

# Concurrent hidden structure & grammar learning

Adeline Tan

University of California, Los Angeles

adelinertan@gmail.com

## Abstract

The concurrent learning of both unseen structures and grammar is an enduring problem in phonological acquisition. The present study develops a joint model of word-UR-SR triples that incorporates a Maximum Entropy model of SRs conditioned on URs. The learner was presented with word-SR frequencies, and successfully learned the hidden structures and grammars that enabled it to generalize well on test data that were withheld during training. When given an option between acquiring a grammar that supported a rich base analysis and one that didn't, the learner always acquired the grammar that supported rich bases. These results suggest that the preference for acquiring a rich base grammar over a non rich base one is an emergent property of the proposed model.

## 1 Introduction

In order to fully acquire language, a child has to acquire both the representations and grammar of her language from observed surface forms. Representations include underlying forms, metrical structures, morphological boundaries within words, *etc.* Such representations are absent from the observed data that the child receives, and are thus termed hidden structure. The current study focuses on the learning of hidden structure(s) concurrently with the grammar.

Multiple approaches have been proposed for the concurrent learning of hidden structure and an accompanying constraint-based grammar (Tesar and Smolensky, 2000; Jarosz, 2015; Boersma and Pater, 2016; Rasin and Katzir, 2016; Nelson, 2019). Following Eisenstat (2009), Pater et al. (2012), Staubs and Pater (2016), Nazarov and Pater (2017), and O'Hara (2017), this study incorporates a Maximum Entropy (MaxEnt) grammar (Goldwater and Johnson, 2003) that governs the mapping between hidden structures and surface forms. The current study

combines the word-hidden structure mapping with the hidden structure-surface form mapping by utilizing the chain rule of probability theory. This produces a joint word-underlying form-surface form (WORD-UR-SR) model that is compatible with a weighted-constraint grammar of UR-SR mappings. While the model is similar to the ones in Staubs and Pater (2016) and Nazarov and Pater (2017), the current study focuses on learning URs that an analyst would posit, with the learned grammar and lexicon subjected to generalization tasks with wug morphemes.

The model and the learner are introduced in §2 and §3 respectively. We then turn to several schematic languages, the first of which is based on English voicing assimilation (§4). This is followed by a set of six stress languages (§5). Two of the languages within the stress set allow multiple analyses, of which only one analysis supports rich bases (Prince and Smolensky, 2004), thus providing the opportunity to determine whether there is a preference for acquiring the rich base grammar over a non rich base one. The final schematic language is based on English velar softening (§6). §7 concludes.

## 2 Model

The knowledge whose acquisition will be investigated is knowledge of a particular distribution over WORD-UR-SR triples (*e.g.* <CROC-PL, /kɪak+z/, [kɪaks]>: 99%; <CROC-PL, /kɪak+z/, [kɪakz]>: .003%; <CROC-PL, /kɪak+s/, [kɪaks]>: .002%; ...). In this paper, WORD<sup>1</sup> represents a sequence of morphemes, and morphemes are represented with uppercase letters.

<sup>1</sup>WORD is also abbreviated WD in this paper.

The probability of a triple can be rewritten as:

$$Pr(WD, UR, SR) = Pr(SR|WD, UR) * Pr(WD, UR) \quad (1)$$

The first term,  $Pr(SR|WD, UR)$ , is the probability of an SR for a given WORD-UR pair, and is determined by the traditional phonological constraint grammar. For instance, if  $Pr([bæŋks]|BANK-PL, /bæŋk+z/) = 0.9$ , then we should interpret it to mean that the WORD-UR pair  $\langle BANK-PL, /bæŋk+z/ \rangle$  is realized as SR  $[bæŋks]$  90% of the time. The model proposed here does not condition the UR-SR mapping on the word. Using the example above, this means that  $Pr([bæŋks]|BANK_1-PL, /bæŋk+z/) = Pr([bæŋks]|BANK_2-PL, /bæŋk+z/) = Pr([bæŋks]|BANK_3-3SG.PRES, /bæŋk+z/)$ , where  $BANK_1$  is the financial institution concept,  $BANK_2$  is the river concept, and  $BANK_3$  is the concept of turning at an angle. Consequently,  $Pr(SR|WD, UR) = Pr(SR|UR)$ , and the probability of the WORD-UR-SR triple can be simplified to equation (2):

$$Pr(WD, UR, SR) = Pr(SR|UR) * Pr(WD, UR) \quad (2)$$

Such probabilistic mappings of SRs conditioned on URs (*i.e.*  $Pr(SR|UR)$ ) are computed by virtually all probabilistic constraint-based grammars (*e.g.* probabilistic OT, probabilistic versions of Harmonic Grammar, *etc.*) The current study uses a MaxEnt model, which is a weighted constraint grammar.

