

Supplementary Information – Morphological convergence as on-line lexical analogy

Péter Rácz, Clay Beckner, Jen Hay, and Janet B. Pierrehumbert

February 18, 2020

This is an Sweave file. The Rnw source file contains the code used to compile the pdf output file. It also contains additional code sections that help replicate the results. The files also contain pointers to further detail in the supplementary depository that accompanies the paper.

In this supplementary information, we detail (1) the structure of our nonce verb stimuli, the setup of (2) the Generalized Context Model (GCM) and (3) the Minimal Generalization Learner (MGL), and (4) how these models compare. We illustrate the models using fits on our *baseline* data.

We also discuss model selection in analysing our regression data (5), how the two learning models compare on our baseline data (6), and how they compare on our ESP post-test data (7).

Sections (2) and (3) also appear in the main text of Rácz et al. (2020), as Appendices A and B. However, this material is included here so that this Supplement can serve as a standalone document for readers who are interested in the technical details of our study. Additional material is reproduced from the main text, revisiting specific regression model summaries and additional methods we used to investigate contributions from the GCM and MGL.

1 Regular and irregular verbs in English

Four irregular verb classes were defined for our stimuli, based on the vowel alternation and affixation processes that apply to the stem:

- DROVE ([aɪ]/[i] → [oʊ])
- SANG ([ɪ] → [æ])
- KEPT ([i] → [ɛ]Ct)
- BURNT ([ɜ]/[ɛ]/[ɪ] → [ɜ]/[ɛ]/[ɪ]Ct)

In the baseline experiment, participants are also tested on monosyllabic verbs that do not change in the irregular past tense (e.g. *cut*, *hit*). These forms are strong outliers and are not reported in the data.

The stimuli for the ESP experiment, 156 nonce verbs across four classes, are shown in Table 1 in the main text. Our formal criteria are based on the behavior of existing

verbs in English and their categorization by Bybee & Slobin (1982), Moder (1992), and Albright & Hayes (2003). We made some adjustments to their categories. For instance, our DROVE class is a generalized version of Moder’s class 4, described using the [aɪ]→[oʊ] alternation but also including *weave*. (In contrast, Albright & Hayes restrict this class to this alternation.) The verb classes could be defined in a number of ways and still be mostly consistent with the work of Moder, for instance. We have trialed a set of possible small changes and found that they have no major effect on our overall results.

2 Implementation of the Generalized Context Model

2.1 Outline

Our implementation of the GCM evaluates the competition between two categories, *regular* and *irregular*, for each nonce verb base form. The framework of Nosofsky (1990) is adapted to morphophonology by using a segmental similarity calculation based on natural classes (Frisch et al., 2004). The same treatment of segmental similarity is used in the implementations of the GCM in Albright & Hayes (2003) and Dawdy-Hesterberg & Pierrehumbert (2014). We build on Dawdy-Hesterberg & Pierrehumbert (2014) in that we define our categories based on formal similarity.

The implementation of the GCM is laid out in `models/gcm/gcm.html`.

2.2 Training data

Participants are presented with a sequence of nonce verb base forms, and have to pick either a regular or an irregular past tense form for each. The irregular past tense form is pre-determined by the class for the stem, so that, for a given verb, the participants can only choose between the regular past tense form and the irregular past tense form we assigned to the verb. (So, for instance, for *splive*, a verb in the DROVE class, they can choose either *splived* or *splove*, but not *splift* or *sploven*, etc.) For a given class (such as DROVE verbs), the GCM has a choice between two sets of verb types.

The irregular set consists of verb types in CELEX that form their past tense according to the pattern captured by the class (such as an {[aɪ],[i]} → [oʊ] alternation). The regular set consists of verb types that have base forms that are similar to these irregular forms but have a regular *-ed* past tense form as well as miscellaneous regular verbs – those that do not belong to our schemata. We narrow the regular set to monosyllabic forms. However, all polysyllabic irregular forms that could serve as a point of comparison are compounds based on monosyllabic forms (an example is *overwrite*, a compound form of irregular *write*). (A compound form might be more regular than a simplex form, but CELEX will list both the regular and the irregular variant in both cases.) Table 1 shows the descriptions of the verb classes using regular expressions.

Our starting point for the training set, following Albright & Hayes (2003), is the list of verbs in the CELEX corpus (Baayen et al. 1993, based on Sinclair 1987) with a token frequency of 10 or above, encompassing 3156 forms. However, similarity requirements restrict the respective training sets. We use the *DISC* transcription in which each contrastive segment of English is represented by a unique character (`dr2v` equals [draɪv]).

