)</i>	J, -----m6
<i>abyste</i>	<i>aby</i>	J, -----e-

Table 73: Examples: Aggregates

## 19 Hyphenated Composites

Words written with a hyphen (e.g. *Praha-Hradčany* ‘Prague-Hradcany’) are tokenized into three tokens in the annotation: the part before the hyphen, the hyphen and the part after the hyphen (see Sect. 21). Each part is analyzed separately. We distinguish:

- hyphenated composite of single words,
- hyphenated compound word,
- hyphenated foreign words.

Wordform	Lemma	Tag
<u>Praha-Hradčany</u>	Praha_ ; G	NNFS1-----A----
<u>Praha-Hradčany</u>	Hradčany_ ; G	NNIP1-----A----
<u>na propan-butanovém hořáku</u>	propan	NNIS1-----A----
<u>Karel-Ferdinandova univerzita</u>	Karel_ ; Y	NNMS1-----A----
<u>Karlo-Ferdinandova univerzita</u>	Karlo	S2-----A----
<u>do e-mailu</u>	e-2_^(e-mail)	S2-----A----
<u>do e-mailu</u>	mail	NNIS2-----A----
<u>hraju ping-pong</u>	ping_^(ping-pong)	S2-----A----
<u>hraju ping-pong</u>	pong_^(ping-pong)	SNIS4-----A----
<u>do Tchaj-peje</u>	Tchaj	S2-----A----
<u>do Tchaj-peje</u>	pej-1	SNFS2-----A----
<u>tae-kwon-do</u>	tae_^(tae-kwon-do)	S2-----A----
<u>tae-kwon-do</u>	kwon_^(tae-kwon-do)	S2-----A----
<u>tae-kwon-do</u>	do-2_^(tae-kwon-do)	SNNS1-----A----
<u>cinéma-vérité</u>	cinéma-77	F%-----
<u>cinéma-vérité</u>	vérité-77	F%-----
<u>v Hanty-Mansijsku</u>	Hanty-77	F%-----
<u>v Hanty-Mansijsku</u>	Mansijsk_ ; G	NNIS6-----A----

Table 74: Examples: Hyphenated composites

**Hyphenated composite of single words** is a composite of two or more hyphenated words with a parataxis relation between them. A typical example is the composite of two or more personal or geographical names (e.g. *Anna-Marie*, *Praha-Hradčany* ‘Prague-Hradcany’, *Clam-Gallasův palác* ‘Clam-Gallas palace’), but also composites like *metyl-alkohol* ‘methyl-alcohol’ or *Hewlett-Packard*. The hyphenated parts are analyzed as if they were separate single words. If the first part (before a hyphen) in an adjectival composite is not an adjective (e.g. *propan* ‘propane’ in *propan-butanový hořák* ‘propane-butane burner’, *Karel* ‘Charles’ in *Karel-Ferdinandova univerzita* ‘Charles-Ferdinand University’), it is captured as a noun. See examples in Tab. 74.

**Hyphenated compound word** consists of two or more parts which are hyphenated to create a new word. The part before the hyphen is (usually) an incomplete word (e.g. *anglicko-česká kniha* ‘Czech-English book’) and it is captured as a prefixal segment. If the part after the hyphen is a meaningful word, it is analysed as that wordform. If a hyphenated compound word cannot be decomposed into meaningful words (wordforms), i.e. a hyphenated word only has meaning in its entirety, all parts are captured as segments - this is typically a case of loanwords (e.g. *sci-fi*, *ping-pong*, *wi-fi*, *Tchaj-pej* ‘Taipei’). There can be more than one segment in a hyphenated compound word (e.g. *tae-kwon-do*). See Tab. 74. See more about capturing segments in Sect. 16.

**Foreign-language composites** (e.g. *cinéma-vérité*, *flos-cuculi*) are represented as foreign words (see Section 17). See examples in Tab. 74.

## 20 Typo, Distortion, Misspelling

Intentional misspellings, typos, distortions, as well as frequent errors caused by ignorance of the rules or new codification rules are analyzed as follows.

In the case of a wordform error (e.g. *vidím jí* instead of correct *vidím ji* ‘I see her’ with a short vowel), the wrong wordform is captured as a special wordform of the corresponding paradigm and it is captured at the 15<sup>th</sup> VAR tag position (by the numbers 5–9, but especially 9; Sect. 5.15). If an error, misspelling or distortion is in the whole paradigm (e.g. *vánoce* ‘christmas’ instead of correct *Vánoce* ‘Christmas’ with the capital letter), there is style label *i* in the paradigm lemma. There can also be a reference of the DS type (see Sect. 4.2.4) to the basic variant. See examples in Tab. 75.

Wordform	Lemma	Tag
<i>vidím jí</i>	on-1	PPFS4--3-----6
<i>odvezl mi sem</i>	já	PH-S4--1-----5
<i>pro Božíkovi hosty</i>	Božíkův_-;Y_-^(*2)	AUMP4M-----9
<i>u Lukášové maminky</i>	Lukášův_-;Y_-^(*2)	AUFS2M-----9
<i>místo lokomotiva řekl lomokotiva</i>	lomokotiva_,i_-^(^DS**lokomotiva)	NNFS1-----A---
<i>o vánočích</i>	vánoce_,i_-^(^DS**Vánoce)	NNFP6-----A---

Table 75: Examples: Typo, distortion, misspelling

## 21 Note on Tokenization

Data in the PDT-C corpora are tokenized from delimiter to delimiter. Delimiters are spaces and all non-alphanumeric characters except for the decimal point and decimal comma. Also wordforms containing numbers are tokenized at seams between a string of characters and a string of numbers.