Following the traditional phonological MaxEnt model, each UR-SR pair  $(x, y)$  is associated with a feature vector,  $\vec{v}(x, y)$ , which captures the pair's properties. For UR-SR pairs, there are two classes of relevant properties. The first class concerns the form that the SR takes. For example, a feature may be used to track how many pairs of adjacent obstruents of an SR have different voicing values. Such features are known as markedness constraints. The second class of features concerns the mapping between the UR and the SR, and are most commonly used to penalize any changes between the two. These features are conventionally known as faithfulness constraints. Each feature has an associated weight, and the feature weights can be organized into the weight vector  $\vec{w}$ . The features

of the UR-SR pair  $(x, y)$  are linearly combined (as in equation (3)<sup>2</sup>) to produce its harmony score,  $h(x, y)$ .  $h(x, y)$  is essentially the weighted sum of the UR-SR pair  $(x, y)$ 's features, and is a scalar (rather than a vector).

$$h(x, y) = -(\vec{w} \cdot \vec{v}(x, y)) \quad (3)$$

The MaxEnt model then maps each pair's harmony score to its probability (equation (4)).

$$Pr(SR = y|UR = x) = \frac{e^{h(x, y)}}{Z(x)} \quad (4)$$

Since the traditional phonological MaxEnt grammar is a conditional ("discriminative") model, the partition function  $Z(x)$  sums over all UR-SR pairs that share the same UR (equation 5).

$$Z(x) = \sum_{y' \in \mathcal{Y}_x} e^{h(x, y')} \quad (5)$$

In equation (5),  $\mathcal{Y}_x$  is the set of all SRs that are compatible with UR  $x$ . This has the effect of normalizing the probability of a particular UR-SR mapping among only all other mappings from the same UR.

The second term in equation (2),  $Pr(WD, UR)$ , is the joint probability of a WORD-UR pair. This implicitly defines a conditional distribution  $Pr(UR|WD)$  (equation (6)).

$$Pr(UR = x|WD = w) = \frac{Pr(WD = w, UR = x)}{\sum_{x'} Pr(WD = w, UR = x')} \quad (6)$$

Under this conditional distribution we would expect  $Pr(/kɹɔk+z/|CROC-PL)$  to be high, and  $Pr(/kɹɔk+s/|CROC-PL)$ ,  $Pr(/kɹɔg+z/|CROC-PL)$ , *etc.* to be low. For the morpheme CROC, the learner needs to choose between 2 possible stem-final segments: voiceless  $/k/$  and voiced  $/g/$ <sup>3</sup>. For the plural morpheme, the learner needs to choose between voiceless  $/s/$  and voiced  $/z/$ . Consequently, there are four potential URs that the learner considers for the word CROC-PL (Table 1). Table 1 also shows the four features for each

<sup>2</sup>A UR-SR pair is active for a phonological constraint when it violates the requirements of that constraint, which in turn reduces the pair's conditional probability. Hence the negative sign in equation (3).

<sup>3</sup>This example is modeled after English voicing assimilation where adjacent obstruents agree in voicing. The surface sequence  $[ks]$  could have arisen from any of the following UR sequences  $\{/k+s/, /k+z/, /g+s/, /g+z/\}$ . Hence, I vary only the stem-final segment, but none of the other stem segments.

of the four variants that the learner has to choose among. These features represent the strength of association between a particular morpheme and an aspect (*e.g.* morpheme-final obstruent voicing) of its UR. Within phonology, such features are also known as UR constraints (Zuraw, 2000; Boersma, 2001). Similar to its UR-SR counterpart, there is a feature vector  $\vec{u}(w, x)$  for each WORD-UR pair  $(w, x)$ . Likewise, the UR constraint weights can be organized into a vector  $\vec{\theta}$ . The harmony score for each WORD-UR pair is computed as per equation (7)<sup>4</sup>.

$$g(x, y) = \vec{\theta} \cdot \vec{u}(x, y) \quad (7)$$

The harmony score of a WORD-UR pair is then mapped to its probability (equation 8).

$$Pr(WD = w, UR = x) = \frac{e^{g(w, x)}}{Z} \quad (8)$$

In contrast to the UR-SR model described above, the WORD-UR model is not conditional. The normalization takes place over all WORD-UR pairs (equation (9)).

$$Z = \sum_{w \in \mathcal{W}} \sum_{x \in \mathcal{X}_w} e^{g(w, x)} \quad (9)$$

In equation (9),  $\mathcal{W}$  is the set of words, and  $\mathcal{X}_w$  is the set of all URs that are compatible with word  $w$ . This normalization produces a generative distribution over WORD-UR pairs, which in turn produces the generative distribution over WORD-UR-SR triples of equation (2). This departs from the models in Staubs and Pater (2016) and Nazarov and Pater (2017), which are discriminative models. A generative model is capable of describing differences in the frequencies of various words, in addition to the relationship between words and their realizations, whereas a discriminative model only does the latter.

### 3 Learning

The model takes in a set of WORD-SR pair frequencies (*e.g.*  $\langle \text{CROC-PL}, [\text{kɪ} \text{ɔ} \text{k} \text{s}] \rangle: 50$ ;  $\langle \text{CROC-PL}, [\text{kɪ} \text{ɔ} \text{k} \text{z}] \rangle: 0$ ;  $\dots$ ), and learns a probability distribution over WORD-UR-SR triples (*e.g.*  $\langle \text{CROC-PL}, / \text{kɪ} \text{ɔ} \text{k} + \text{z} /, [\text{kɪ} \text{ɔ} \text{k} \text{s}] \rangle: 99\%$ ;  $\langle \text{CROC-PL}, / \text{kɪ} \text{ɔ} \text{k} + \text{z} /, [\text{kɪ} \text{ɔ} \text{k} \text{z}] \rangle: .003\%$ ;  $\langle \text{CROC-PL},$

<sup>4</sup>A WORD-UR pair is active for a particular UR constraint when it contains the morpheme, segment, *etc.*, required by that constraint, which in turn increases the pair’s probability. Hence the sign difference between equations (3) and (7).