Table 2 shows the number of verbs in CELEX that were used as training sets for our verb classes. The irregular set consists of forms that match the schema and are irregular.

	input		alternation
class	regular expression	IPA	
DROVE	[i2] [zvd1tnk]\$	{i,ai} + {z,v,d,l,t,n,k} ##	{i,ai} → ou
SANG	I(m N Nk)\$	{i} + {m,ŋ,ŋk} ##	i → æ
KEPT	i[lpnm]\$	{i} + {l,p,n,m} ##	i → εCt
BURNT	[3EI]n1\$	{3,ε,i} + {n,l} ##	{3,ε,i} → {3,ε,i}Ct

Table 1: Descriptions of verb classes in the GCM. ‘C’ marks any consonant.

The regular set contains schema matches which are regular, in addition to miscellaneous regulars. The miscellaneous set consists of monosyllabic verbs that do not belong to any of the schemata and are regular. These are included in the regular training set of EACH verb class.

The model calculates the similarity of a given nonce verb to the regular and the irregular set. Comparisons to stems in other classes are not calculated, as past tense markings for these classes were not available to participants in the forced-choice tasks.

	verb class	irregular set	regular set	miscellaneous regular verbs
1	BURNT	6	42	1218
2	DROVE	14	83	1218
3	KEPT	12	31	1218
4	SANG	8	13	1218
5	sum	40	169	1218

Table 2: Number of forms in each verb class, GCM training data

2.3 Estimation

To calculate the similarity between two words, we first compute their dissimilarity. This is achieved using the string-edit (Levenshtein) distance, which is the smallest number of changes needed to transform one word into the other. Costs range from 0 (the corresponding segments are identical) to 1 (inserting or deleting an entire segment). Following Albright & Hayes (2003) and Dawdy-Hesterberg & Pierrehumbert (2014), costs between 0 and 1 are assigned to corresponding segments that are not identical, based on how much the segments differ.

All parts of the word are weighted equally, because although there is evidence that past tense formation in English is predominantly driven by overlaps in word endings, onsets also play a role. (cf. the predominance of *s*(+stop) onsets in irregular verbs forming the past tense with a vowel change, e.g. *stink*, *sink*, etc. – see Bybee & Moder 1983).

The following transformation, originating with Nosofsky (1990), is used to convert dissimilarity into similarity:

$$\eta_{ij} = \exp(-d_{ij}/s)^p$$

In the equation above, η_{ij} represents the similarity between form i and form j , while d_{ij} is the dissimilarity between the two forms. s and p are free parameters.

We explored a range of parameter settings and use $s = 0.9$ and $p = 1$ which provide the best model fit on the baseline data. (In contrast, Albright and Hayes use $s = 0.4$ and $p = 1$.)

When p is set to 1, as here, the similarity function is an exponential, rather than a Gaussian, function of the dissimilarity. The weighting parameter s controls how quickly the similarity decreases as the difference (or distance) between the forms increases. When s is small, the behavior of the model will be dominated by the small group of instances that differ very little from any given novel form. As it becomes larger, instances that differ more increase their influence on the overall model behavior (Nosofsky, 1990; Nakisa et al., 2001; Albright & Hayes, 2003; Dawdy-Hesterberg & Pierrehumbert, 2014). Thus, s effectively controls the size of the set of verbs that will be taken into account in determining the support for the regular versus the irregular outcome.

The overall similarity S_{iC_J} of a test form i to a set C_J is calculated by summing the similarity η_{ij} of each member j of class C_J to the test form i , and dividing by the summed similarity η_{ik} of each member k of class C_K (the class of all stored forms) to the test form i . This calculation is summarized in the following equation.

$$S_{iC_J} = \frac{\sum_{j \in C_J} \eta_{ij}}{\sum_{k \in C_K} \eta_{ik}}$$

2.4 Output format

The overall score used in our analyses is the *regularity score*, which is the complement to the *irregularity score* and reaches a maximum of 1.0 when the output is most likely to be regular. Unlike Dawdy-Hesterberg & Pierrehumbert (2014), there is no decision rule on top of the scoring, such that any form that is more likely than not to be regular is predicted to surface as regular all the time. This specific decision rule is statistically optimal, and was imposed in Dawdy-Hesterberg & Pierrehumbert (2014) in order to determine the ceiling performance for a computational model. The present paper, in contrast, analyzes data aggregated across human participants with differing decision thresholds. As discussed in Schumacher et al. (2014) and Schumacher & Pierrehumbert (2017), the input-output relationship in such aggregated data are typically reported to be nearly probability-matching.

We standardize the regularity score to match the range of participant responses: $[0,1]$. The modified score is interpretable as the probability that the outcome will be regular in aggregated data. It is also appropriate to attribute this type of gradience to people’s initial expectations about other people’s behavior, on the assumption that people realistically encode the variability they have encountered.

2.5 Example: *splive*

The nonce form *splive* belongs to the DROVE class in our model. The two past forms of *splive* in the experiment are regular *splived* and irregular *splove*. It is compared to 1301 regular verbs – these are 83 verbs that match the DROVE schema (e.g. *side*, *hive*, *line*) and 1218 miscellaneous verbs. It is also compared to 14 irregular verbs (e.g. *drive*, *stride*,

smite) in this class. Overall, it is more similar to the regular set: its *regularity score* is 0.57.