All numbers written with digits are assigned the tag: C=-----.

Non-alphanumeric characters are assigned the tag: Z:-----.

From this simple rule it follows that even units that we normally understand as one word break down into more tokens. E.g. the single quote (') or hyphen (-) splits the string such as *C'tung* or *wi-fi* into three tokens. Because of the number inside the word, the word *12bodový* ‘12-point’ is divided into tokens *12* and *bodový* ‘point’ or the non-standard name of the car *V3ska* is divided into tokens *V*, *3* and *ska*. No attempt is made to put together these separate tokens within morphological annotation. Each token is subject to morphological annotation, although determining the lemma and morphological tag of these “pieces of words” can be difficult (see also Sect. 16 about segments and Sect. 19 about hyphenated composites).

Cf. examples in Tab. 76.

Wordform	Lemma	Tag
5	5	C=-----
35	35	C=-----
3.5	3.5	C=-----
?	?	Z:-----
)	)	Z:-----
%	%	Z:-----
<i>Mao C'tung</i>	C-1_-;Y	NNMXX-----A---
<i>Mao C'tung</i>	tung-1	SNMS1-----A---
<i>12bodový</i>	12	C=-----
<i>12bodový systém</i>	bodový	AAIS1---1A---
<i>V3ska</i>	V-33	Q3-----
<i>V3ska</i>	3	C=-----
<i>V3ska</i>	ska-2	SNFS1-----A---

Table 76: Examples: Tokenization

## 22 Appendix

### 22.1 Detailed Part of Speech (SUBPOS): Quick Reference

- # Sentence boundary (for the “virtual” word ###) MM: JH JH: To by chtelo vysvetlit MM: neumím. Honzo?
- \* (POS: J) Binary mathematical operations as a conjunction (e.g. *plus*; *krát* ‘times’)
- , (POS: J) Conjunction subordinate (e.g. *protože* ‘because’; *že* ‘that’; incl. *aby* ‘in order to’, *kdyby* ‘if’ in all forms)
- : (POS: Z) Punctuation, non-alphanumeric character (e.g. , %)
- = (POS: C) Number written using digits (e.g. 38, 3.5)
- ~ (POS: J) Conjunction connecting main clauses (*a* ‘and’; *ale* ‘but’)
- % (POS: F) Foreign word (e.g. *home made*)
- © (POS: X) Unrecognized word form, unknown (used only by automatic tagger, not in the dictionary and manual annotation)
- } (POS: C) Numeral, written using Roman numerals (e.g. *XIV*)
- 1 (POS: P) Relative possessive pronoun *jehož* and *jehožto* ‘whose’ including wordforms *jehož*, *jejíž*, *jejížto*, *jejichž* ‘whose’, etc.
- 2 (POS: S) Prefixal segment (e.g. *černo-* ‘black-’)
- 3 (POS: Q) Isolated letter
- 4 (POS: P) Relative/interrogative pronoun with agreement GENDER (*jaký* ‘what’, *který* ‘which’, *čí* ‘whose’, *jenž* ‘who’)
- 5 (POS: P) Clitical form of personal pronoun *on* ‘he’ (only *mu*, *ho* ‘him’)
- 6 (POS: P) Personal reflexive pronoun *se* in its long forms (only *sebe*, *sobě*, *sebou* ‘myself’ /‘yourself’ /‘herself’ /‘himself’ in various cases)
- 7 (POS: P) Personal reflexive pronoun *se* in its short (clitic) forms (only wordforms *se*, *si*, *ses*, *sis* ‘myself’ /‘yourself’ /‘herself’ /‘himself’)
- 8 (POS: P) Personal pronoun reflexive possessive *svůj* ‘my’ / ‘your’ / ‘her’ / ‘his’, the POSSGENDER and POSSNUMBER positions are not filled
- 9 (POS: P) Personal pronoun possessive for the 3rd person *jeho* ‘his’, including wordforms *její* ‘her’, *jejich* ‘their’ etc. with the POSSGENDER and POSSNUMBER positions filled
- A (POS: A) Adjective, general (e.g. *velký* ‘big’, *dlouhý* ‘long’)
- B (POS: V) Verb, present (e.g. *pracuje* ‘he-works’) or future form (e.g. *bude* ‘will’, *pojedu* ‘I-will-go’)
- C (POS: A) Adjective, nominal (short, participial) form (e.g. *rád* ‘pleased’, *schopen* ‘able’)
- D (POS: P) Demonstrative pronoun (*ten* ‘this’, ‘that’, *onen* ‘that over there’)
- E (POS: P) Personal pronoun *on* ‘he’ for the 3<sup>rd</sup> person (including wordforms *ona* ‘she’, *jim* ‘them’ etc.), for which the GENDER position is filled
- F (POS: R) Preposition, part of; never appears isolated, always in a phrase (e.g. *nehledě (na)* ‘regardless’, *vzhledem (k)* ‘because of’)