$/ \text{kɪ} \text{ɔ} \text{k} + \text{s} /, [\text{kɪ} \text{ɔ} \text{k} \text{s}] \rangle: .002\%$ ;  $\dots$ ). The triple probability defined in Section 2 in fact implicitly defines a distribution over WORD-SR pairs as well. More concretely, the probability of pairs can be computed from the probability of triples via this summation:

$$\begin{aligned} Pr(WD = w, SR = y) &= \sum_x Pr(WD = w, UR = x, SR = y) \\ &= \sum_x Pr(SR = y | UR = x) * Pr(WD = w, UR = x) \\ &= \sum_x \frac{e^{h(x, y)}}{Z(x)} * \frac{e^{g(w, x)}}{Z} \end{aligned} \quad (10)$$

The likelihood  $Pr(WD, SR)$  can be understood as a function of the parameters  $\vec{w}$  and  $\vec{\theta}$ . Experimentation showed that regularization terms did not improve performance in fitting to test data that was withheld from training, so the learner’s objective is to seek the values of  $\vec{w}$  and  $\vec{\theta}$  that maximize this likelihood. In order to assess which values of  $\vec{w}$  and  $\vec{\theta}$  will be found by the learner, I use the Expectation-Maximization (EM) algorithm (Dempster et al., 1977). Notice that  $Pr(WD, SR)$  is a marginal distribution<sup>5</sup>. The likelihood function of marginal distributions is not guaranteed to be convex, so each EM run finds a local maximum. I take the highest of these local maxima to identify the predicted outcome of learning.

## 4 English Voicing Assimilation

In English voicing assimilation, adjacent obstruents with different voicing values are resolved with suffixes assimilating their voicing value to that of the stem. The underlying voicing value of stem-final obstruents and suffixes constitute the hidden structures.

### 4.1 Experimental setup

The first language had the words {CROC-PL, DOG-PL, COW-PL}. Its grammar had the constraints {AGREE(voice), IDENT<sub>stem</sub>, IDENT<sub>general</sub>}, so  $\vec{w}$  was 3-dimensional for this language. In addition, six potential UR variants { $/ \text{kɪ} \text{ɔ} \text{k} /, / \text{kɪ} \text{ɔ} \text{ɔ} /, / \text{d} \text{ɔ} \text{k} /, / \text{d} \text{ɔ} \text{ɔ} /, / - \text{s} /, / - \text{z} /$ } were considered, making  $\vec{\theta}$  6-dimensional. These nine dimensions correspond to the first nine rows in Table 3.

What constitutes successful learning? First, we can check whether the UR learned for each morpheme of the training data matches what would

<sup>5</sup>The summation in equation (10) produces marginal probabilities.

WORD	UR <sub>WORD</sub>	(CROC, /kɪɑk/)	(CROC, /kɪɑg/)	(PL, /-s/)	(PL, /-z/)
CROC-PL	/kɪɑk+s/	1	0	1	0
	/kɪɑg+s/	0	1	1	0
	/kɪɑk+z/	1	0	0	1
	/kɪɑg+z/	0	1	0	1

Table 1: UR constraints for the word CROC-PL.

be predicted via traditional phonological analyses. For example, a phonologist would posit that a child learns the URs /kɪɑk/, /dɑg/ and /-z/ for the CROC, DOG and -PL morphemes respectively. Recall that the model produces a distribution over WORD-UR-SR triples. In order to find the probability of the word-sized UR containing /-z/ in the appropriate position, given that the WORD (*i.e.* sequence of morphemes) has -PL in that position, we apply the following equation:

$$p(/-z/|-PL) = \frac{[Pr(UR=/dɑk-z/, WD=DOG-PL) + Pr(UR=/dɑg-z/, WD=DOG-PL) + Pr(UR=/kɪɑk-z/, WD=CROC-PL) + Pr(UR=/kɪɑg-z/, WD=CROC-PL) + Pr(UR=/kɑw-z/, WD=COW-PL)]}{[Pr(WD=DOG-PL) + Pr(WD=CROC-PL) + Pr(WD=COW-PL)]} \quad (11)$$

The probability of a particular UR for each morpheme is calculated in the same manner for each of the other morpheme-UR pairs.

Second, the model must be able to generalize to unseen data in the way that humans do. Unseen data are WORD-SR pairs that the model wasn't provided with in the training set. For English voicing assimilation, the words in Table 2 provide a good test set for generalizability. Each test word is com-

WORD
WUG-PL
HEAK-PL
CRA-PL
DOG-D
CROC-D

Table 2: Test set words for English voicing assimilation.

posed of a new morpheme and an old morpheme. This isolates each of the old morphemes, so that only one old morpheme appears in each word. This allows us to test what the model knows between /-s/ vs. /-z/ as the UR of the -PL morpheme, as well as what the model knows about how /-s/ & /-z/ are realized on the surface. The new nouns WUG /wɑg/, HEAK /hɪk/ & CRA /kɪɑ/ will illuminate what the model had learned about the -PL suffix. I

introduce a novel suffix with UR /-d/ (representing some morpheme I will write as -D), to test what the model had learned about the roots DOG and CROC. For example, if  $Pr((/hɪk-z/, [hɪks])|HEAK-PL)$  is high, then we'd know that the model generalizes in the same way that English speakers do, as evidenced by wug tests. The probability of a UR-SR pair for a given word is calculated as per equation (12).

$$Pr(UR = x, SR = y | WD = w) = \frac{Pr(WD = w, UR = x, SR = y)}{\sum_{x'} \sum_{y'} Pr(WD = w, UR = x', SR = y')} \quad (12)$$

## 4.2 Results

The training data consisted of 10 logically possible WORD-SR pairs, of which only three were observed. Each of the three observed pairs {(CROC-PL, [kɪɑks]), (DOG-PL, [dɑgz]), (COW-PL, [kɑwz])} was only observed once. The learner sought the parameter values that maximized the likelihood of the training data. Five settings of the parameters are shown in Table 3. I found these by running the EM algorithm from 20 randomly initialized<sup>6</sup> starting points. The likelihood of the training data for each of the five parameter settings is  $0.33 \times 0.33 \times 0.33 = 0.33^3$ . These five settings have already hit the maximum likelihood of training data – there isn't another parameter setting that would provide a much better likelihood since these settings have already matched the empirical relative frequencies almost perfectly.