3 Implementation of the Minimal Generalization Learner

3.1 Outline

The Minimal Generalization Learner is an algorithm for forming input-output rules of varying generality, which then compete to generate the output.

The Minimal Generalization Learner is implemented here from materials made available by Albright and Hayes (Albright & Hayes, 2003). These include their Segmental Similarity Calculator, implementing the natural class based similarity metric due to Frisch et al. (2004), also used in the GCM implementation. Due to issues with the MGL code, we had to fit the MGL separately for participants.

The implementation of the MGL is laid out in `models/mgl/mgl.html`.

3.2 Training data

For our model fitted on our baseline nonce word stimuli, the MGL is trained on regular and irregular English verbs with a minimum frequency cutoff of 10 in CELEX (Baayen et al., 1993), encompassing 4160 past/present verb transcriptions.

The MGL builds rules based on all verb forms in CELEX with a token frequency of 10 or above. However the structural descriptions of the resulting rules do not cover all these forms. Table 3 shows the number of unique forms covered by the structural descriptions of the ‘regular’ and ‘irregular’ rules that are relevant to each class.

The MGL generates multiple possible past tense forms for each nonce verb. We only consider those rules *relevant* that generate the past tense forms that appear in the experiment (e.g. *splive* : *splived* / *splove*). There is at most one relevant regular rule and one relevant irregular rule for one verb, but multiple rules can generate the (ir)regular forms for each verb class. We return to this in the next section.

Note that the sets of exceptions and related forms of each rule can overlap, both respectively and with each other. As a consequence, the MGL rules apply to fewer forms than apparent from the table: 456 (instead of 617) in total.

	category	rule.type	related forms	exceptions
1	burnt	irregular	24	41
2	kept	irregular	10	9
3	sang	irregular	11	77
4	burnt	regular	38	21
5	drove	regular	29	27
6	kept	regular	67	41
7	sang	regular	47	38

Table 3: Number of forms in each verb class, MGL training data

3.3 Estimation

The MGL begins by considering the relationship between each verb and its past tense as a ‘rule’. For each pair of verbs in the training data, it then attempts to create a more general rule. It does so by aligning the wordforms and analyzing shared phonetic features. For example, merging the word-specific rules for *ring/rang* and for *stink/stank* yields a more general rule that expresses the information that they share: [ɪ] → [æ] / [+coronal, - cont] ____ [ŋ]. Each rule inferred in this way is then further generalized on the basis of more comparisons; for instance, taking note of *swim/swam* expands the [ɪ] → [æ] rule to specify that it occurs before all [+nasal] consonants.

The structural description for each rule has a *scope*, which is the number of verbs conforming to the description, to which the rule might apply. The number of *hits* is the number of such verbs where the rule generates the correct output. In our example, *think* and *blink* fall in the scope of the rule, but they are not hits, because their past tenses display other patterns (*thought* and *blinked*). The *raw confidence* of the rule is the ratio of hits to scope:

$$\text{Raw confidence} = \frac{\text{hits}}{\text{scope}}$$

The raw confidence is 1.0 if the rule applies to all forms that meet its structural description. It is less than 1.0 if some forms meeting its structural description have past tenses other than that predicted by the structural change. Raw confidence values of 0 are not found, because a rule needs to apply to two or more examples to be posited in the first place.

The MGL raw confidence metric is adjusted on the basis of user-specified confidence limits, to generate an *adjusted confidence score* that takes into account the amount and distribution of available data. The MGL’s lower limit affects how much confidence is assigned to rules that have a small number of instances; generalizations that are based on a smaller number of word types are penalized. The MGL’s upper limit curtails the application of seemingly general rules which are in fact driven by a more specific rule (Albright & Hayes, 2002). The MGL is implemented here with its default settings, with the exception of the algorithm’s confidence limits. We implement the MGL with lower and upper confidence limits of 55% and 95%, respectively, since these values afford the best fit to English verb data in Albright & Hayes (2003).

Note that the MGL algorithm automatically groups together verbs on the basis of shared phonological properties; thus, verbs are most likely to form strong generalizations with other verbs that share the same onset or rhyme. Attempts to merge diverse wordforms under a single generalization would be more likely to incur penalties (i.e. exceptions). This feature of the MGL is important for comparing with the methods of the GCM. Both algorithms allow for category-specific similarities to play a role in rule formation.

3.4 Output format

Recall that in both our baseline and ESP experiment the trial task is prompted by a stem and offers a choice between a regular form and a specific irregular form, presented orthographically. In order to model this choice, we take the MGL rule for the stem that

outputs the regular form (the *relevant regular rule*) and the rule that outputs the specific irregular form (the *relevant irregular rule*). If several regular / irregular rules generate the same form, we take the one with the highest *adjusted confidence*, following Albright & Hayes (2003). We use these rules to calculate the form’s *relative (adjusted) confidence*.