- G (POS: A) Adjective derived from present transgressive form of a verb (e.g. *dělající* ‘working’)
- H (POS: P) Clitical (short) form of personal pronouns *já* ‘I’ and *ty* ‘you’, for which GENDER position is not filled (e.g. *mě* ‘me’, *mi* ‘me’, *ti* ‘you’)
- I (POS: I) Interjection (e.g. *ach*)
- K (POS: P) Indefinite pronoun for which GENDER position is not filled (e.g. *někdo* ‘somebody’, *něco* ‘something’, *bůhvíkdo* ‘whoever’, *cosi* ‘something’)
- L (POS: P) Delimiting pronoun (e.g. *všechnen* ‘all’, *sám* ‘alone’)
- M (POS: A) Adjective derived from verbal past transgressive form (e.g. *udělavší* ‘done’)
- N (POS: N) Noun (e.g. *dům* ‘house’, *Jan* ‘John’)
- O (POS: A) Adjective *svůj* ‘own self’, *nesvůj* ‘not-in-mood’ and *tentam* ‘gone’ with agreement GENDER
- P (POS: P) Personal pronoun *já* ‘I’, *ty* ‘you’, *my* ‘we’ and *vy* ‘you’, for which GENDER position is not filled
- Q (POS: P) Relative/interrogative pronoun for which GENDER position is not filled (*kdo* ‘who’, *co* ‘what’, *cožpak* ‘isn’t-it-true-that’)
- R (POS: R) Preposition (general, without vocalization; e.g. *v* ‘in’, *pod* ‘under’)
- S (POS: P) Personal pronoun possessive *můj* ‘my’, *tvůj* ‘your’, *náš* ‘our’ and *váš* ‘your’ for which the POSSGENDER position is not filled
- T (POS: T) Particle (e.g. *ano* ‘yes’)
- U (POS: A) Adjective possessive, with the masculine ending *-ův* (e.g. *otcův* ‘father’s’) as well as *-in* (e.g. *matčin* ‘mother’s’)
- V (POS: R) Preposition with vocalization *-e* (e.g. *ve* ‘in’, *pode* ‘under’) or *-u* (e.g. *ku* ‘to’)
- W (POS: P) Negative pronoun with agreement GENDER (e.g. *nijaký* ‘no/none’, *ničí* ‘nobody’s’, *žádný* ‘no/none’)
- Y (POS: P) Negative pronoun for which GENDER position is not filled (e.g. *nic* ‘nothing’, *nikdo* ‘nobody’)
- Z (POS: P) Indefinite pronoun with agreement GENDER (e.g. *nějaký* ‘some’, *některý* ‘some’, *něčí* ‘somebody’s’, *číkolik* ‘anybody’s’)
- a (POS: C) Cardinal numeral indefinite for which GENDER and NUMBER position is not filled, incl. interrogative numeral *kolik* ‘how much’ (e.g. *mnoho* ‘much/many’, *málo* ‘little/few’, *tolik* ‘that much/many’, *několik* ‘some (number of)’, *kdovíkolik* ‘who-knows-how-much/many’, *pár* ‘some (number of)’)
- b (POS: D) Adverb without a possibility to form negation and degrees of comparison (e.g. *pozadu* ‘behind’, *naplocho* ‘flatly’, *včera* ‘yesterday’); i.e. positions of the NEGATION and GRADE are not filled
- c (POS: V) Conditional of the verb *být* ‘to be’ (e.g. *by*, including *bych*, *bys*, *byste* ‘would’) that are treated as aggregates
- d (POS: C) Generic numeral definite with agreement GENDER (e.g. *jedny* ‘one-kind’, *dvojí* ‘two-kinds’, *desaterý* ‘ten-kinds’, *patery* ‘five-kinds’)