<sup>6</sup>Initial weights for all eight languages in the present study were drawn from a uniform distribution with range=[0.1, 5) for phonological constraints & range=[0, 100) for UR constraints.

<sup>7</sup>In this paper, negative weights were only allowed for UR constraints. Weights for regular phonological constraints were not allowed to be negative. For voicing assimilation, the UR constraint (COW, /kɑw/) was excluded from the set of features, since the morpheme COW had only one underlying form under consideration. This resulted in (DOG, /dɑg/) attaining 0 weight, which pushed (DOG, /dɑk/) to a negative weight. Because it is the difference between weights rather than the actual value of the weights that matter, the negative weights do not have any meaningful impact on the results.

	1	2	3	4	5
AGREE(voice)	24.7	40.9	24.4	33.3	26.5
IDENT <sub>stem</sub>	15.0	11.9	12.6	29.4	11.6
IDENT <sub>general</sub>	11.9	18.5	12.0	18.5	13.6
(DOG, /dak/)	-43.2	-11.9	-30.0	-26.4	-11.6
(DOG, /dag/)	0.0	0.0	0.0	0.0	0.0
(CROC, /kɪak/)	0.0	0.0	0.0	0.0	0.0
(CROC, /kɪag/)	-16.3	-23.7	-20.5	-33.1	-13.3
(-PL, /-s/)	0.3	14.4	9.5	20.7	24.6
(-PL, /-z/)	91.0	56.8	76.6	83.0	77.4
$Pr(\text{CROC-PL, [kɪaks]})$	0.33	0.33	0.33	0.33	0.33
$Pr(\text{DOG-PL, [dagz]})$	0.33	0.33	0.33	0.33	0.33
$Pr(\text{COW-PL, [kawz]})$	0.33	0.33	0.33	0.33	0.33
<i>Likelihood of training data</i>	0.33 <sup>3</sup>	0.33 <sup>3</sup>	0.33 <sup>3</sup>	0.33 <sup>3</sup>	0.33 <sup>3</sup>
<i>Negative log-likelihood of training data</i>	-3.29585	-3.29585	-3.29585	-3.29584	-3.29586

Table 3: Feature weights<sup>7</sup>, probability of observed data, & likelihood of training data from the best five runs (English voicing assimilation).

Recall the first criterion of successful learning: the learner has to learn the very same morpheme-sized URs that human learners are posited by phonologists to learn. Applying the equation in (11), we see that the lexicon learned is indeed the one that matches with traditional phonological analysis (Table 4<sup>8</sup>). A further examination of the UR constraints confirms that the UR constraints associated with expected URs (*i.e.* DOG has underlying root-final /g/, CROC has underlying root-final /k/, the plural morpheme is underlying /-z/) have higher weights than their counterparts (Table 3).

	1	2	3	4	5
$p(/dag/ DOG)$	1.0	1.0	1.0	1.0	1.0
$p(/kɪak/ CROC)$	1.0	1.0	1.0	1.0	1.0
$p(/-z/ -PL)$	1.0	1.0	1.0	1.0	1.0

Table 4: Lexicon (English voicing assimilation).

To fulfill the second criterion of successful learning, the learned models had to generalize in the same way that English speakers generalize. The generalization task considered 8 UR-SR combinations for each word<sup>9</sup>. For nonce word HEAK /hik/, an English speaker would produce [hiks] for HEAK-PL via the UR /hik-z/. Thus, successful generalization for the word HEAK-PL required assigning high probability to the UR-SR

<sup>8</sup>All probabilities in Table 4 were very close to or at 100%. The lowest was 99.9991% for  $p(/dag/|DOG)$  of the fifth parameter setting.

<sup>9</sup>To illustrate the 8 combinations, consider the nonce word WUGS. 2 variations are available via the UR: /-s, -z/, 2 via the SR of the stem-final consonant: [wɪk, wɪg], and 2 via the SR of the suffix consonant: [-s, -z]. This produces  $2^3 = 8$  UR-SR combinations.

pair (/hik-z/, [hiks]), and low probability to the seven other pairs. Table 5 presents, in the top five rows,  $Pr((UR, SR)|WD)$  for the UR-SR pairs that are expected to have high probability for their respective words, given what we know about how English speakers behave on wug tests. The results in Table 5<sup>10</sup> indicate that all five models very successfully generalized in a manner that mimicked speakers, with probabilities close to or at 100%. A look at the learned phonological constraint weights in Table 3 shows why all five parameter settings mirrored speakers so well in the generalization task. The models all learned the two crucial weight-inequalities required for English voicing assimilation:  $AGREE(voice) > IDENT_{general}$  as well as  $IDENT_{stem} > 0$ .

For voicing assimilation, the model did well on both the UR-learning of morphemes & wug-test mirroring tasks because both the lexicon & grammar learned by the learner mirrored what English speakers are believed to have learned for voicing assimilation.

