

## **Learning covert phonological interaction: an analysis of the problem posed by the interaction of stress and epenthesis \***

**John Alderete** alderete@rucss.rutgers.edu

**Bruce Tesar** tesar@rucss.rutgers.edu

**Rutgers University, September 2002**

### **Abstract**

The purpose of this report is to contribute to formal learning theory in Optimality Theory by providing an analysis of the problem posed by covert phonological interaction. Building on standard theories of syllable structure and metrical stress, we analyze a typological system in which regular processes of epenthesis interact in non-transparent ways with metrical stress. The nature of this interaction, which implies several complex learning decisions, is shown to support the following conclusions relevant to general mechanisms of language learning:

1. A process of lexical acquisition in which surface phonological structure is directly incorporated into the lexicon (sometimes called 'the identity map') leads to a state in which the learner is committed to a grammar that over-generates.
2. Acquisition of some aspects of a lexical representation (LR) may take place in the absence of morpho-phonemic alternations; acquisition of LRs that are not identical to the surface form is necessary to solve the over-generation problem.
3. Over-generation problems may arise in learning from constraint interactions that are distinct from those that exist among faithfulness constraints that stand in a special/general relation.

## **1. Learning overt versus covert phonological interaction**

### **1.1 Implications of the Subset Problem for lexical acquisition**

A familiar problem in formal learning theory is the Subset Problem ((Angluin, 1980) and (Baker, 1979)). Linguistic theory provides typologies that are rich in logical structure: for example, the structural inventories produced by some grammars are a proper subset of the inventories of other grammars in the same typological space. Because of these relations between inventories, the linguistic input available in learning may be consistent with more than one grammar. This makes it possible for a learner to select a grammar that has greater generative capacity than necessary, while remaining consistent with the available data. Moreover, learning takes place on the basis of positive evidence, so if a learner does postulate a grammar that over-generates, no amount of additional evidence will be inconsistent with this incorrect grammar. Approaches to the Subset Problem therefore provide methods for computing the most restrictive grammar given the available data.

The standard analysis of the Subset Problem in OT models of early learning (i.e., before morphological analysis) involves placing an inductive bias for grammars in which markedness constraints dominate faithfulness constraints. This preference can be attempted in different ways, either as a condition on the 'initial state' ((Gnanadesikan, to appear), (Hayes, to appear), (Levelt, 1995), (Smolensky, 1996)), or as a durative principle that persists throughout grammar learning (Prince and Tesar, to appear), but the

---

\* This report has benefited greatly from conversations with Graham Horwood, I-Ju Sandra Lai, Koichi Nishitani, and Alan Prince. It is supported in part by NSF BCS-0083101 and a NIH NRSA training grant awarded to the Rutgers Center for Cognitive Science (NIH 1-T32-MH-19975-05). If any errors remain, despite this help, we alone are responsible.

utility of the markedness-over-faithfulness (**M>>F**) bias in grammar learning is clear. Markedness constraints characteristically increase restrictiveness by asserting phonological requirements of various types. For example, the markedness constraint ‘if heavy then stressed’ (i.e., the Weight-to-Stress Principle of (Prince, 1990)), when active, eliminates the potential for lexical stress to emerge in words of uneven quantity, effectively reducing the inventory of stress patterns. Ranking markedness constraints as high as possible is therefore one way of preferring the most restrictive grammar consistent with the data.

In the standard analysis, the **M>>F** bias works in tandem with a specific hypothesis for the incorporation of surface phonological structure into the lexicon. We call this proposal the Identity Map Hypothesis, after (Hayes, to appear) and (Prince and Tesar, to appear).

(1) The Identity Map (IM) Hypothesis

The phonological content of surface forms is mapped directly into lexical representations.

The identity map is essential to discovering the correct system of contrast. Contrast maintenance in OT is the job of faithfulness constraints, so for these constraints to be correctly inserted in the constraint hierarchy, lexical representations must contain relevant structure to be faithful to. Since surface forms can be relied upon (absolute neutralization aside) to exhibit these contrasts, IM-style lexical acquisition ensures that the lexical representations will have this structure. Furthermore, the behavior of non-contrasting structure is handled with the grammar, since the **M>>F** imperative presupposes that the distribution of non-contrasting phonological structure is under grammatical control, i.e., governed by markedness.

The potential for contrast sponsorship created by the identity map, and mitigated by the **M>>F** bias, has been shown to be successful in modeling a host of examples of early learning ((Hayes, to appear), (Prince and Tesar, to appear), and (Tesar, to appear)). In these test cases, the phonological structures to be learned are directly accessible in the output form, so they represent a kind of *overt* phonological interaction. As we will show below, however, when crucial aspects of phonological structure are not overt in the surface form, the program for grammar learning outlined above fails to converge on the most restrictive grammar consistent with the available data. Of particular interest is that such covert interaction can be exhibited independently of evidence from phonological alternations. The chief aim of this report is to clarify the problems for learning implicated by this covert phonological interaction.

## 1.2 BCD and overt phonological interaction

Consider the interaction between stress and vowel quality illustrated below, a well-attested kind of phonological interaction.<sup>1</sup>

(2) Illustration of overt phonological interaction

- |                                  |      |      |      |      |      |      |
|----------------------------------|------|------|------|------|------|------|
| a. Lg A. Free Stress             | páka | paká | pákə | paké | pəka | pəká |
| b. Lg B. Lexical Stress with *ə́ | páka | paká | pákə |      |      | pəká |

Lg B differs from A in having a significant interaction between stress and vowel quality: stress never appears on a schwa in Lg B. This interaction is significant because it poses limitations on the set of surface forms occurring in Lg B. The repulsion of stress from schwas in (2b) entails that the forms allowed in B are a proper subset of those allowed in A.