- e (POS: V) Verb, transgressive present (endings *-e/-ě, -íc, -íce*; e.g. *dělaje* ‘doing’), also archaic present transgressive of perfective verbs (e.g. *udělaje* ‘(he-)having-done’; VAR: 4)
- f (POS: V) Verb, infinitive (e.g. *dělat* ‘to do’)
- g (POS: D) Adverbs forming negation and degrees of comparison (e.g. *dobře* ‘well’, *zajímavě* ‘interestingly’); positions of the GRADE and NEGATION are filled
- h (POS: C) Generic numeral indefinite with agreement GENDER incl. interrogative numeral *kolikery* ‘how-many-kinds’ (e.g. *nejedny* ‘not-only-one-kind’, *tolikery* ‘that-many-kinds’, *několikery* ‘several-kinds’, *několikery* ‘several-kinds’)
- i (POS: V) Verb, imperative form (e.g. *dělej* ‘do!’)
- j (POS: C) Generic numeral definite used as a syntactic noun, with lexical GENDER (e.g. *čtvero* ‘four-kinds-of’, *desatero* ‘ten-kinds-of’)
- k (POS: C) Generic numeral indefinite used as a syntactic noun, with lexical GENDER incl. interrogative numeral *kolikero* ‘how-many-kinds’ (e.g. *několikero* ‘several-kinds-of’, *tolikero* ‘that-many-kinds’)
- l (POS: C) Cardinal numeral definite for which GENDER position is not filled (e.g. *tři* ‘three’, *čtyři* ‘four’, *pět* ‘five’, *půl* ‘half’)
- m (POS: V) Verb, past transgressive (e.g. *udělav* ‘(he-)having-done’)
- n (POS: C) Cardinal numeral definite with agreement GENDER (only *jeden* ‘one’, *dva* ‘two’ and *oba* ‘both’)
- o (POS: C) Multiplicative numeral indefinite incl. interrogative numeral *kolikrát* ‘how-many-times’ (e.g. *mnohokrát* ‘many-times’, *tolikrát* ‘that-many-times’, *několikrát* ‘several-times’, *nejednou* ‘not-only-one-time’)
- p (POS: V) Verb, past participle, active (e.g. *pracoval* ‘(he)-worked’)
- q (POS: V) Verb, past participle, active with the archaic enclitic *-ť* (e.g. *pracovalť* ‘(he)-worked-could-you-imagine-that?’)
- r (POS: C) Ordinal numeral definite, adjective declension without degrees of comparison (e.g. *třetí* ‘third’, *pátý* ‘fifth’)
- s (POS: V) Verb, past participle, passive (e.g. *udělán* ‘done’)
- t (POS: V) Verb, present or future tense, with the archaic enclitic *-ť* with meanings (perhaps) “could-you-imagine-that?” or “but-because” (e.g. *dělámet* ‘(we)-do-could-you-imagine-that?’)
- v (POS: C) Multiplicative numeral definite (e.g. *pětkrát* ‘five-times’, *sedmkrát* ‘seven-times’)
- w (POS: C) Ordinal numeral indefinite, adjective declension without degrees of comparison, incl. interrogative numeral *kolikáty* ‘at-what-position-in-a-sequence’ (e.g. *tolikáty* ‘at-that-position-in-a-sequence’, *několikáty* ‘umpteenth’)
- y (POS: C) Cardinal numeral indefinite with agreement GENDER (only *nejeden* ‘not-only-one’)
- z (POS: C) Cardinal numeral definite with lexical GENDER (e.g. *sto* ‘hundred’, *milion* ‘million’, *nula* ‘zero’, *čtvrt* ‘quarter’)

## 22.2 Categories Relevant for POS and SUBPOS Combinations

The tables of applicability/non-applicability of the tag categories related to SUBPOS value are presented here. For each value of the SUBPOS category (see the list in Sect. 22.1 above and also Sect. 5.2), the major part-of-speech value (POS) is given to which the current subcategory value uniquely belongs. Only abbreviation (B) and segment (S) POS can potentially be associated with any SUBPOS value. The combinations that have occurred in the current version of the dictionary are also listed here.

The SUBPOS category serves as an indicator of applicability/non-applicability of other tag categories (i.e. the categories GENDER, NUMBER, CASE, etc. up to the last category, VAR). Thus each subsequent row of the table shows the applicable tag categories for the given SUBPOS. The values of such categories that occur in the dictionary are listed in the right-hand column.

If a tag category is used for a given SUBPOS, then the tag position never bears the non-applicable value (-). There are only two exceptions: the AGGREGATE and VAR categories (see Sect. 5.14 and 5.15). These two categories are optional for all SUBPOS values. The AGGREGATE and VAR tag positions can bear the non-applicable value (-). [MM: poznamka k formátu tech tabulek. Mohly by být bez tech popisku "Category co-occurrence ... pod kazdou tabulkou. Zdá se nám to tam zbytečné a moc to ten seznam natahuje](#)

# Sentence boundary (for the “virtual” word ###)

Category	Values used
POS	Z

Category co-occurrence for SUBPOS = #

: Punctuation, non-alphanumeric character (e.g. , %)

Category	Values used
POS	Z

Category co-occurrence for SUBPOS = :

© Unrecognized word form, unknown (used only by automatic tagger, not in the dictionary and manual annotation)

Category	Values used
POS	X

Category co-occurrence for SUBPOS = ©

\* Binary mathematical operations as a conjunction (e.g. plus; krát ‘times’)

Category	Values used
POS	J
VAR	- 1

Category co-occurrence for SUBPOS = \*

% Foreign word (e.g. home made)

Category	Values used
POS	F

Category co-occurrence for SUBPOS = %

= Number written using digits (e.g. 38, 3.5)

Category	Values used
POS	C

Category co-occurrence for SUBPOS = =

} Numeral, written using Roman numerals (e.g. *XIV*)

Category	Values used
POS	C
VAR	- 1 2

Category co-occurrence for SUBPOS = }

, Conjunction subordinate (e.g. *protože* ‘because’; *že* ‘that’; incl. *aby* ‘in order to’, *kdyby* ‘if’ in all forms)

Category	Values used
POS	J
AGGREGATE	- c e m s
VAR	- 6 7 8

Category co-occurrence for SUBPOS = ,

~ Conjunction connecting main clauses (*a* ‘and’; *ale* ‘but’)

Category	Values used
POS	J
AGGREGATE	- s
VAR	- 2

Category co-occurrence for SUBPOS = ~

1 Relative possessive pronoun *jehož* and *jehožto* ‘whose’ including wordforms *jehož*, *jejíž*, *jejížto*, *jejichž* ‘whose’, etc.