## 5 PAKA stress languages

### 5.1 Experimental setup

The next 6 languages were generated using similar morphemes and constraints that generated the PAKA World dataset in Tesar et al. (2003). There were two roots: (PA & BA), as well as two suffixes (-KA & -GA). The URs of PA and KA were always unstressed: /pa/ & /-ka/. In contrast, the

<sup>10</sup>All probabilities in Table 5 were very close to or at 100%. The lowest was 99.9986% for  $Pr((/dag-d/, [dagd])|DOG-D)$  of the fifth parameter setting.

	1	2	3	4	5
$Pr((/wAg-z/, [wAgz]) WUG-PL)$	1.0	1.0	1.0	1.0	1.0
$Pr((/hik-z/, [hiks]) HEAK-PL)$	1.0	1.0	1.0	1.0	1.0
$Pr((/kɪɑ-z/, [kɪɑz]) CRA-PL)$	1.0	1.0	1.0	1.0	1.0
$Pr((/dɑg-d/, [dɑgd]) DOG-D)$	1.0	1.0	1.0	1.0	1.0
$Pr((/kɪɑk-d/, [kɪɑkt]) CROC-D)$	1.0	1.0	1.0	1.0	1.0
$Pr((/wAg-s/, [wAgz]) WUG-PL)$	0.0	0.0	0.0	0.0	0.0
$Pr((/wAg-z/, [wAkz]) WUG-PL)$	0.0	0.0	0.0	0.0	0.0
⋮	⋮	⋮	⋮	⋮	⋮

Table 5: Probability of UR-SR pair for a given test set word (English voicing assimilation).

URs of BA & GA could bear stress, so the model considered the potential-URs  $\{/ba/, /ba/, /-'ga/, /-ga/\}$ . Accordingly, the relevant UR constraints for this dataset were: (BA, /ba/), (BA, /ba/), (GA, /-'ga/) & (GA, /-ga/). The four features that went into  $Pr(SR|UR)$  were:

- MAINLEFT (*ML*)
  - Stress the leftmost syllable.
- MAINRIGHT (*MR*)
  - Stress the rightmost syllable.
- MAX<sub>general</sub>-STRESS (*F*)
  - If a syllable is stressed in the UR, retain its stress in the SR.
- MAX<sub>root</sub>-STRESS (*FR*)
  - If a root syllable is stressed in the UR, retain its stress in the SR.

These morphemes & constraints produced six logically possible languages<sup>11</sup>, which are the six sets of observed SRs shown in Table 6. Languages 3, 4 & 6 were each compatible with only 1 lexicon. Language 5 was compatible with both /-ga/ & /-'ga/, but compatible with only a single grammar. For Languages 1 & 2, four lexicon-grammar combinations were available for each language.

## 5.2 Results

For each language, the training data had 12 logically possible WORD-SR pairs<sup>13</sup>, of which four pairs were each observed once. The four SRs for

<sup>11</sup>Since MaxEnt generates probabilistic languages, there are technically an infinite number of possible languages. However, I’m restricting the set of languages to only those where there is effectively only one winning SR per UR.

<sup>12</sup>Since the model is MaxEnt rather than non-probabilistic Harmonic Grammar, the difference between the terms on both sides of an inequality need to be sufficiently large in order to generate categorical outcomes. Determining exactly how large a difference is needed for each inequality is difficult. Nevertheless, the test task provides a way to check that the trained weights indeed produce sufficient difference between the two terms of an inequality. If the difference were not sufficiently large, the test task would fail to produce categorical outcomes.

<sup>13</sup>There were three SRs per word – left-stress, right-stress, and no stress at all.

each language can be read off column “Observed SRs” of Table 6. As with English voicing assimilation, I did 20 EM runs per language. For all six languages, the learner succeeded in finding multiple parameter settings that hit the maximum likelihood of training data  $0.25 \times 0.25 \times 0.25 \times 0.25 = 0.25^4$ . For the sake of brevity, only one of these parameter settings is presented for each of the six languages (Table 7).

To test generalizability, two new roots  $\{SO /so/, ZE /ze/\}$  and two new suffixes  $\{-FO /-fo/, -VE /-ve/\}$  were introduced to form test set words (Table 8). As expected, all parameter settings that attained the maximum likelihood of training data generalized to test words at near 100% probability. A sample of the probabilities of UR-SR pairs for test word BA-FO is shown for one simulation of Language 3, where the combination of trained morpheme BA with an unaccented suffix like /-ka/ produced stress on the first syllable. Likewise, when BA combines with unstressed /-fo/, successful generalization requires a UR-SR pair with [bafo] to have high probability (Table 9).

## 5.3 Rich base supporting grammars

According to Prince and Smolensky (2004), the role of a constraint-based grammar is to assign an output to each input<sup>15</sup>. In the case of absolute ill-formedness (*e.g.* absence of right-stressed SRs in a left-stressed language), the grammar (*i.e.* the constraint interactions that govern the UR-SR mapping) must ensure that no input ever leads to ill-formed outputs (*e.g.* not even a UR with rightmost

<sup>14</sup>UR constraints for the morphemes PA & -KA that do not have multiple URs under consideration were included in this feature set. Hence negative UR constraint weights do not make an appearance here.

<sup>15</sup>Prince and Smolensky (2004) were writing about Optimality Theory, where the grammar consisted of ranked constraints picking a sole output for each input. Nevertheless, the grammar’s role in mapping inputs to outputs still holds for probabilistic constraint-based grammars.