The difference between Lg A and B resides in an observable property of their output forms: Lg A has stressed schwas but B does not. The overt interaction between stress and vowel quality in Lg B can thus be attributed to the markedness constraint \*ə́, which simply bans stressed schwas (Cohn and McCarthy, 1998). Stressed schwa markedness is therefore in conflict with the faithfulness constraint, FAITHACCENT, which requires every lexically specified stress to surface in the output form on its lexical sponsor, i.e., on the syllable in the input that supports the prosodic structure for stress (see (Alderete, 2001)). In this analysis, the formal relationship between stress and vowel quality is directly characterized by markedness.

We assume that grammar learning, given overt data like that in (2), involves two separate processes: interpretive parsing, for assigning full structural descriptions to the overt data (Tesar, 1995) (Tesar and Smolensky, 2000), and Biased Constraint Demotion (BCD, (Prince and Tesar, to appear)). BCD instantiates the desired **M>>F** bias: it will insert a markedness constraint in the hierarchy above a

---

<sup>1</sup> See for example (Cohn, 1989), (Cohn and McCarthy, 1998) for an interaction of this type in Indonesian, and (Kenstowicz, 1994) for a host of additional examples.

faithfulness constraint whenever there is a free choice between the two, and given a choice of more than one faithfulness constraint, it ranks the one that frees up the most markedness constraints for ranking in a subsequent pass of the algorithm (see (Prince and Tesar, to appear) for the details of the algorithm). Interpretive parsing involves selection of a lexical representation, which is the main topic of this paper. The Identity Map hypothesis should be seen as a claim about the lexical dimension of interpretive parsing.<sup>2</sup>

Given the overt data in (2b), BCD will always rank \* $\acute{\text{e}}$  over FAITHACCENT, since there are no overt forms that contradict it and this ordering respects the **M>>F** bias. Effectively, BCD presupposes that there is a significant interaction between stress and vowel quality, unless given positive evidence to the contrary.

The only kind of evidence that can motivate FAITHACCENT >> \* $\acute{\text{e}}$  are forms with a stressed schwa, like [pak $\acute{\text{e}}$ ] from the free stress language (2a). Since this form violates \* $\acute{\text{e}}$ , it can only win out over the plausible competitor [p $\acute{\text{a}}$ k $\acute{\text{e}}$ ] by satisfying FAITHACCENT. Importantly, the latter constraint will only have force if there is a lexically-specified accent in the lexical representation; without accent, faithfulness is silent. Because the loser [p $\acute{\text{a}}$ k $\acute{\text{e}}$ ] satisfies \* $\acute{\text{e}}$  (hence the ‘L’ in the column for this constraint below) regardless of the assumed LR, the only way to make [pak $\acute{\text{e}}$ ] the winner is to rank FAITHACCENT above \* $\acute{\text{e}}$ , and to have this constraint prefer the winner. But FAITHACCENT only has this role if surface stress is projected back into the lexical representation.

(3) Faithfulness solution requires the identity map

| /pak $\acute{\text{e}}$ /  | * $\acute{\text{e}}$ | FAITHACCENT |
|--|----------------------|-------------|
| pak $\acute{\text{e}}$ ~ p $\acute{\text{a}}$ k $\acute{\text{e}}$ | L                    | W           |

The IM is a natural solution to the problem of incorporating marked structure in a developing system of contrast. When confronted with markedness-violating structure in the output, the IM prescribes incorporating this structure in the lexical representation. A fundamental assumption of this analysis is therefore that any faithfulness-based solution to attested marked structure (e.g., stressed schwas) must involve preservation of that structure in the input. It is exactly this assumption that leads to an instance of the Subset Problem in learning covert interaction, and consequently, this finding will lead us in section 3 to suggest other mechanisms for learning the lexicon.

### 1.3 Consequences of the identity map for covert phonological interaction

We choose to exemplify our point with a very common type of covert phonological interaction, the interaction between regular processes of vowel insertion and stress assignment. Many languages avoid assigning surface stress on a vowel that is not present in the lexical representation of a word (see (Broselow, 1982) and (Alderete, 1999) for a variety of examples). We illustrate covert phonological interaction with the language types given below.<sup>3</sup>

(3) Illustration of covert phonological interaction

- a. Lg A. Free Stress                      p $\acute{\text{a}}$ kat    pak $\acute{\text{a}}$ t    p $\acute{\text{a}}$ kit    pak $\acute{\text{i}}$ t    p $\acute{\text{i}}$ kat    pik $\acute{\text{a}}$ t    p $\acute{\text{i}}$ kit    pik $\acute{\text{i}}$ t
- b. Lg B. Stress-Epenthesis                      pak $\acute{\text{a}}$ t    p $\acute{\text{a}}$ k $\acute{\text{i}}$ t    pak $\acute{\text{i}}$ t                      pik $\acute{\text{a}}$ t    p $\acute{\text{i}}$ k $\acute{\text{i}}$ t    pik $\acute{\text{i}}$ t
- c. Lg C. Final Stress                      pak $\acute{\text{a}}$ t                      pak $\acute{\text{i}}$ t                      pik $\acute{\text{a}}$ t                      pik $\acute{\text{i}}$ t

The phonological interaction of interest here is exemplified in Lg B, which avoids stressing epenthetic vowels. In this typology, the quality of epenthetic  $\acute{\text{i}}$  is identical phonetically to lexical high front vowels, so, while they are identified with underlining above for concreteness, this information is not available in the acoustic signal. Because of this identity, the forms of Lg B are a proper subset of those of Lg A and a superset of those of Lg C. Lg B does not have fully contrastive stress, like Lg A, because the stress contrast above is only in words that end in high front vowels, i.e., the vowel inserted by epenthesis. Lg B

<sup>2</sup> This paper is concerned with hidden lexical structure. Previous work ((Tesar, 1997), (Tesar, 2001)) has focused on hidden structure that is not specifically lexical but missing from the overt forms.