Category	Values used
POS	P
GENDER	F I M N X Z
NUMBER	D P S X
CASE	1 2 3 4 6 7 X
POSGENDER	F X Z
POSSNUMBER	P S
PERSON	3
VAR	- 2

Category co-occurrence for SUBPOS = 1

2 Prefixal segment (e.g. *černo-* ‘black-’)

Category	Values used
POS	S
NEGATION	A N

Category co-occurrence for SUBPOS = 2

3 Isolated letter

Category	Values used
POS	Q

Category co-occurrence for SUBPOS = 3

4 Relative/interrogative pronoun with agreement GENDER (*jaký* ‘what’, *který* ‘which’, *čí* ‘whose’,

*jenž* ‘who’)

Category	Values used
POS	P
GENDER	F I M N X Y Z
NUMBER	D P S X
CASE	1 2 3 4 5 6 7 X
AGGREGATE	- s
VAR	- 1 2 3 4 6 7 8 9

Category co-occurrence for SUBPOS = 4

- 5 Clitical form of personal pronoun *on* ‘he’ (only *mu, ho* ‘him’)

Category	Values used
POS	P
GENDER	Z
NUMBER	S
CASE	2 3 4
PERSON	3

Category co-occurrence for SUBPOS = 5

- 6 Personal reflexive pronoun *se* in its long forms (only *sebe, sobě, sebou* ‘myself’ /‘yourself’ /‘herself’ /‘himself’ in various cases)

Category	Values used
POS	P
CASE	2 3 4 6 7

Category co-occurrence for SUBPOS = 6

- 7 Personal reflexive pronoun *se* in its short (clitic) forms (only wordforms *se, si, ses, sis* ‘myself’ /‘yourself’ /‘herself’ /‘himself’)

Category	Values used
POS	P
CASE	3 4
AGGREGATE	- s

Category co-occurrence for SUBPOS = 7

- 8 Personal pronoun reflexive possessive *svůj* ‘my’ / ‘your’ / ‘her’ / ‘his’, the POSSGENDER and POSSNUMBER positions are not filled

Category	Values used
POS	P
GENDER	F H I M N X Y Z
NUMBER	D P S
CASE	1 2 3 4 5 6 7
VAR	- 1 6 7

Category co-occurrence for SUBPOS = 8

- 9 Personal pronoun possessive for the 3rd person *jeho* ‘his’, including wordforms *její* ‘her’,

*jejich* ‘their’ etc. with the POSSGENDER and POSSNUMBER positions filled

Category	Values used
POS	P
GENDER	F I M N X Z
NUMBER	D P S X
CASE	1 2 3 4 5 6 7 X
POSSGENDER	F X Z
POSSNUMBER	P S
PERSON	3
VAR	- 6

Category co-occurrence for SUBPOS = 9

A Adjective, general (e.g. *velký* ‘big’, *dlouhý* ‘long’)

Category	Values used
POS	A
GENDER	F I M N X
NUMBER	D P S X
CASE	1 2 3 4 5 6 7 X
GRADE	1 2 3
NEGATION	A N
VAR	- 1 3 5 6 7 8 9 a b

Category co-occurrence for SUBPOS = A

B Verb, present (e.g. *pracuje* ‘he-works’) or future form (e.g. *bude* ‘will’, *pojedu* ‘I-will-go’)

Category	Values used
POS	V
NUMBER	P S
PERSON	1 2 3
TENSE	F P
NEGATION	A N
VOICE	A
ASPECT	B I P
VAR	- 1 2 3 4 5 6 7 8 9

Category co-occurrence for SUBPOS = B

C Adjective, nominal (short, participial) form (e.g. *rád* ‘pleased’, *schopen* ‘able’)

Category	Values used
POS	A
GENDER	F M N Q T Y
NUMBER	P S W
CASE	- 4
NEGATION	A N

Category co-occurrence for SUBPOS = C

- D Demonstrative pronoun (*ten* ‘this’, ‘that’, *onen* ‘that over there’)

Category	Values used
POS	P
GENDER	F I M N X Y Z
NUMBER	D P S X
CASE	1 2 3 4 6 7 X
AGGREGATE	- s
VAR	- 1 2 4 5 6 7 b

Category co-occurrence for SUBPOS = D

- E Personal pronoun *on* ‘he’ for the 3<sup>rd</sup> person (including wordforms *ona* ‘she’, *jim* ‘them’ etc.), for which the GENDER position is filled

Category	Values used
POS	P
GENDER	F I M N X Y Z
NUMBER	P S
CASE	1 2 3 4 6 7
PERSON	3
AGGREGATE	- d n o p z
VAR	- 1 2 6 7

Category co-occurrence for SUBPOS = E

- F Preposition, part of; never appears isolated, always in a phrase (e.g. *nehledě (na)* ‘regardless’, *vzhledem (k)* ‘because of’)

Category	Values used
POS	R

Category co-occurrence for SUBPOS = F

- G Adjective derived from present transgressive form of a verb (e.g. *dělající* ‘working’)

Category	Values used
POS	A
GENDER	F I M N
NUMBER	D P S
CASE	1 2 3 4 5 6 7
NEGATION	A N
VAR	- 6

Category co-occurrence for SUBPOS = G

- H Clitical (short) form of personal pronouns *já* ‘I’ and *ty* ‘you’, for which GENDER position is not filled (e.g. *mě* ‘me’, *mi* ‘me’, *ti* ‘you’)

Category	Values used
POS	P
NUMBER	S
CASE	2 3 4
PERSON	1 2
VAR	- 5 6

Category co-occurrence for SUBPOS = H

I Interjection (e.g. *ach*)