<i>Lg</i>	<i>Observed SRs</i>	<i>Description</i>	<i>Lexicon</i>	<i>Required weight inequalities</i> <sup>12</sup>
1	[ˈpaka, ˈpaga, ˈbaka, ˈbaga]	predictable left-stress	/ba, -ga/ /ˈba, -ga/	ML > MR
			/ba, -ˈga/ /ˈba, -ˈga/	ML > MR + F
2	[paˈka, paˈga, baˈka, baˈga]	predictable right-stress	/ba, -ga/ /ba, -ga/	MR > ML
			/ˈba, -ga/ /ˈba, -ˈga/	MR > ML + F + FR
3	[ˈpaka, paˈga, ˈbaka, ˈbaga]	full accentual contrast, default left	/ˈba, -ˈga/	MR + F > ML > MR
4	[paˈka, paˈga, ˈbaka, baˈga]	full accentual contrast, default right	/ˈba, -ˈga/	ML + F + FR > MR > ML + FR
5	[paˈka, paˈga, ˈbaka, ˈbaga]	contrast in roots only, default right	/ˈba, -ga/ /ˈba, -ˈga/	ML + FR > MR > ML
6	[ˈpaka, paˈga, ˈbaka, baˈga]	contrast in suffixes only, default left	/ba, -ˈga/	MR + F > ML > MR

Table 6: PAKA languages and respective logically-possible lexicon-grammar combinations.

	<i>Lg 1</i>	<i>Lg 2</i>	<i>Lg 3</i>	<i>Lg 4</i>	<i>Lg 5</i>	<i>Lg 6</i>
MAINLEFT (ML)	19.6	0.0	16.9	8.4	0.0	44.7
MAINRIGHT (MR)	0.0	20.0	0.0	23.3	19.4	0.0
MAX <sub>general</sub> -STRESS (F)	1.5	1.7	32.4	37.4	21.7	111.4
MAX <sub>root</sub> -STRESS (FR)	3.1	0.0	7.5	0.0	23.2	2.2
(BA, /ba/)	69.9	47.3	15.5	36.4	27.8	81.1
(BA, /ˈba/)	0.5	77.2	61.4	80.5	66.7	34.2
(GA, /-ga/)	59.6	28.1	11.1	25.8	71.4	59.3
(GA, /-ˈga/)	6.2	65.2	61.2	85.0	35.2	76.8
<i>Likelihood</i> <small>training data</small>	0.25 <sup>4</sup>	0.25 <sup>4</sup>	0.25 <sup>4</sup>	0.25 <sup>4</sup>	0.25 <sup>4</sup>	0.25 <sup>4</sup>
<i>Negative log-likelihood</i> <small>training data</small>	-5.545177	-5.545177	-5.545178	-5.545178	-5.545177	-5.545177

Table 7: Feature weights<sup>14</sup> & likelihood of training data from the best runs for each PAKA language.

WORD	WORD	UR-SR pair	<i>Pr</i> (UR, SR WD)
SO -FO	BA -FO	/ba-fo/, [bafo]	$2.7 \times 10^{-12}$
SO -GA		/ba-fo/, [ˈbafo]	$4.6 \times 10^{-5}$
BA -FO		/ba-fo/, [baˈfo]	$2.7 \times 10^{-11}$
ZE -GA		/ˈba-fo/, [bafo]	$1.8 \times 10^{-21}$
BA -VE		/ˈba-fo/, [ˈbafo]	$9.9995 \times 10$
		/ˈba-fo/, [baˈfo]	$1.7 \times 10^{-20}$

Table 8: Test set words for the six PAKA languages.

Table 9: Generalization to BA-FO in Language 3 (one run shown).

stress can produce an SR with rightmost stress). Within models that feature probabilistic UR-SR mappings, this translates to the grammar ensuring that no inputs ever map to ill-formed outputs with anything other than a vanishingly small probability. In other words, the grammar should be fail-safe; it should be able to map all URs (even implausible ones like a right-stressed UR in a left-stressed language) to SRs with appropriate probability values. This concept is known as the Richness of the Base (Prince and Smolensky, 2004).

Language 1 & Language 2 are languages with predictable left- and right-stress respectively. Each of these two languages is compatible with two

grammars (Table 6). In Language 1, the two possible grammars are Grammar 1 ( $ML > MR$ ) & Grammar 2 ( $ML > MR + F$ ). The lexicon that includes /-ga/ minimally requires Grammar 1, while the lexicon that includes /-ˈga/ minimally requires Grammar 2. Since the weights of phonological constraints could not be negative, Grammar 2 entails Grammar 1. It follows that Grammar 2 is compatible with both /-ga/ & /-ˈga/ while Grammar 1 is compatible with only /-ga/. Grammar 2 is thus a grammar that supports rich bases because it is capable of producing SRs with the right proba-

bilities even with an implausible UR (underlying stressed suffix /-'ga/ in a left-stressed language where unstressed root /pa/ also exists). In contrast, Grammar 1 is the less restrictive grammar because it requires only  $ML > MR$ , thus allowing  $MR + F > ML$ . The gang effect of right-stress ( $MR$ ) and general faithfulness ( $F$ ) over left-stress ( $ML$ ) results in /pa-'ga/ surfacing with rightmost stress \*[pa-'ga]. For Language 2, this entailment relation also holds amongst its two grammars, with Grammar 4 ( $MR > ML + F + FR$ ) being the rich base supporting grammar, and Grammar 3 ( $MR > ML$ ) the less restrictive one. These two languages thus provide useful test cases on whether there is a preference for a rich base supporting grammar over its less restrictive rival or vice versa.