Out of 156 test verbs in the ESP post-test, the CELEX-trained MGL generates a *relevant regular rule* for every verb. It does not generate a *relevant irregular rule* for 28 verbs. These are all nonce verbs in the KEPT category (see Section 4.1 in the main text). In this category, irregular forms are derived from the stem through a vowel change (e.g. *greel* → *greelt*). The verbs missing the relevant irregular rule all have bases ending in <m>, <n>, or <l>. This is because, in our implementation, there is an insufficient number of verb types in the training set to support these irregular rules. Decreasing the cutoff criterion for the model leads to the generation of more of the currently ‘missing’ irregular rules, but the overall model fit becomes worse (cf. below). Therefore, we keep the cutoff criterion and assume that the adjusted confidence of the irregular rule for these 28 verbs is zero. We then take the relative confidence of the regular rule as compared to the regular *and* the irregular rule for each verb and took this as the adjusted regular confidence of the given verb. (If the irregular rule is missing, the value of this adjusted regular confidence is 1.) This is given by the following equation:

$$\text{Relative (adjusted) confidence} = \frac{\text{adjusted confidence of relevant regular rule}}{\text{adj. conf. reg. rule} + \text{adj. conf. relevant irreg. rule}}$$

This relative adjusted confidence represents the MGL regularity score for an item, to be compared against the regularity score from the GCM (see Section 2).

3.5 Example: *splive*

The two past forms of *splive* in the experiment are regular *splived* and irregular *splove*. The *relevant regular rule* that generates the regular past tense is ‘ $\emptyset \rightarrow [d] / \{\delta, \int, \theta, \mathfrak{z}, f, s, v, z\} __$ ’. The *structural description* indicates that this is a suffixation rule that can apply to forms that *end in an anterior fricative* (a natural class in our feature system). The *raw confidence* of this rule is 0.98. This is because this rule applies to most forms in its scope (698/712). The *adjusted confidence* is very similar: 0.968. This is because this rule applies to a large number of forms overall. The *relevant irregular rule* that generates the irregular form is ‘ $[aɪ] \rightarrow [oʊ] / \{\delta, \mathfrak{z}, d\mathfrak{z}, d, l, n, r, z\} __v$ ’. It applies to *[aɪ]* in the nucleus *preceded by voiced anterior consonant and followed by [v]*. In CELEX, the rule applies to three forms (*drive*, *strive*, *dive*) and fails to apply to five (*arrive*, *thrive*, *contrive*, *rive*, *connive*). Its *raw confidence* is 0.375. Its *adjusted confidence* is slightly lower (0.366). This is because it applies to a smaller number of forms overall. The *relative (adjusted) confidence* of the predicted regularity of *splive* is $0.98 / (0.98 + 0.37) = 0.73$. The respective rescaled value is 0.363.

4 Further notes on the GCM and the MGL

The training set for the GCM is focused initially; it is grouped into four verb classes and based on formal similarity with base forms of existing irregular verbs playing a role within

each class. This, in effect, assumes lexical gangs as a starting point (Alegre & Gordon, 1999). In contrast, the MGL starts establishing rules across all forms in the starting dictionary.

As we see, however, the MGL also focuses the training set. The structural descriptions of the forms only cover a fraction of all verbs in the starting dictionary, organized by shared groups of segments. While disjunction increases the power of a theory (increasing the set of acceptable classes), disjunct classes frequently serve as input for phonological processes (Mielke, 2008). In the MGL, the disjunction emerges despite the lack of an initial specification: while one rule accounts for the specific behavior of a given form, multiple rules cover a given class. The number of regular and irregular rules posited by the MGL for each verb class can be seen in Table 4.

	category	rule.type	related forms	exceptions
1	burnt	irregular	24	41
2	kept	irregular	10	9
3	sang	irregular	11	77
4	burnt	regular	38	21
5	drove	regular	29	27
6	kept	regular	67	41
7	sang	regular	47	38

Table 4: The number of unique rules posited by the MGL in each verb class

Table 5 shows the list of (irregular and regular) rules for one of the classes, the KEPT class. We choose this class as it provides a good illustration of the estimation process. An MGL rule describes an $A \rightarrow B$ alternation in an XAY environment. Since the KEPT rules are all suffixing, the following environment (Y) is empty. For the preceding environment (X), we give a list of segments that are matched by the structural description if it applies to four segments or fewer. Otherwise, we give a phonological description. So, for instance, the first rule turns [il] to [ɛlt] following any obstruent. Note that in some cases, the MGL generates multiple rules for the same alternation (such as $ip \rightarrow \epsilon pt$). As discussed in Section 3.4, in such instances, we choose the candidate rule with the highest adjusted confidence in order to calculate the MGL regularity score. (Fricatives and nasal consonants constitute a natural class in the segmental feature system used by Albright & Hayes 2003.)