Category	Values used
POS	I
VAR	- 1 6

Category co-occurrence for SUBPOS = I

K Indefinite pronoun for which GENDER position is not filled (e.g. *někdo* ‘somebody’, *něco* ‘something’, *bůhvíkdo* ‘whoever’, *cosi* ‘something’)

Category	Values used
POS	P
CASE	1 2 3 4 5 6 7
VAR	- 6

Category co-occurrence for SUBPOS = K

L Delimiting pronoun (e.g. *všechnen* ‘all’, *sám* ‘alone’)

Category	Values used
POS	P
GENDER	F I M N X Y Z
NUMBER	D P S
CASE	1 2 3 4 5 6 7
VAR	- 1 3 4 5 6 7 8

Category co-occurrence for SUBPOS = L

M Adjective derived from verbal past transgressive form (e.g. *udělavší* ‘done’)

Category	Values used
POS	A
GENDER	F I M N
NUMBER	D P S
CASE	1 2 3 4 5 6 7
NEGATION	A N
VAR	- 1 6 7

Category co-occurrence for SUBPOS = M

N Noun (e.g. *dům* ‘house’, *Jan* ‘John’)

Category	Values used
POS	N
GENDER	F I M N X
NUMBER	D P S X
CASE	1 2 3 4 5 6 7 X
NEGATION	- A N
VAR	- 1 2 3 4 5 6 7 8 9 a b c

Category co-occurrence for SUBPOS = N

O Adjective *svůj* ‘own self’, *nesvůj* ‘not-in-mood’ and *tentam* ‘gone’ with agreement GENDER

Category	Values used
POS	A
GENDER	F I M N Y
NUMBER	P S
VAR	- 1 6

Category co-occurrence for SUBPOS = O

P Personal pronoun *já* ‘I’, *ty* ‘you’, *my* ‘we’ and *vy* ‘you’, for which GENDER position is not filled

Category	Values used
POS	P
NUMBER	P S
CASE	1 2 3 4 5 6 7
PERSON	1 2
AGGREGATE	- s
VAR	- 6 9

Category co-occurrence for SUBPOS = P

Q Relative/interrogative pronoun for which GENDER position is not filled (*kdo* ‘who’, *co* ‘what’, *cožpak* ‘isn’t-it-true-that’)

Category	Values used
POS	P
CASE	1 2 3 4 6 7
AGGREGATE	- N O V Z n o s v z
VAR	- 1 6 9

Category co-occurrence for SUBPOS = Q

R Preposition (general, without vocalization; e.g. *v* ‘in’, *pod* ‘under’)

Category	Values used
POS	R
CASE	1 2 3 4 6 7 X
VAR	- 6 7 a b c

Category co-occurrence for SUBPOS = R

S Personal pronoun possessive *můj* ‘my’, *tviž* ‘your’, *náš* ‘our’ and *váš* ‘your’ for which the POSSGENDER position is not filled

Category	Values used
POS	P
GENDER	F H I M N X Y Z
NUMBER	D P S
CASE	1 2 3 4 5 6 7
POSSNUMBER	P S
PERSON	1 2
VAR	- 1 6 7 9 b

Category co-occurrence for SUBPOS = S

T Particle (e.g. *ano* ‘yes’)

Category	Values used
POS	T
AGGREGATE	- s
VAR	- 1 2 3 6 7 a b

Category co-occurrence for SUBPOS = T

U Adjective possessive, with the masculine ending *-uv* (e.g. *otcív* ‘father’s’) as well as *-in* (e.g.

*matčin* ‘mother’s’)

Category	Values used
POS	A
GENDER	F I M N X
NUMBER	D P S X
CASE	1 2 3 4 5 6 7 X
POSGENDER	F M
VAR	- 1 2 5 6 7 8 9 a b

Category co-occurrence for SUBPOS = U

V Preposition with vocalization *-e* (e.g. *ve* ‘in’, *pode* ‘under’) or *-u* (e.g. *ku* ‘to’)

Category	Values used
POS	R
CASE	2 3 4 6 7
VAR	- 1

Category co-occurrence for SUBPOS = V

W Negative pronoun with agreement GENDER (e.g. *nijaký* ‘no/none’, *ničí* ‘nobody’s’, *žádný* ‘no/none’)

Category	Values used
POS	P
GENDER	F I M N X Y Z
NUMBER	D P S
CASE	1 2 3 4 5 6 7
VAR	- 6 7

Category co-occurrence for SUBPOS = W

Y Negative pronoun for which GENDER position is not filled (e.g. *nic* ‘nothing’, *nikdo* ‘nobody’)

Category	Values used
POS	P
CASE	1 2 3 4 6 7
VAR	- 2 6

Category co-occurrence for SUBPOS = Y

Z Indefinite pronoun with agreement GENDER (e.g. *nějaký* ‘some’, *některý* ‘some’, *něčí* ‘somebody’s’, *číkolik* ‘anybody’s’)

Category	Values used
POS	P
GENDER	F I M N X Y Z
NUMBER	D P S
CASE	1 2 3 4 5 6 7
VAR	- 1 5 6 7

Category co-occurrence for SUBPOS = Z

a Cardinal numeral indefinite for which GENDER and NUMBER position is not filled, incl. interrogative numeral *kolik* ‘how much’ (e.g. *mnoho* ‘much/many’, *málo* ‘little/few’, *tolik* ‘that much/many’, *několik* ‘some (number of)’, *kdovíkolik* ‘who-knows-how-much/many’, *pár*