<sup>3</sup> Concrete examples exhibiting the pattern of stress-epenthesis interaction sketched in (3b) include Mohawk and Selayarese, whose rightward oriented stress is shifted back by epenthesis, and Yimas, in which leftward oriented stress is shifted forward by epenthesis.

does not have fully predictable stress either, as in Lg C, because stress-epenthesis interaction creates surface minimal pairs like [pákit] versus [pakít].

In modeling these systems, two new constraints are needed. One constraint, MAINRIGHT, describes the surface true generalization in Lg C that stress is on the rightmost syllable. This constraint is dominated in the grammar for Lg A by FAITHACCENT to account for the free distribution of stress, and in Lg B by the positional faithfulness constraint, HEADDEP, defined below.

(4) HEADDEP (Alderete, 1999)

Nonlexical vowels are not allowed in prosodic heads (=stressed syllables)

This constraint accounts for the observed avoidance of stressing epenthetic vowels. It is a faithfulness constraint because its evaluation requires access to the lexical representation, a conclusion that is consistent with its behavior in typology, because elevating HEADDEP in the constraint hierarchy introduces surface contrast (see section 2).

Given that both HEADDEP and FAITHACCENT are faithfulness constraints, when presented with just forms with final stress, BCD will rank MAINRIGHT above these constraints. This state obtains in the language with stress-epenthesis interaction (3b) until the learner confronts forms with non-final stress, like [pákit], which are not optimal with top-ranked MAINRIGHT. There are two ways of making this form optimal: through lexical specification of stress and promotion of FAITHACCENT (5), or via epenthesis and consequent activation of HEADDEP (6).

(5) The faithfulness solution 1: identity between surface and lexical forms

| /pákit/       | HEADDEP | FAITHACCENT | MAINRIGHT |
|---------------|---------|-------------|-----------|
| pákit ~ pakít |         | W           | L         |

(6) The faithfulness solution 2: non-identical lexical form

| /pakt/        | HEADDEP | FAITHACCENT | MAINRIGHT |
|---------------|---------|-------------|-----------|
| pákít ~ pakít | W       |             | L         |

Both approaches to non-final stress involve domination of a markedness constraint by a faithfulness constraint. In solution 1, faithfulness directly preserves accent, in accord with the assumptions supporting the identity map. In solution 2, by contrast, faithfulness says nothing about lexical accents; accent is not present in the LR. The introduction of marked non-final stress comes somewhat indirectly, via epenthesis, and the imperative to avoid stress on an epenthetic vowel. Solution 2 is a faithfulness solution because it involves promoting the faithfulness constraint HEADDEP, but it differs from solution 1 in that faithfulness is not obviously engaged in preserving phonological structure that is specified lexically.

A second important difference between the two approaches is that solution 2 requires an input that is non-identical with the surface form. Why should the learner ever prefer the second solution, then, since it violates the identity map hypothesis? The answer to this question is clear: because the grammar produced by the first strategy is less restrictive than the grammar resulting from the second. Solution 1 yields a grammar with free stress, whereas solution 2 produces the limited stress contrast stemming from stress-epenthesis interaction. Respecting the identity map, therefore, results in the grammar of the superset language (Lg A) on the basis of data from the subset language (Lg B) alone. Furthermore, Solution 1 cannot be saved by a weaker interpretation of the identity map in which only selected phonological structures are mapped into the lexicon. If, for example, stress is not lexically specified but consonant and vowel structure is, then the observed form becomes unavoidable suboptimal because there is no way to resolve the unchecked loser marks incurred by MAINRIGHT.

## 2. A concrete illustration of the Subset Problem

This section fills in the details missing from the argument presented above in outline form. In particular, we give a concrete illustration of the Subset Problem by providing explicit analyses of stress and syllabification in the language types in (3). These analyses reveal important distributional differences and confirm that the overt forms of a language with regular stress and stress-epenthesis interaction (3b) form a proper subset of the overt forms possible in a free stress language (3a).

### 2.1 Properties of a language with epenthesis

For concreteness, our test case is modeled after the system of syllabification in Iraqi Arabic ((Broselow, 1982), (Itô, 1989)). Though there are a few characteristics of Arabic not represented,<sup>4</sup> the properties of ‘Pseudo-Arabic’ presented below are predicted by a particular ranking of generally espoused constraints on segment and syllable structure. In particular, the markedness and faithfulness constraints given in (7) below are ranked as shown in (8) to produce Pseudo-Arabic syllables.

(7) Constraints implicated by epenthesis (see (McCarthy and Prince, 1995), (Prince and Smolensky, 1993))

- a. MAXC: every consonant in the input has a corresponding element in the output
- b. DEP<sub>V</sub>: every vowel in the output has a corresponding element in the input
- c. \*COMPLEX: sequences of two tautosyllabic consonants are not allowed
- d. NOCODA: syllable-final consonants are not allowed
- f. RIGHTMOSTCLOSED: closed syllables are more harmonic closer to the right edge of the word

(8) Constraint rankings for Pseudo-Arabic syllables

- a. MAXC >> DEP<sub>V</sub>: vowel insertion is better than consonant deletion
- b. \*COMPLEX >> DEP<sub>V</sub> >> NOCODA: CV(C) syllable template
- c. DEP<sub>V</sub> >> RIGHTMOSTCLOSED: right-to-left directional syllabification; favors open syllables closer to the beginning of the word (after (Mester and Padgett, 1994))

The syllabification system above resolves tautosyllabic consonant clusters with epenthesis, building syllables from right to left. Since vowel structure and the possibility of onsetless syllables are orthogonal to our study, they are ignored here.