Inspection of the table shows that the MGL fails to generate some expected irregular rules; in the KEPT class, the MGL omits rules for any stems ending in [in] or [im]. Under the MGL, there is not enough lexical support for irregular rules of these types, and yet participants do irregularize nonce forms like *cheen* and *kleem* in the forced choice task. In this respect, the variegated similarity of the non-natural GCM class (which can be used to categorize *cheen*) allows for a better approximation of participant behavior. Note also that in some cases, the MGL generates an overly-specific rule; the rule for stems ending in [il] applies only after obstruents, whereas participants also irregularize KEPT stems with a sonorant in the onset (*greel/greelt*).

The rule structure of the MGL likely follows from its parameter settings, which were tuned to optimize its performance on our fairly irregular-looking nonce word stimuli. It requires further work to determine the extent to which the MGL fit on individuals would improve if the MGL parameters were customized for individual participants.

alternation	preceding environment
il → εlt	obstruent
ip → εpt	{j, l, r, w}
ip → εpt	dorsal consonant
ip → εpt	l
ip → εpt	any consonant
il → εlt	fricative or nasal consonant
∅ → d	alveolars
∅ → d	{b, m}
∅ → d	alveolar fricative or nasal consonant
∅ → d	{b, m, p}

Table 5: The unique rules posited for the KEPT class

5 Regression modeling procedures

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.108	0.161	-0.672	0.501
verb_baseline_average	8.259	0.686	12.040	0.000
bot_lexical_typicalityrandom	0.043	0.169	0.253	0.800
bot_lexical_typicalityreversed	0.035	0.173	0.204	0.839
participant_prereg_average	7.125	0.563	12.667	0.000
bot_regularisation_rate-40%	-0.456	0.185	-2.468	0.014
bot_regularisation_rate+40%	0.549	0.166	3.307	0.001
verb_baseline_average:bot_lexical_typicalityrandom	-2.402	0.927	-2.591	0.010
verb_baseline_average:bot_lexical_typicalityreversed	-2.127	0.955	-2.226	0.026

Table 6: Experimental Data, Post-test Regression Model Summary. Model formula: post-test regular response ~ verb baseline mean x lexical typicality + participant pretest mean + regularization shift + (1 + verb.baseline.mean | participant) + (1 | verb)

Our regression model selection in this article proceeded as follows. Except where stated otherwise, we fit logistic mixed-effects regression models, using R’s lme4 package (Bates et al., 2015). Non-significant predictors are removed from our models, as determined by a series of comparisons of nested models, using likelihood ratio tests in R’s anova function. However, our model tables (e.g., Tables 5 and 6 in the main text) include the significance levels generated by lme4 summaries.

We investigate maximal random effects structures, allowing for random intercepts for participants and stimulus items (verbs), as well as random slopes for all within-unit factors. We compare nested random effects structures as above, and retain any random slopes and intercepts which are supported by likelihood ratio tests (Baayen et al., 2008).

In the model in Table 5 in the main text, repeated here as Table 6, we elect to mean center continuous variables to counteract collinearity that becomes evident when including the interaction between and *lexical typicality*¹. Although some debate exists regarding the practice of transforming variables (Belsley, 1991; Echambadi & Hess, 2007), centering variables is common practice to remove nonessential ill conditioning, that is, relationships

¹In the SI, all p values rounded to 0.00 should be taken to indicate $p < 0.001$.

that inevitably exist between main and product terms (Aiken et al., 1991; Jaccard & Turrisi, 2003; Jaeger, 2008). We note that if *verb baseline mean* is *not* centered, collinearity is high (VIF=15.79) in this model. A non-centered model also yields a significant main effect for *lexical typicality*, with increased post-test regularization in the random and reversed conditions. However, we elect to take a more conservative approach by not including a potentially spurious main effect. Note that other than this main effect difference, models with and without centering yield results that are qualitatively indistinguishable; all other main effects and interactions are present in either case. The collinearity for the centered model (Table 5 in the main text) is acceptable, with a maximum VIF score of 2.9.

Since it is unclear whether the absence of a difference between random and reversed conditions is due to a lack of data, we refit the model with the baseline average : typicality interaction on the random and reversed conditions only. Then we calculated the Bayes Factor of the null model (with no interaction) over the full model (with interaction) using Bayesian Information Criteria calculated by lme4. The result strongly supports the null model over the full model, indicating that the preference for the null model is *not* due to the lack of evidence to reject it.

6 Fits of the two categorization models on the baseline data

We fit the CELEX-GCM and the CELEX-MGL on the baseline data to see how well they predict regularization responses by the baseline participants. The CELEX-MGL outperforms the CELEX-GCM, if slightly. Both models contribute independently to explaining variation in the data.