‘some (number of)’

Category	Values used
POS	C
CASE	1 2 3 4 5 6 7 X
AGGREGATE	- s
VAR	- 1

Category co-occurrence for SUBPOS = a

- b Adverb without a possibility to form negation and degrees of comparison (e.g. *pozadu* ‘behind’, *naplocho* ‘flatly’, *včera* ‘yesterday’); i.e. positions of the NEGATION and GRADE are not filled

Category	Values used
POS	D
AGGREGATE	- s
VAR	- 1 2 3 4 6 7 8 a b

Category co-occurrence for SUBPOS = b

- c Conditional of the verb *být* ‘to be’ (e.g. *by*, including *bych*, *bys*, *byste* ‘would’) that are treated as aggregates

Category	Values used
POS	V
ASPECT	I
AGGREGATE	- c e m s
VAR	- 6

Category co-occurrence for SUBPOS = c

- d Generic numeral definite with agreement GENDER (e.g. *jedny* ‘one-kind’, *dvojí* ‘two-kinds’, *desaterý* ‘ten-kinds’, *patery* ‘five-kinds’)

Category	Values used
POS	C
GENDER	F I M N X Y
NUMBER	D P S
CASE	1 2 3 4 5 6 7
VAR	- 1 2 6 7

Category co-occurrence for SUBPOS = d

- e Verb, transgressive present (endings *-e/-ě, -íč, -íče*; e.g. *dělaje* ‘doing’), also archaic present transgressive of perfective verbs (e.g. *udělaje* ‘(he-)having-done’; VAR: 4)

Category	Values used
POS	V
GENDER	H X Y
NUMBER	P S
NEGATION	A N
ASPECT	B I P
VAR	- 1 2 4 6

Category co-occurrence for SUBPOS = e

- f Verb, infinitive (e.g. *dělat* ‘to do’)

Category	Values used
POS	V
NEGATION	A N
ASPECT	B I P
VAR	- 1 2 3 4 6 7 b

Category co-occurrence for SUBPOS = f

- g Adverbs forming negation and degrees of comparison (e.g. *dobře* ‘well’, *zajímavě* ‘interestingly’); positions of the GRADE and NEGATION are filled

Category	Values used
POS	D
GRADE	1 2 3
NEGATION	A N
VAR	- 1 2 3 4 6 7 b

Category co-occurrence for SUBPOS = g

- h Generic numeral indefinite with agreement GENDER incl. interrogative numeral *kolikery* ‘how-many-kinds’ (e.g. *nejedny* ‘not-only-one-kind’, *tolikery* ‘that-many-kinds’, *několikery* ‘several-kinds’, *několikery* ‘several-kinds’)

Category	Values used
POS	C
GENDER	F I M N X Y
NUMBER	D P S
CASE	1 2 3 4 5 6 7
VAR	- 1 6 7

Category co-occurrence for SUBPOS = h

- i Verb, imperative form (e.g. *dělej* ‘do! ’)

Category	Values used
POS	V
NUMBER	P S
PERSON	1 2 3
NEGATION	A N
ASPECT	B I P
VAR	- 1 2 3 4 5 6 7 8 9 b

Category co-occurrence for SUBPOS = i

- j Generic numeral definite used as a syntactic noun, with lexical GENDER (e.g. *čtvero* ‘four-kinds-of’, *desatero* ‘ten-kinds-of’)

Category	Values used
POS	C
GENDER	N
NUMBER	P S X
CASE	1 2 3 4 5 6 7 X
VAR	- 1

Category co-occurrence for SUBPOS = j

- k Generic numeral indefinite used as a syntactic noun, with lexical GENDER incl. interrogative numeral *kolikero* ‘how-many-kinds’ (e.g. *několikero* ‘several-kinds-of’, *tolikero* ‘that-many-kinds’)

Category	Values used
POS	C
GENDER	N
NUMBER	P S X
CASE	1 2 3 4 5 6 7 X

Category co-occurrence for SUBPOS = k

- l Cardinal numeral definite for which GENDER position is not filled (e.g. *tři* ‘three’, *čtyři* ‘four’, *pět* ‘five’, *půl* ‘half’)

Category	Values used
POS	C
NUMBER	D P S X
CASE	1 2 3 4 5 6 7 X
VAR	- 1 2 6

Category co-occurrence for SUBPOS = l

- m Verb, past transgressive (e.g. *udělav* ‘(he-)having-done’)

Category	Values used
POS	V
GENDER	H X Y
NUMBER	P S
NEGATION	A N
ASPECT	B I P
VAR	- 1 2 6

Category co-occurrence for SUBPOS = m

- n Cardinal numeral definite with agreement GENDER (only *jeden* ‘one’, *dva* ‘two’ and *oba* ‘both’)

Category	Values used
POS	C
GENDER	F H I M N X Y Z
NUMBER	D P S X
CASE	1 2 3 4 5 6 7 X
VAR	- 1 6 8

Category co-occurrence for SUBPOS = n

- o Multiplicative numeral indefinite incl. interrogative numeral *kolikrát* ‘how-many-times’ (e.g. *mnohokrát* ‘many-times’, *tolikrát* ‘that-many-times’, *několikrát* ‘several-times’, *nejednou* ‘not-only-one-time’)

Category	Values used
POS	C
VAR	- 1

Category co-occurrence for SUBPOS = o

- p Verb, past participle, active (e.g. *pracoval* '(he)-worked')

Category	Values used
POS	V
GENDER	F M N Q T Y
NUMBER	P S W
TENSE	R
NEGATION	A N
VOICE	A
ASPECT	B I P
AGGREGATE	- s
VAR	- 1 2 3 6 7 8 9

Category co-occurrence for SUBPOS = p

- q Verb, past participle, active with the archaic enclitic *-ť* (e.g. *pracovalť* '(he)-worked-could-you-imagine-that?)