### 2.2 Factoring stress into the combinatorics

Now we can examine the influence of epenthesis on stress and distinguish it empirically from languages that do not have this type of interaction. An additional edgemost constraint, MAINLEFT, is added to the constraint set from section 1 to illustrate the ranking for rightmost stress. With these constraints in hand, the three language types discussed in section 1, namely Lg A (3a), Lg B (3b), and Lg C (3c), can be described with the constraint rankings given below.

(9) A three-way typology of stress systems

- a. Lg A: free stress  
FAITHACCENT >> {MAINRIGHT, MAINLEFT} >> HEADDEP
- b. Lg B: rightmost stress, with stress-epenthesis interaction  
HEADDEP >> MAINRIGHT >> MAINLEFT >> FAITHACCENT
- c. Lg C: rightmost stress, no stress-epenthesis interaction  
MAINRIGHT >> {MAINLEFT, FAITHACCENT, HEADDEP}

Lg A has top-ranked faithfulness, and, accordingly, has emergent lexical accent in all positions. Lg B, by contrast, has a default position for stress at the right edge of a word, except in words that have final epenthetic syllables, in which case stress appears on the rightmost syllable that contains a lexical vowel. Lg C differs from both of these in that both FAITHACCENT and HEADDEP are dominated by an edgemost constraint, namely MAINRIGHT, so this language has fixed final stress. These languages do not differ, however, in the possibility of epenthesis; they all have the rankings in (8).

---

<sup>4</sup> The chief characteristic missing here is the availability of final superheavy syllables, which must be absent in order to study the impact of epenthesis on final stress.

With the analysis of epenthesis and stress in hand, we may ask, what are the available surface forms in each language, i.e., the inventory of word types based on the phonological structures that can be used to distinguish words overtly? For concreteness, we have fixed the number of syllable types to two (i.e., CV and CVC, the syllable types predicted by the grammar in (8)), and the number of lexical vowels in the demonstration below to five, a typologically common vowel inventory (Maddison, 1984), but nothing crucial to our argument hinges on these choices. By factoring stress, syllable type, syllable count, and vowel type into a set of formulas given in the appendix, we calculate the quantities of distinct surface forms in the Pseudo-Arabic typology, shown below.

(10) Quantities of phonologically distinct surface forms in stress typology

| Syllable Count | Language A | Language B | Language C |
|----------------|------------|------------|------------|
| 1              | 10         | 10         | 10         |
| 2              | 200        | 110        | 100        |
| 3              | 3,000      | 1,110      | 1,000      |
| 4              | 40,000     | 11,110     | 10,000     |
| Total          | 43,210     | 12,340     | 11,110     |

These numbers alone do not confirm the claimed subset relations asserted in section 1. To show that Lg B is a subset of Lg A and Lg C is a subset of Lg B, we must also show that all the forms of B are inside the inventory of A, and likewise for C and B.

Since the range of syllable and vowel types is held constant in all three languages, the inventory differences shown above must stem from differences in the range of possible stress patterns. Every form of Lg C has final stress. Each form of C is therefore achievable in the other two languages by specifying all vowels and accenting the final syllable. Since both A and B have surface forms not present in C (i.e., those with non-final stress), C is a proper subset of A and B. Furthermore, every form of Lg B is attainable in Lg A by full specification of vowel types and lexical stress. Again, since Lg A has a wide range of forms not represented in B (i.e., those with non-final stress followed by non-epenthetic vowels), B is a proper subset of A. Thus, a language with stress-epenthesis interaction (Lg B above) has a proper subset of the forms of a language with free stress (Lg A). However, BCD and the identity map approach to lexical acquisition allow Lg A, the superset language, to be selected by the learner on the basis of data from Lg B alone. The next question is, what are some possible approaches to learning covert phonological interaction that do not have this problem?

### 3. Possible approaches to learning

#### 3.1 The target grammar and lexicon

The correct grammar for Lg B is one in which a top-ranked HEADDEP induces a limited contrast in stress. For this constraint to be active, however, the inputs that lead to vowel insertion must be learned.

(11) Target grammar and lexicon representations

a. Grammar: HEADDEP >> MAINRIGHT >> MAINLEFT >> FAITHACCENT

b. Lexicon: /pakit/ → [pakít] and /pakt/ → [pákít], etc.

The challenge for the learner is to somehow motivate the acquisition of inputs that will lead to the correct distribution for epenthetic vowels, which in turn, implicates HEADDEP in grammar learning. Below, we sketch some approaches for achieving this result.

#### 3.2 Waiting for alternations?

The covert structure in stress-epenthesis interaction is the lexical representation (LR). A standard view on learning non-identical LRs is that they are learned from alternations. Allomorphy in general complicates the identity map, since it is not clear how to map multiple morpheme shapes onto a single lexical entry. Evidence from alternations showing epenthesis might therefore subvert the expectations of the IM hypothesis, and, once the correct mechanisms are known for sorting out allomorphy, help the learner select the correct LRs shown in (11) above.

We believe that such a solution may be viable for some languages, in particular, languages in which there is sufficient evidence from alternations. Though the learner may have to postpone some phonotactic learning to a much later stage in the development process, i.e., past the stage at which the correct morphological segmentation has been made, alternations provide a tractable way of learning the correct LR<sub>s</sub>, and, in turn, the correct grammar, which crucially depends on these LR<sub>s</sub>.

However, waiting for alternations will not be viable for all languages, because in some languages with documented stress-epenthesis interaction, there is not good evidence for alternations in the tokens of interest, so postponing grammar-learning will not be guided later by the right LR<sub>s</sub>. For example, in the Papuan language Yimas (Foley, 1991), epenthesis is motivated statically by the phonotactics of the language, so there are forms with epenthetic vowels that shun stress, but no morphologically related forms signaling the  $V \sim \emptyset$  alternation. Indeed, the only sign that these vowels are epenthetic is their impact on the stress system. This case, and others like it, represent a kind of absolute neutralization, where the only evidence for the non-identical LR<sub>s</sub> comes from system-wide observations about the distribution of the epenthetic vowel. These are observations that may encompass stress and other overt structure, but crucially lack dynamic alternations. To make some headway with these cases, therefore, it is necessary to have some general strategy for learning non-surface true LR<sub>s</sub> in phonotactic learning, without LR-motivating alternations.