All 20 runs for Languages 1 & 2 always learned the rich base grammar. Trained weights for an example run of Language 1 is shown in Table 7, where the rich base grammar,  $ML (19.6) > MR + F (0 + 1.5)$ , is learned. To test this further, I ran 200 more simulations for both languages. All 200 runs for both languages always learned the rich base supporting grammar, sometimes with the lexicon that minimally required the rich base grammar, and sometimes with the lexicon that minimally required the less restrictive one. This indicated a strong preference for learning the rich base grammar over its less restrictive counterpart.

The preference for acquiring the rich base grammar is an emergent property of the model. EM finds the local maximum by hill-climbing from a randomly initialized point within the solution space. The solution space is the likelihood function of the marginal distribution (equation (10)) of the model defined in §2. Hill-climbing (*i.e.* gradient ascent) is guided by the gradients of the solution space at the current point. The preference for converging at maxima corresponding to the rich base grammar indicates the following: within the solution space, there are more points with gradients pointing towards maxima corresponding to the rich base grammar and fewer points with gradients pointing towards maxima corresponding to the less restrictive grammar. Since the solution space is a property of the model (rather than that of a particular learner), models with similar architecture (*e.g.* [Staubs and Pater \(2016\)](#); [Nazarov and Pater \(2017\)](#)) are likely to also favor the acquisition of rich base grammars.

## 6 Velar Softening

In English velar softening, /k/ → [s] before a high front vowel when a morpheme-boundary intervenes (*e.g.* *electri*[k]~*electri*[s]-*ity*). Velar softening is an instance of the derived environment effect (DEE) because its triggering environment requires the presence of a morpheme boundary. DEEs are a puzzle because both the alternation and the segmentation into morphemes must be acquired simultaneously. In an additional wrinkle, DEEs are often only triggered by specific morphemes. For example, the *-ity* morpheme triggers velar softening, but the *-ish* morpheme does not. Velar softening thus has three sources of hidden structure – presence of a morpheme boundary, whether a particular suffix is exceptional in triggering velar softening, and the usual UR-segment-learning (/k/ or /s/) that we’ve already seen in the preceding test cases. I use the \*-symbol to indicate the exception tagged UR variant.

### 6.1 Experimental setup

There were eight observed words {ELECTRIC, ELECTRICITY, ELECTRICISH, KITTY, SECURE, SECURITY, SMALL, SMALLISH}. The three sources of hidden structure were combined to produce URs for these eight words. The URs for ELECTRICITY are shown with their relevant UR constraint features (Table 10). The UR /*electrik*-*\*ity*/ contained the morphemes ELECTRIC<sup>16</sup> & -ITY, so it was active for those two features. These features represent a new class of UR constraint, namely those that indicate the presence of certain morphemes. Such features were not required in the preceding cases as the morpheme boundary was not in question. The UR /*electrik*-*\*ity*/ had underlying /k/ for morpheme ELECTRIC, and the exception-tagged version of the -ITY suffix, so it was also active for features (ELECTRIC, k)<sup>17</sup> & (-ITY, -\*ity) respectively. These UR constraints are the same kind that we’ve seen before.

Three phonological constraints controlled the UR-SR mapping – a general markedness constraint against [ki] sequences ( $M$ ), an exception-tagged version that was active only when an exception-tagged morpheme was part of the [ki] sequence ( $ME$ ), and a general IDENT constraint ( $F_{gen}$ ).

<sup>16</sup>Abbreviated as EL...C in Table 10.

<sup>17</sup>Abbreviated as (EL...C, k) in Table 10.



UR <sub>ELECTRICITY</sub>	EL...C	EL...CITY	-ITY	(EL...C, k)	(EL...C, s)	(EL...CITY, k)	(EL...CITY, s)	(-ITY, -*ity)	(-ITY, -ity)
/elektrik-*ity/	1	0	1	1	0	0	0	1	0
/elektrik-ity/	1	0	1	1	0	0	0	0	1
/electris-*ity/	1	0	1	0	1	0	0	1	0
/electris-ity/	1	0	1	0	1	0	0	0	1
/elektrikity/	0	1	0	0	0	1	0	0	0
/electricity/	0	1	0	0	0	0	1	0	0

Table 10: URs under consideration for the word ELECTRICITY shown with their UR constraint features (English velar softening).

## 6.2 Results

The training data had 12 logically possible WORD-SR pairs, of which eight were each observed once. Of 125 EM runs, multiple parameter settings were found to have reached the maximum likelihood of training data  $0.125^8 = 5.9605 \times 10^{-5}$ . 90.9% of these parameter settings learned the very same hidden structures that matched the standard phonological analysis of velar softening. This included learning that the -ITY morpheme was exception-tagged but that -ISH wasn't. All of these hidden features were learned at probabilities<sup>18</sup> greater than 97%, with the lowest going to the probability of a morpheme boundary in SECURITY at 97.3%.