The concordance indices of the CELEX-MGL (0.586) and the CELEX-MGL (0.581) on the baseline data are similar, with some edge retained by the MGL, as in the case of (Albright & Hayes, 2003). Our concordance indices are calculated using the somers2 function in R’s *Hmisc* package (Harrell Jr, 2008).

Since the regularity scores of the GCM and the MGL for the nonce verbs are correlated ($r = 0.478$), we use residualization to see whether the two models make separate significant contributions to explaining variation in the baseline dataset (Section 6.3 in the main text.)

We fit two simple linear regression models; no random effects are explored in these models, since we are merely predicting algorithmic scores for each item. We rescale and center GCM and MGL scores for the model fits.

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.525	0.124	-4.219	0.000
rescaled_CELEX_mgl_feature_distance	1.347	0.093	14.482	0.000
rescaled_CELEX_gcm_feature_distance_residual	1.423	0.128	11.100	0.000

Table 7: Model 3: MGL scores and GCM residuals, baseline data

Model 1 predicts the GCM scores from the MGL scores (M1: GCM regular score \sim MGL regular score). Model 2 predicts the MGL scores from the GCM scores (M2: MGL regular score \sim GCM regular score). In two additional models, we then use the residuals from each model (M1, M2) as an estimate of the variation present in one score that cannot be attributed to the other. First, we use the residuals of Model 1 in another, mixed-effects logistic regression model along with the MGL scores to predict regular responses in the

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-1.225	0.160	-7.682	0.000
rescaled_CELEX_gcm_feature_distance	2.391	0.149	16.063	0.000
rescaled_CELEX_mgl_feature_distance_residual	1.008	0.090	11.154	0.000

Table 8: Model 4: GCM scores and MGL residuals, baseline data

baseline task (M3: regular response \sim MGL score + GCM residual + (1 + MGL score + GCM residual|participant) + (1 | verb)). We then residualize for the MGL. We use the residuals of Model 2 along with the GCM scores in a similar logistic regression model (M4: regular response \sim GCM score + MGL residual + (1 + GCM score + MGL residual |participant)).

The model summaries can be seen in Table 7 (M3) and Table 8 (M4). The residuals remain significant predictors in both models. Note that while this paper reports values fit to the baseline experiment data, the same pattern of results also holds for models fit to the pre-test data.

Although residualization is sometimes critiqued when used as a stopgap for collinearity or sign changes in a regression model (Wurm & Fisicaro, 2014), these concerns do not apply in the present case. We use residualization as a tool to test whether one metric contributes beyond what is covered by a different metric of the same attribute, as in Baayen et al. (2006), and as also acknowledged in Wurm & Fisicaro (2014).

7 Model fits on the post-test data

We fit the CELEX-GCM, the CELEX-MGL, as well as the Individual-GCM and the Individual-MGL on the post-test data. Recall that the CELEX-based models were trained on verbs in CELEX. The individual-based models were trained for each participant, on CELEX verbs as well as the co-player’s responses in the ESP test. They were tested on the post-test for each participant, individually.

In the paper, we list four techniques of evaluating model fit. We go through these in turn.

7.1 Comparing concordance indices

First, we calculated concordance indices for the CELEX- and individual-models on the post-test data of the ESP experiment.

Table 9 summarizes the concordance indices between our data and the various categorization models. The CELEX-MGL outperforms the CELEX-GCM on the post-test data to some extent, much like in the case of our baseline data. The individual-GCM considerably improves on the CELEX-GCM in terms of explaining the post-test data. The individual-MGL also improves on the CELEX-MGL, but to a lesser extent.

The issue here is that a concordance index is not a statistical test, and does not provide a margin of error. The trends are apparent: the individual-GCM is more of an improvement than the individual-MGL. But the strength of these trends is less clear.

	training data	method	test data	C
CELEX-GCM	real verbs	‘analogy’	verbs in baseline	0.581
CELEX-MGL	real verbs	‘rules’	verbs in baseline	0.586
CELEX-GCM	real verbs	‘analogy’	verbs in post-test	0.602
individual-GCM	real verbs + verbs seen in ESP	‘analogy’	verbs in post-test	0.681
CELEX-MGL	real verbs	‘rules’	verbs in post-test	0.61
individual-MGL	real verbs + verbs seen in ESP	‘rules’	verbs in post-test	0.625

Table 9: Summary of concordance indices between our data and the categorization models. The values above the line show the concordance indices for the baseline data, as reported in Section 6.3 of the main text. The values below the line show the indices for the ESP post-test data.

7.2 Extra information provided by the individual model over the CELEX-model

Second, we considered the amount of extra variation explained by the individual-model over the CELEX-model, for the MGL and the GCM, respectively.

Section 6.4 in the main text explains that we calculate the extra information given by the individual models by subtracting the CELEX model prediction from the individual model prediction for each item for each participant.