### 3.3 A suggestion for revising the r-measure

A language with stress-epenthesis interaction (3b) was shown above to be more restrictive than a language with free stress (3a) in that the latter allows more surface forms and it has all of the forms of the former. However, a proposed measure of restrictiveness, Prince & Tesar's r-measure, does not currently account for the difference between these two cases. The r-measure of a language is calculated by counting, for each faithfulness constraint **F**, the number of markedness constraints that dominate **F**. The r-measures for (3a) and (3b) are computed below, holding constant the shared rankings for epenthesis (8).<sup>5</sup>

(12) R-measures of languages A (3a) and B (3b)

a. Rankings for Lg A (r=2): FAITHACCENT >> {MAINRIGHT, MAINLEFT} >> HEADDEP

b. Rankings for Lg B (r=2): HEADDEP >> MAINRIGHT >> MAINLEFT >> FAITHACCENT

Despite the observational differences between A and B, they are equal in r-measure. The chief difference between the rankings in (12a) and (12b) is the relative ranks of FAITHACCENT and HEADDEP. Perhaps a way of addressing the Subset Problem clarified above is to revise the r-measure such that all faithfulness constraints are not evaluated uniformly.

With this in mind, the alignment constraints MAINRIGHT and MAINLEFT seem to have more significant consequences for restrictiveness when they dominate FAITHACCENT than when they dominate HEADDEP. Put crudely, alignment seems to have more 'mashing power' in terms of suppressing surface contrast in the former case. After all, FAITHACCENT can license a stress contrast anywhere in the word, whereas HEADDEP only introduces a stress contrast if the system of syllabification happens to posit an epenthetic vowel in the canonical position for stress. Working with the numbers in (10), demoting FAITHACCENT to a position below alignment reduces the inventory by 74% (43,210 forms to 11,110), while the demotion of HEADDEP only shrinks the inventory by 10% (12,340 to 11,110). A formalization of these differences, i.e., one that accounts for the actual content of the constraints and the units they refer to, may provide a better empirical basis for quantifying restrictiveness. We conjecture that when more accurate measures are available, such measures will make it possible to instantiate a bias for the more restrictive grammar of (3b). Finally, we note that revising the r-measure in this way still leaves the problem of learning the correct LR<sub>s</sub> that activate HEADDEP, but perhaps awareness of a more restrictive faithfulness constraint could also motivate LR<sub>s</sub> that would capitalize on it. How exactly this would be achieved, we leave as a problem for future research.

---

<sup>5</sup> Though (3a) does not in fact require the epenthesis ranking, it is necessary that some appropriate faithfulness constraint dominate NOCODA, and that \*COMPLEX dominate either MAXC or DEPV, to derive the same syllable structures.

### 3.3 Phonologically based contrast classes

Recently, a number of researchers have proposed mechanisms for examining the linguistic structures of a language globally and inferring certain characteristics about the language in this way. For example, (Frisch et al., 1997) provide a similarity metric for making generalizations over the lexicon, and (Flemming, 1995) gives structure to a dispersion theory for maximizing the perceptual distance among phonemes by making system-wide comparisons (see also (Padgett, to appear)). If one admits this kind of power, either as part of the grammar itself or the mechanisms involved in language learning, this may give the learner a background for setting up the correct LRs for stress-epenthesis interaction. Consider for example the idealized structures below.

(24) Idealized structures for languages A (3a) and B (3b)

|      |       |       |                |       |       |       |                |       |
|------|-------|-------|----------------|-------|-------|-------|----------------|-------|
| Lg A | CáCaC | CaCáC | CáCiC          | CaCiC | CíCaC | CiCáC | CíCiC          | CiCiC |
| Lg B |       | CaCáC | CáC <i>í</i> C | CaCiC |       | CiCáC | CíC <i>í</i> C | CiCiC |

The boxes above represent the stress contrasts in the logically possible CVCVC sequences. If allowed to arrange the data in this way, the learner may observe that Lg B only has a contrast in words with the following skeletal profile: CVCiC. In terms of surface stress contrast, Lg B only distinguishes words that end in *i*. This observation may enable the learner to figure out that a contrast can be achieved in ways other than representing stress lexically, since there is some regularity in the system of contrast. For example, if the final *i* was removed, the system can rely on a phonological process to insert it in the correct position, and therefore account for the limited contrast with lexical specification of segmental structure alone. Considerations such as these may provide a means of arriving at the target lexicon, without appeal to evidence from alternations.

### 4. Conclusion

The interaction of two or more phonological structures may be covert in the sense that the analysis of these structures is not immediately apparent from the overt data available in the acoustic output. In this report, the problem of learning the grammar of stress was studied in systems where surface stress interacts significantly with the hidden lexical structure implied in the analysis of vowel epenthesis, leading to the following conclusions:

1. *The Subset Problem and stringency*: it is not just a consequence of special/general relations among faithfulness constraints.
2. *Restricting identity mappings*: lexical acquisition needs to be constrained such that not all aspects of a surface form are directly incorporated into the lexicon.
3. *Non-identical LRs without evidence from alternations*: non-surface true lexical representations must be posited, even without alternations to support them.