Category	Values used
POS	V
GENDER	M N Q T Y
NUMBER	P S W
TENSE	R
NEGATION	A N
VOICE	A
ASPECT	B I P
VAR	2 3 4 5 6

Category co-occurrence for SUBPOS = q

- r Ordinal numeral definite, adjective declension without degrees of comparison (e.g. *třetí* ‘third’, *pátý* ‘fifth’)

Category	Values used
POS	C
GENDER	F I M N
NUMBER	D P S
CASE	1 2 3 4 5 6 7
VAR	- 6 7

Category co-occurrence for SUBPOS = r

- s Verb, past participle, passive (e.g. *udělán* 'done')

Category	Values used
POS	V
GENDER	F M N Q T Y
NUMBER	P S W
CASE	- 4
TENSE	H X
NEGATION	A N
VOICE	P
ASPECT	B I P
AGGREGATE	- s
VAR	- 1 2 5 6 7 8

Category co-occurrence for SUBPOS = s

- t Verb, present or future tense, with the archaic enclitic *-ť* with meanings (perhaps) “could-you-imagine-that?” or “but-because” (e.g. *dělámet* '(we)-do-could-you-imagine-that?)

Category	Values used
POS	V
NUMBER	P S
PERSON	1 2 3
TENSE	F P
NEGATION	A N
VOICE	A
ASPECT	B I P
VAR	1 2 3 4 5 6 7 8 9

Category co-occurrence for SUBPOS = t

- v Multiplicative numeral definite (e.g. *pětkrát* ‘five-times’, *sedmkrát* ‘seven-times’)

Category	Values used
POS	C
VAR	- 1 7

Category co-occurrence for SUBPOS = v

- w Ordinal numeral indefinite, adjective declension without degrees of comparison, incl. interrogative numeral *kolikáty* ‘at-what-position-in-a-sequence’ (e.g. *tolikáty* ‘at-that-position-in-a-sequence’, *několikáty* ‘umpteenth’)

Category	Values used
POS	C
GENDER	F I M N X Y Z
NUMBER	D P S
CASE	1 2 3 4 5 6 7
VAR	- 6 7

Category co-occurrence for SUBPOS = w

- y Cardinal numeral indefinite with agreement GENDER (only *nejeden* ‘not-only-one’)

Category	Values used
POS	C
GENDER	F I M N Y Z
NUMBER	S
CASE	1 2 3 4 5 6 7

Category co-occurrence for SUBPOS = y

- z Cardinal numeral definite with lexical GENDER (e.g. *sto* ‘hundred’, *milion* ‘million’, *nula* ‘zero’, *čtvrt* ‘quarter’)

Category	Values used
POS	C
GENDER	F I N
NUMBER	P S X
CASE	1 2 3 4 5 6 7 X
VAR	- 1 2 6 b

Category co-occurrence for SUBPOS = z

### SUBPOS co-occurrence tables for POS B (Abbreviations)

- A Abbreviation of adjective (AA)

Category	Values used
POS	B
GENDER	X
NUMBER	X
CASE	X
GRADE	1
NEGATION	A

Category co-occurrence for SUBPOS = A

- N Abbreviation of noun (NN)

Category	Values used
POS	B
GENDER	F I M N X
NUMBER	X
CASE	X
NEGATION	A

Category co-occurrence for SUBPOS = N

- b Abbreviation of adverb (Db)

Category	Values used
POS	B

Category co-occurrence for SUBPOS = b

- ~ Abbreviation of conjunction (J~)

Category	Values used
POS	B

Category co-occurrence for SUBPOS = ~

### SUBPOS co-occurrence tables for POS S (Segments)

**A Postfixal segment of adjective (AA)**

Category	Values used
POS	S
GENDER	F I M N X
NUMBER	D P S X
CASE	1 2 3 4 5 6 7 X
GRADE	1 2
NEGATION	A
VAR	- 6 7

Category co-occurrence for SUBPOS = A

**N Postfixal segment of noun (NN)**

Category	Values used
POS	S
GENDER	F I M N X
NUMBER	P S X
CASE	1 2 3 4 5 6 7 X
NEGATION	A
VAR	- 1 6

Category co-occurrence for SUBPOS = N

**b Postfixal segment of adverb (Db)**

Category	Values used
POS	S

Category co-occurrence for SUBPOS = b

**1 Postfixal segment of numeral (C1)**

Category	Values used
POS	S
NUMBER	X
CASE	X

Category co-occurrence for SUBPOS = 1

**2 Prefixal segment**

Category	Values used
POS	S
NEGATION	A N

Category co-occurrence for SUBPOS = 2