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-2.951	0.295	-10.007	0.000
rescaled_CELEX_gcm_feature_distance	4.369	0.425	10.284	0.000
gcm_score_diff	2.612	0.470	5.563	0.000
participant_prereg_average	5.289	0.602	8.788	0.000

Table 10: Model 5: Celex-GCM and individual score difference, post-test data

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-1.318	0.152	-8.663	0.000
rescaled_CELEX_mgl_feature_distance	1.944	0.176	11.045	0.000
mgl_score_diff	0.304	0.139	2.181	0.029
participant_prereg_average	7.203	0.591	12.190	0.000

Table 11: Model 6: Celex-MGL and individual score difference, post-test data

We do this for both the GCM (Model 5: regular response \sim CELEX GCM score + extra GCM information + participant pre-test mean + (1 + CELEX GCM + extra GCM information|participant) + (1|verb) and the MGL (Model 6: regular response \sim CELEX MGL score + extra MGL information + participant pre-test mean + (1 + CELEX MGL + extra MGL information|participant) + (1|verb)).

The extra information is a significant predictor of post-test responses in both models. However, this effect is much more robust for the GCM than for the MGL, as seen in Table 10 (M5) and Table 11 (M6). The individual-GCM has some success in modeling post-test behavior when trained with real English verbs and information from the ESP test, whereas adjustments made to the individual-MGL do not vastly improve predictive power.

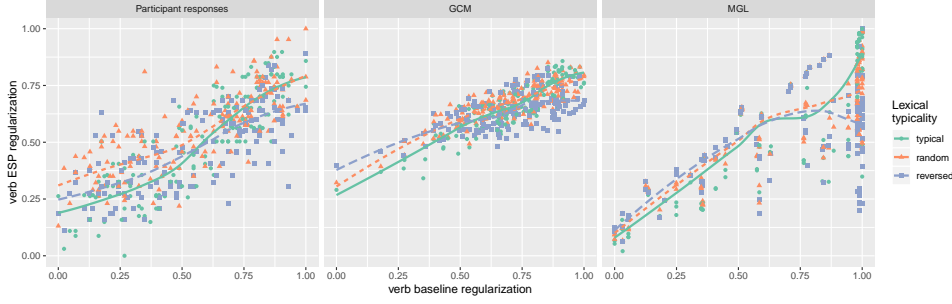


Figure 1: Baseline and individual regularization weights: participants, GCM, MGL

These techniques are indicative of the relative strengths of the individual-MGL and -GCM respective to their CELEX-counterparts and each other. However, they do not compare them directly.

7.3 Using combined model predictions to explain variation in participant post-test responses

We used the CELEX-GCM, CELEX-MGL, individual-GCM, and individual-MGL together to predict participant post-test responses. The best model, reported in the paper, contains the CELEX-MGL and the individual-GCM. The summary can be seen in Table 12.

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-2.942	0.270	-10.896	0.000
r_baseline_mgl_features	1.600	0.171	9.361	0.000
r_individual_gcm_features	2.786	0.378	7.369	0.000
participant_prereg_average	5.437	0.585	9.298	0.000

Table 12: Celex-MGL and Individual-GCM predict variation in the post-test

Goodness-of-fit tests show that model fit is significantly worse if we exclude CELEX-MGL ($\chi^2 = 201.86$) or individual-GCM ($\chi^2 = 81.82$) as predictors. It is not significantly better if we include individual-MGL ($\chi^2 = 0.09$) or CELEX-GCM ($\chi^2 = 2.03$).

7.4 Visualising model contributions to explaining variation in post-test responses

Figure 1 shows how baseline and post-test regularization means of words compare for participant responses, the GCM, and the MGL.

The Figure shows the average baseline regularity (x axis) and the post-test regularity (y axis) of the verbs in the ESP experiment. The post-test values vary across the dimension of co-player lexical typicality. Loess smooths show the general trajectory of the effects of typicality. In the first panel, regularity translates to aggregate participant responses in the baseline and in the ESP post-test. In the second and third panels, regularity translates to regularity scores provided by the CELEX- and individual-GCM and -MGL. Values are rescaled for better comparison.

While the loess smooths are likely slightly overfit, they show the general trend in the raw data, which is echoed in the regression analysis in the main text: the ‘typical’ trend

of baseline-post-test difference is stronger than either the "random" or "reversed" trends. These are not very much different from each other. The differences are gradual. The GCM, while generally overestimating regularity, captures this gradual shift, since verb behavior is shifted individually in the post-test. The MGL also overestimates regularization. Since rules affect groups of verbs, rather than individual verbs, the post-test shifts are abrupt. The "typical" trend continues climbing after a longer plateau. However, the 'random' and 'reversed' trends trail off, where the "reversed" trend turns negative towards the higher tail of the distribution.