Recent work in OT on the Subset Problem has largely focused on phonotactics in early learning, identifying a set of problems that result from decisions that can be made in ranking faithfulness constraints. These problems were initially found with positional faithfulness constraints ((Hayes, to appear) and (Smith, 2000)), because they characteristically represent set-superset relations among faithfulness constraints, and so their insertion in the constraint hierarchy has the potential to over-generate. (Prince and Tesar, to appear) enriches this discussion by showing that special/general relations can in fact be derived for faithfulness constraints that don't naturally represent a set-superset relation, through the ordering of other constraints in the grammar. This paper extends the discussion of the problem further by showing that the Subset Problem rears its head even in contexts in which the relevant constraints to be ranked are not in a special/general relation, either intrinsically or derived through other orderings. The violations profiles of HEADDEP and FAITHACCENT do not have the entailment relations of constraints that stand in a special/general relation. The subset relation between the languages results not just from the relation between the two constraints, but the different lexical hypotheses as well.

The Subset Problem exhibited by languages with stress-epenthesis interaction also leads to a conclusion about the nature of mappings from overt forms to lexical forms. When stress is directly incorporated in the lexicon, it leads to a learning *faux pas*, effectively activating FAITHACCENT and leading to the acquisition of the superset language with free stress. In order to avoid this incorrect outcome, it

seems prudent to allow for processes of lexical acquisition that do not directly represent all potentially phonemic properties of the acoustic signal. In other words, a strict identity mapping does not seem to be tenable for all languages.

This conclusion converges nicely with some of the results and conclusions of (Pater, to appear). In this work, Pater argues for two sets of faithfulness constraints to account for the well-known divergences between linguistic comprehension and production. One type of faithfulness regulates the mapping from overt linguistic data to the forms stored by the learner in the lexicon. Yet another governs the faithfulness between lexical forms and forms produced by the learner. One of the basic arguments for this division of labor is that the same markedness constraints are active in both receptive competence and production. The influence from markedness on comprehension has the consequence of producing non-strict mappings from overt data to lexical forms, which is exactly the conclusion arrived at here, on the basis of different data. We believe, therefore, that our analysis provides a learnability basis for Pater's conclusion that early learners may posit non-identical lexical forms.

Finally, perhaps one of the most interesting findings of the study is that it provides a learning theoretic motivation for non-identical LRs, even in the absence of evidence from alternations. Pater's work aside, it has been standardly assumed in OT formal learning theory that alternations are the only type of data that could possibly motivate LRs which are distinct from surface forms. In phonotactics, for example, as long as the faithfulness constraints are ranked as low as possible, there appears to be no harm in including the phonological detail available in output forms (Prince and Tesar, to appear), hence, the identity map. However, the case of stress-epenthesis interaction shows that perfect replication of the overt form can in fact have negative consequences. As shown in section 1, full specification of stress leads to the Subset Problem. In particular, the only way to motivate the correct path for grammar-learning is to first acquire the LRs consistent with epenthesis, which in turn activates HEADDEP. Moreover, some languages, like Yimas, show that alternations cannot be relied upon to motivate the correct LRs, so lexicon-learning must embody inferences from purely distributional data.

As discussed in section 3, we believe that two ideas may prove useful in providing tractable lines of analysis for learning the correct LRs without evidence from alternations. First, we conjecture that a revised r-measure may enable the learner to find LRs that, through grammar learning, maximize the r-measure, one of the intrinsic properties of BCD. Second, we suggest that phonologically-based contrast classes, classes that assess the extent of a phonological contrast syntagmatically by evaluating global properties of the system, may enable the learner to construct a lexicon that leads to grammar-learning that also avoids the Subset Problem illustrated above.

### ***Appendix: possible words in the Pseudo-Arabic typology***

How does one calculate the number of possible word types in each of the language types from (3) above? Because the effect of stress placement is different in each system, separate formulas are called for in counting the surface forms in each case.

For Lg C, stress is predictable, so stress does not factor in the formula. The number of distinct surface forms is thus the number of distinct sequences of vowels and syllable types, stated below as a function of form length in syllables (*fmLen*):

$$\text{Surface forms in Lg C: } \text{allForms}_{LgC}(fmLen) = (2lexVow)^{fmLen}$$

In disyllabic forms, for example, a language with five lexical vowels and the two syllable types of Pseudo-Arabic (CV and CVC) has  $(2 * 5)^2 = 100$  different surface forms, as shown in the table in (10) from section 2.

Lg A, on the other hand, has contrastive stress, so the calculation of surface forms must factor in stress differences. Each sequence of distinct vowel and syllable types can be lexically specified for stress, so the number of distinct forms per sequence is the number of syllables ( $=fmLen$ ):

$$\text{Surface forms in Lg A: } \text{allForms}_{LgA}(fmLen) = fmLen * (2lexVow)^{fmLen}$$

Thus, the number of surface forms in Lg A differs from Lg C by a factor of *fmLen*: Lg A has, for example, 200 distinct disyllabic forms, as opposed to the 100 forms of Lg C. As a result, the number of forms in Lg A grows faster than it does in Lg C as the number of syllables increases, as shown in (10).

In Lg B, stress appears on the rightmost lexical vowel, so the impact of epenthesis on stress is that it may cause non-final stress if an epenthetic vowel occurs to the right of the rightmost lexical vowel. Concretely, the word type CV.CVC will cause initial stress, because of the top-ranked status of HEADDEP. The number of different surface forms therefore involves summing the distinct vowel/syllable type sequences for all of these stress permutations:

Surface forms in Lg B:

$$lexForms(fmLen) = (2lexVow)^{fmLen}$$

$$allForms_{LgB}(fmLen) = \sum_{x=0}^{fmLen-1} lexForms(fmLen - x) = \sum_{x=0}^{fmLen-1} (2lexVow)^{fmLen-x}$$

The factor notation above describes the fact that the number of distinct surface forms with stress moved  $x$  syllables from the right ( $x = 0$  if stress is final) is equal to the number of lexical forms of length  $fmLen - x$ , because all of the syllables following the stressed one have the same vowel (the epenthetic one) and the same shape, namely CVC.<sup>6</sup> Since, by assumption, there is only one epenthetic vowel and the quality of the epenthetic vowel is inside the inventory of lexical vowels, epenthetic vowels only introduce distinct stress patterns (not new syllable types or vowel types), but this assumption is qualified below in languages that have more than one epenthetic vowel or epenthetic vowels that are outside the set of lexical vowels. Note that the quantities in (10) only go to  $fmLen - 1$ , because if  $x = fmLen$ , then all vowels are epenthetic, and the grammar in (9b) predicts default rightmost stress in such a case, so these structures are identical to forms with final stress and lexical vowels identical to the epenthetic vowel.