To test whether the learned grammars generalized in a way that mimicked human learners, the following three morphemes were introduced: a new root with a morpheme-final-/k/ (CLEMIC /klɛmɪk/), an exception-tagged suffix \*-ISM /\*-izm/ and a non-exception-tagged suffix -Y /-i/ to create the test words in Table 11. For each

WORD	Expected UR-SR
ELECTRIC-ISM	electri/k/*ism, electri[s]ism
ELECTRIC-Y	electri/k/-y, electri[k]y
CLEMIC-ITY	clɛmi/k/*ity, clɛmi[s]ity
CLEMIC-ISH	clɛmi/k/-ish, clɛmi[k]ish

Table 11: Test set words for English velar softening & UR-SR pairs expected to be learned by human learners.

test word, its expected UR-SR pair arose from what traditional phonological analysis would have a child positing as its UR and SR. For instance, the word ELECTRIC-ISM would be posited to have morpheme-final /k/ as opposed to /s/ for ELECTRIC, with that /k/ surfacing as [s]. The expected UR-SR pair for each word is shown in column “Expected UR-SR” of Table 11. The same 90.9% of parameter settings that hit the maximum likelihood of training data generalized well to the test set with probabilities of the expected UR-SR pair for each

<sup>18</sup>These probability values were calculated using the same method shown in eq (11).

word approaching 100%<sup>19</sup>. Examination of the weights learned for the parameter settings that successfully generalized confirmed that they had each learned the grammar necessary for velar softening. That is, the alternation applied when the suffix was exception-tagged ( $ME + M > F_{gen}$ ), but did not take place when the suffix wasn't exception-tagged ( $F_{gen} > M$ ).

What would this high-but-not-100% rate of acquisition of the velar softening grammar mean for human learners? Perhaps 10% of people fail to learn the velar softening grammar, and instead rely on memorized forms for existing words. These people are predicted to not apply velar softening in a wug test. Interestingly, in a wug test with nonce stems and the -ity suffix, Pierrehumbert (2006) found that 2 in 10 subjects did not have productive velar softening.

## 7 Conclusion

The present study produced a domain-general model that concurrently learned both hidden structure and a weighted-constraint grammar. The model was trained on eight languages, and generalized well to test data on all of them. Two languages in particular presented a choice between acquiring a grammar that supported rich bases versus one that didn't. This study found a strong preference for acquiring the rich base grammar, which I argued was an emergent property of the model. The present study thus presented a way in which a rich base grammar may be acquired when URs are not known in advance.

## Acknowledgments

I would like to thank Tim Hunter and Kie Zuraw for stimulating discussion which made this paper better. I also thank the three anonymous reviewers for helpful feedback.

<sup>19</sup>The lowest probability was  $Pr(/klɛmɪk-ɪf/, [klɛmɪktʃ]) | \text{CLEMIC-ISH} = 99.983\%$ .

## References

- Paul Boersma. 2001. Phonology-semantics interaction in OT, and its acquisition. In Robert Kirchner, Wolf Wikeley, and Joe Pater, editors, *Papers in Experimental and Theoretical Linguistics*, volume 6, pages 24–35. University of Alberta, Edmonton.
- Paul Boersma and Joe Pater. 2016. Convergence properties of a gradual learning algorithm for harmonic grammar. In John J. McCarthy and Joe Pater, editors, *Harmonic Grammar and Harmonic Serialism*. Equinox Press, London.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. 1977. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B (Methodological)*, 39(1):1–38.
- Sarah Eisenstat. 2009. Learning underlying forms with maxent. Master’s thesis, Brown University.
- Sharon Goldwater and Mark Johnson. 2003. Learning OT constraint rankings using a maximum entropy model. In Jennifer Spenader, Anders Eriksson, and Östen Dahl, editors, *Proceedings of the Stockholm Workshop on Variation within Optimality Theory*, pages 111–20.
- Gaja Jarosz. 2015. Expectation driven learning of phonology, ms. University of Massachusetts, Amherst.
- Aleksei Nazarov and Joe Pater. 2017. Learning opacity in stratal maximum entropy grammar. *Phonology*, 34:299–324.
- Max Nelson. 2019. [Segmentation and UR acquisition with UR constraints](#). *Proceedings of the Society for Computation in Linguistics*, 2:60–68.
- Charlie O’Hara. 2017. How abstract is more abstract? learning abstract underlying representations. *Phonology*, 34:325–345.
- Joe Pater, Robert Staubs, Karen Jesney, and Brian Smith. 2012. Learning probabilities over underlying representations. *Proceedings of the Twelfth Meeting of the Special Interest Group on Computational Morphology and Phonology*, pages 62–71.
- Janet Pierrehumbert. 2006. The statistical basis of an unnatural alternation. In Louis Goldstein, D. H. Whalen, and Catherine T. Best, editors, *Laboratory Phonology VIII, Varieties of Phonological Competence*, pages 81–107. Mouton de Gruyter, Berlin.
- Alan Prince and Paul Smolensky. 2004. *Optimality Theory: Constraint interaction in generative grammar*. Blackwell Publishing, Malden, MA.
- Ezer Rasin and Roni Katzir. 2016. On evaluation metrics in Optimality Theory. *Linguistic Inquiry*, 47(2):235–282.
- Robert Staubs and Joe Pater. 2016. Learning serial constraint-based grammars. In John McCarthy and Joe Pater, editors, *Harmonic Grammar and Harmonic Serialism*. Equinox Press.
- Bruce Tesar, John Alderete, Graham Horwood, Nazarré Merchant, Koichi Nishitani, and Alan Prince. 2003. Surgery in language learning. In G. Garding and M. Tsujimura, editors, *WCCFL 22 Proceedings*, pages 477–90. Cascadia Press, Somerville, MA.
- Bruce Tesar and Paul Smolensky. 2000. *Learnability in Optimality Theory*. MIT Press, Cambridge, MA.
- Kie Zuraw. 2000. *Patterned exceptions in phonology*. Ph.D. thesis, UCLA.