References

- AIKEN, LEONA S.; STEPHEN G. WEST; and RAYMOND R. RENO. 1991. *Multiple regression: Testing and interpreting interactions*. Newbury Park, CA: Sage.
- ALBRIGHT, ADAM, and BRUCE HAYES. 2002. Modeling English past tense intuitions with minimal generalization. *Proceedings of the ACL-02 Workshop on Morphological and Phonological Learning*, vol. 6, 58–69. Stroudsburg, PA.
- ALBRIGHT, ADAM, and BRUCE HAYES. 2003. Rules vs. analogy in English past tenses: A computational/experimental study. *Cognition* 90.119–161.
- ALEGRE, MARIA, and PETER GORDON. 1999. Rule-based versus associative processes in derivational morphology. *Brain and Language* 68.347–354.
- BAAYEN, R. HARALD; DOUGLAS J. DAVIDSON; and DOUGLAS M. BATES. 2008. Mixed-effects modeling with crossed random effects for subjects and items. *Journal of Memory and Language* 59.390–412.
- BAAYEN, R. HARALD; LAURIE B. FELDMAN; and ROBERT SCHREUDER. 2006. Morphological influences on the recognition of monosyllabic monomorphemic words. *Journal of Memory and Language* 55.290–313.
- BAAYEN, R. HARALD; RICHARD PIEPENBROCK; and HEDDERIK VAN RIJN. 1993. *The CELEX lexical database on CD-ROM*. Philadelphia: Linguistic Data Consortium.
- BATES, DOUGLAS; MARTIN MÄCHLER; BEN BOLKER; and STEVE WALKER. 2015. Fitting linear mixed-effects models using lme4. *Journal of Statistical Software* 67.1–48, DOI: 10.18637/jss.v067.i01.
- BELSLEY, DAVID A. 1991. *Conditioning diagnostics: Collinearity and weak data in regression*. New York: Wiley.
- BYBEE, JOAN L., and CAROL LYNN MODER. 1983. Morphological classes as natural categories. *Language* 59.251–270.
- BYBEE, JOAN L., and DAN I. SLOBIN. 1982. Rules and schemas in the development and use of the English past tense. *Language* 58.265–289.
- DAWDY-HESTERBERG, LISA, and JANET B. PIERREHUMBERT. 2014. Learnability and generalisation of Arabic broken plural nouns. *Language, Cognition and Neuroscience* 29.1268–1282.
- ECHAMBADI, RAJ, and JAMES D. HESS. 2007. Mean-centering does not alleviate collinearity problems in moderated multiple regression models. *Marketing Science* 26.438–445.
- FRISCH, STEFAN A.; JANET B. PIERREHUMBERT; and MICHAEL BROE. 2004. Similarity avoidance and the OCP. *Natural Language and Linguistic Theory* 22.179–228.
- HARRELL JR, ET AL, FRANK E. 2008. Hmisc: A package of miscellaneous R functions. *R package version 3*. Online: <http://biostat.mc.vanderbilt.edu/Hmisc>.

- JACCARD, JAMES, and ROBERT TURRISI. 2003. *Interaction effects in multiple regression*. Thousand Oaks, CA: Sage.
- JAEGER, T. FLORIAN. 2008. Categorical data analysis: Away from ANOVAs (transformation or not) and towards logit mixed models. *Journal of Memory and Language* 59.434–446.
- MIELKE, JEFF. 2008. *The emergence of distinctive features*. Oxford: Oxford University Press.
- MODER, CAROL LYNN. 1992. *Productivity and categorization in morphological classes*. Buffalo, NY: State University of New York dissertation.
- NAKISA, RAMIN C.; KIM PLUNKETT; and ULRIKE HAHN. 2001. A cross-linguistic comparison of single and dual-route models of inflectional morphology. *Models of language acquisition: Inductive and deductive approaches*, ed. by Peter Broeder and Jaap Murre, 201–222. Cambridge, MA: MIT Press.
- NOSOFSKY, ROBERT M. 1990. Relations between exemplar-similarity and likelihood models of classification. *Journal of Mathematical Psychology* 34.393–418.
- SCHUMACHER, R. ALEXANDER, and JANET B. PIERREHUMBERT. 2017. Prior expectations in linguistic learning: A stochastic model of individual differences. *Proceedings of the 39th Meeting of the Cognitive Science Society (CogSci2017)*. London, UK.
- SCHUMACHER, R. ALEXANDER; JANET B. PIERREHUMBERT; and PATRICK LASHELL. 2014. Reconciling inconsistency in encoded morphological distinctions in an artificial language. *Proceedings of the 36th Annual Meeting of the Cognitive Science Society (CogSci2014)*, 2895–2900. Quebec City, CA.
- SINCLAIR, JOHN M. (ed.) 1987. *Looking up: An account of the COBUILD project in lexical computing*. London: Collins.
- WURM, LEE H., and SEBASTIANO A. FISICARO. 2014. What residualizing predictors in regression analyses does (and what it does not do). *Journal of Memory and Language* 72.37–48.