It is important to point out that many of the formal properties of systems studied here are logically independent of the subset relations confirmed in section 2. For example, these relations are independent of language particular vowel inventories and syllable structures. Thus, increasing or decreasing the number of vowels or syllable shapes does not effectively change the logical structure of the disparities among the three languages above. To be concrete, if the set of lexical vowels in all of these languages is three instead of five, as chosen in section 2, the differences in disyllabic word types is 36 (Lg C), versus 42 (Lg B), and 72 (Lg A), which has the same pattern of increased number of forms as we go from C to B, and B to A. Likewise for syllables, if there is just one syllable type, the unmarked CV syllable, there are still important differences stemming from the lack of lexical vowels in the default position for stress. Indeed, the only context where Lg B is no longer a proper subset of A is one in which there is only one lexical vowel and one syllable type. In this scenario, however, there is no correlation to be made between vowel type and stress, since there are no distinctions in vowel type, so there is really no interesting interpretation to be made for such a case.

We also made a simplifying assumption that the languages of the Pseudo-Arabic typology all have a single epenthetic vowel. Importantly, if the number of epenthetic vowels increases, this change does not affect the subset relations, as long as the inserted vowels are identical to one of the vowels of the lexical inventory. For example, if a language, Lg B', has two epenthetic vowels, the number of surface forms is increased, because both epenthetic vowels can cause 'irregular stress' and therefore introduce additional stress patterns in words with epenthetic vowels in the right positions. But this will never create new forms that are not matched by the superset language Lg A', because stress is free in Lg A', so it has all logical stress-vowel quality permutations. The only cases we need to allow for are languages in which the set of epenthetic vowels is identical to the set of lexical vowels. In such a case, the forms of Lg B' are no longer a *proper* subset of Lg A'; the two inventories are identical. There are no adverse consequences to the learner selecting Lg A'.

---

<sup>6</sup> This limitation to post-tonic closed syllables is a consequence of right-to-left syllabification, which prefers open syllables closer to the beginning of the word, and the ranking DEP V >> NOCODA, which prohibits multiple epenthetic vowels when one would do to satisfy the syllable canons. The formula for *allForms* is therefore only limited to Lg B, since different patterns of epenthesis would result in different syllabification possibilities and therefore change the number of syllable types.

It is worth mentioning that, in languages with more than one epenthetic vowel, it is most often the case that there is one vowel that is inserted for phonological reasons, and other ‘epenthetic’ vowels that are not purely phonological and are in fact motivated morphologically or phonetically.<sup>7</sup> In Mohawk, for example, there is the phonological epenthetic vowel e, and two other vowels, i.e., prothetic i and ‘stem-joiner’ a that are called upon in certain morphologically defined environments (Michelson, 1988). For example, prothetic i is used only in verbs to augment the word to satisfy Word Binariness. The fact that prothetic i is restricted to verbs is symptomatic of a morphological analysis, and indeed, a coherent analysis has been proposed for an entirely parallel case of ‘the peg element’ in Athabaskan languages (see (Hargus and Tuttle, 1997)). The important point is that if these ‘other vowels’ are inserted for non-phonological reasons, it is not clear that HEADDEP says anything about them, since they are either present underlyingly (morphologically motivated), so they do not trigger a violation of HEADDEP, or they are phonetic structures that this constraint doesn’t refer to. Thus, these cases seem to involve a certain degree of complexity that the constraints that define our typology do not make any predictions about.

A final variation on the above typology does actually break down the logical relation between set and superset languages. If the epenthetic vowel inserted by the language with stress-epenthesis interaction (Lg B’’) is not inside the inventory of lexical vowels, and the superset language (Lg A’’) does not contain this vowel, then Lg B’’ is no longer a subset of Lg A’’ at all, since it has surface forms that A’’ doesn’t. For example, suppose that schwa is inserted phonologically to resolve biconsonantal clusters, but it is not used contrastively in the lexicon. There is a spurious superset language Lg A’’ that has free stress and all of the vowels of Lg B’’ except schwa; Lg A’’ is not a superset language anymore because it doesn’t have words with schwas. This class of cases presents an interesting learning challenge, namely, learning the distribution of schwas via epenthesis, but it is conceptually different, we believe, from the cases above, because in this case, there is overt structure available to the learner that can help in learning. The existence of extra-phonemic schwa is overt structure that can be related to irregular stress in some way. This case is more complicated than the learning setting presented in section 1 for languages with lexical schwa that can’t be stressed, since the problem is to learn the behavior of non-phonemic schwa. However, since it is conceptually distinct from learning non-overt phonological interaction, we will set this problem aside for the time being.

## References

N.b., ROA = Rutgers Optimality Archive, <http://roa.rutgers.edu/>

- Alderete, John. 1999. Head dependence in stress-epenthesis interaction. In: *The derivational residue in phonological Optimality Theory*, ed. by Ben Hermans and Marc van Oostendorp, 29-50. Amsterdam: John Benjamins. ROA-453.
- . 2001. Morphologically governed accent in Optimality Theory: Outstanding dissertations in Linguistics. New York: Routledge. ROA-309.
- Angluin, Dana. 1980. Inductive inference of formal languages from positive data. *Information and Control*, 45.117-35.
- Baker, C.L. 1979. Syntactic theory and the projection problem. *Linguistic Inquiry*, 10.533-81.
- Broselow, Ellen. 1982. On predicting the interaction of stress and epenthesis. *Glossa*, 16.115-32.
- Cohn, Abigail. 1989. Stress in Indonesian and bracketing paradoxes. *Natural Language and Linguistic Theory*, 7.167-216.
- Cohn, Abigail and McCarthy, John. 1998. Alignment and parallelism in Indonesian phonology. In: *Working Papers of the Cornell Phonetics Laboratory* 12, 53-137. Ithaca, NY: Cornell University. ROA-25.
- Flemming, Edward. 1995. Auditory representations in phonology. UCLA, Doctoral dissertation.
- Foley, William A. 1991. *The Yimas language of New Guinea*. Stanford, CA: Stanford University Press.
- Frisch, Stefan, Broe, Michael and Pierrehumbert, Janet. 1997. Similarity and phonotactics in Arabic. Unpublished manuscript, Northwestern University, Evanston, Illinois. ROA-223.

---

<sup>7</sup> This assessment is made on the basis of a survey of 18 languages reported to have more than one epenthetic vowel, compiled from the responses to a query made in 1996 by Deborah Schmidt (see *Linguist List* 7.1379). Thanks to Deborah for sharing this information with us.

- Gnanadesikan, Amalia. to appear. Markedness and faithfulness in child phonology. In: *Fixing priorities: Constraints in phonological acquisition*, ed. by René Kager, Joe Pater and Wim Zonneveld. Cambridge: Cambridge University Press. ROA-67.
- Hargus, Sharon and Tuttle, Siri G. 1997. Augmentation as affixation in Athabaskan languages. *Phonology*, 14.177-220. ROA-191.
- Hayes, Bruce. to appear. Phonological acquisition in Optimality Theory. In: *Fixing priorities: Constraints in phonological acquisition*, ed. by René Kager, Joe Pater and Wim Zonneveld. Cambridge: Cambridge University Press. ROA-327.
- Itô, Junko. 1989. A prosodic theory of epenthesis. *Natural Language and Linguistic Theory*, 7.217-59.
- Kenstowicz, Michael. 1994. Sonority-driven stress. Unpublished manuscript, MIT, Cambridge, MA. ROA-33.
- Levelt, Claartje. 1995. Unfaithful kids: Place of articulation patterns in early child language. Paper presented at Department of Cognitive Science, The Johns Hopkins University, Baltimore, MD.
- Maddison, Ian. 1984. *Patterns of sounds*. Cambridge: Cambridge University Press.
- McCarthy, John J. and Prince, Alan S. 1995. Faithfulness and reduplicative identity. In: *University of Massachusetts Occasional 18, Papers in Optimality Theory*, ed. by Jill Beckman, Suzanne Urbanczyk and Laura Walsh, 249-384. Amherst, MA: Graduate Linguistic Student Association. ROA-60.
- Mester, Armin and Padgett, Jaye. 1994. Directional syllabification in Generalized Alignment. In: *Phonology at Santa Cruz 3*, ed. by Jason Merchant, Jaye Padgett and Rachel Walker, 79-85. Santa Cruz, CA: University of California. ROA-1.
- Michelson, Karen. 1988. *A comparative study of Lake-Iroquoian accent*. Dordrecht: Reidel.
- Padgett, Jaye. to appear. Contrast and post-velar fronting in Russian. *Natural Language and Linguistic Theory*. ROA-504.
- Pater, Joe. to appear. Bridging the gap between receptive and productive development with minimally violable constraints. In: *Fixing priorities: Constraints in phonological development*, ed. by René Kager, Joe Pater and Wim Zonneveld. Cambridge: Cambridge University Press. ROA-296.
- Prince, Alan. 1990. Quantitative consequences of rhythmic organization. In: *Parasession on the Syllable in Phonetics and Phonology*, ed. by M. Ziolkowski, M. Noske and K. Deaton, 355-98. Chicago: Chicago Linguistic Society.
- Prince, Alan and Smolensky, Paul. 1993. *Optimality theory: constraint interaction in generative grammar*. Report no. RuCCS-TR-2. Piscataway, NJ: Rutgers Center for Cognitive Science. ROA-537.
- Prince, Alan and Tesar, Bruce. to appear. Learning phonotactic distributions. In: *Fixing priorities: Constraints in phonological acquisition*, ed. by René Kager and Joe Pater. Cambridge: Cambridge University Press. ROA-353.
- Smith, Jennifer. 2000. Positional faithfulness and learnability in Optimality Theory. In: *Proceedings from the Eastern States Conference on Linguistics*, ed. by Rebecca Daly and Anastasia Riehl, 203-14. Ithaca, NY: CLC.
- Smolensky, Paul. 1996. On the comprehension/production dilemma in child language. *Linguistic Inquiry*, 27.720-31. ROA-118.
- Tesar, Bruce. 1995. *Computational Optimality Theory*. University of Colorado, Doctoral dissertation. ROA-90.
- . 1997. An iterative strategy for learning metrical stress in Optimality Theory. In: *Proceedings of the twenty first annual Boston University Conference on Language Development*, ed. by Elizabeth Hughes, Mary Hughes and Annabel Greenhill, 615–26. Somerville, MA: Cascadilla Press.
- . 2001. Using inconsistency detection to overcome structural ambiguity in language learning. Report no. RuCCS-TR-58. Piscataway, NJ: Rutgers Center for Cognitive Science. ROA-426.
- . to appear. Enforcing grammatical restrictiveness can help resolve structural ambiguity. In: *Proceedings of the twenty first West Coast Conference on Formal Linguistics*, 443-56. University of California, Santa Cruz.
- Tesar, Bruce and Smolensky, Paul. 2000. *Learnability in Optimality Theory*. Cambridge, MA: MIT Press. ROA-155.