

Learning hidden structure in paradigms*

- Speakers have a rich and detailed knowledge of their lexicon, which they evidence in their treatment of novel words (“wug-testing”). I will show that this knowledge is biased by *naturalness*: The same kinds of relations that cause regular processes in some languages, regulate irregular processes in other languages. This means that this lexical knowledge is mediated by the grammar.
- I propose an OT-based model in which regular and irregular morpho-phonology is derived from the same set of universal constraints, CON.
- This theory requires derivations to proceed “inside-out” (Hayes 1995, 1998, 1999). It adds the benefits of OT-based work to the single surface base hypothesis (Albright 2002, 2008a).

1 The naturalness of lexical trends

1.1 Turkish (Becker, Ketrez & Nevins 2008)

Famously, Turkish final stops are predominantly voiceless. When a vowel-initial affix is added, some words keep the stop faithfully voiceless, while others alternate (Lees 1961, Zimmer & Abbott 1978, Kaisse 1986, Inkelas & Orgun 1995, Inkelas et al. 1997, Avery 1996, Kallestinova 2004, Petrova et al. 2006, among others).

(1)	bare stem	possessive	
	sop	sop-u	‘clan’
	ɟop	ɟob-u	‘nightstick’

*Ideas presented today owe much to discussions with Adam Albright, Wendell Kimper, John McCarthy, Joe Pater, and Matt Wolf. Thanks also to the audience at MUMM 2, especially Edward Flemming, John Kingston, and Donca Steriade, and the audience at the UCSC linguistics department, especially Junko Ito, Grant McGuire, Armin Mester, and Matt Tucker. I assume the responsibility for any remaining errors, here and elsewhere.

1.2 The lexicon and speakers’ knowledge of it

Given a noun like *sop*, Turkish speakers have to remember whether the possessive is *sop-u* or *sob-u*. But it helps that *sop-u* is a better guess than *sob-u*...

We searched TELL (Inkelas et al. 2000), and found that final stops in mono-syllables mostly don’t alternate, but in poly-syllables they mostly do.

(2)	Size	<i>n</i>	% alternating
	Monosyllabic, simplex coda	137	12%
	Monosyllabic, complex coda	164	26%
	Polysyllabic	2701	59%

Most final *t*’s don’t alternate, other stops mostly do.

(3)	Place	<i>n</i>	% alternating
	Labial (p)	294	84%
	Coronal (t)	1255	17%
	Palatal (tʃ)	191	61%
	Dorsal (k)	1262	85%

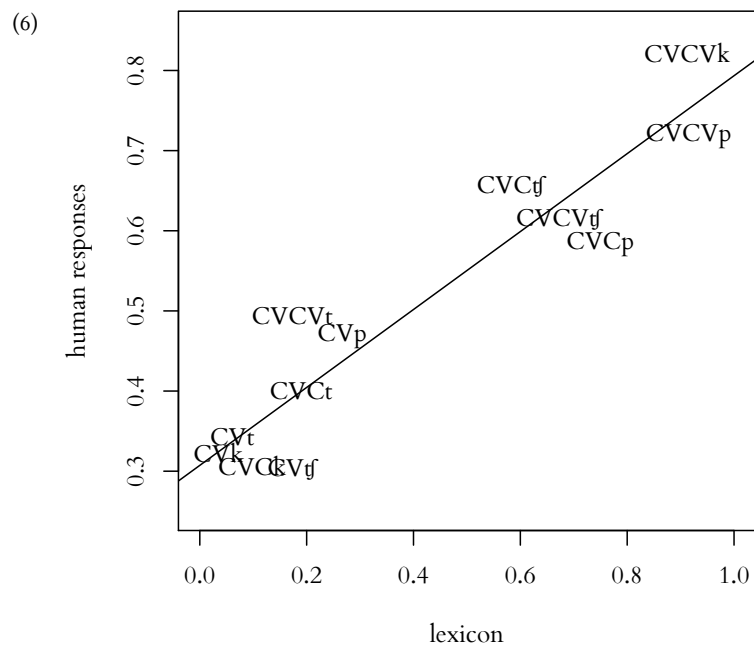
Two other factors that partially predict alternation: The height and backness of the final vowel of the stem.

(4)	Height of stem’s final vowel	<i>n</i>	% alternating
	–high	1690	42%
	+high	1312	72%

(5)	Backness of stem’s final vowel	<i>n</i>	% alternating
	–back	1495	50%
	+back	1507	60%

We gave 24 Turkish speakers a novel noun task (“wug-test”, Berko 1958) with 72 novel nouns of four places (p, t, tʃ, k), three sizes (CVC, CVCC, CVCVC), and eight vowels (a, i, e, i, o, u, ø, y).

The speakers replicated the size and place effects from the lexicon, as in (6), but not the vowel quality effects (not shown, see stats and more detail in paper).



- (7) It’s natural to treat alternations in mono-syllabic stems separately from polysyllabic stems via initial syllable faithfulness (Beckman 1997, 1998, Casali 1998).
- (8) It’s natural to treat the propensity of different stops to voice differently.
- (9) However, no language is known to change the voicing of a consonant based on the height or backness of a neighboring vowel.

In other words, Turkish speakers only learned the natural (=typologically supported) aspects of their lexicon, and ignored the unnatural ones. UG acts as a filter on the kinds of generalizations that speakers learn.²

²In Hayes et al. (to appear), unnatural trends in the data are learned, but they are attenuated relative to natural trends.

1.3 Dutch (Ernestus & Baayen 2003)

Essentially, the same as Turkish:

(10)

Imperative	Past tense	
stɔp	stɔp-tə	‘stop’
tɔp	tɔb-də	‘worry’

Speakers replicated the lexical trends for the different final obstruents (p, t, s, f, χ). They also replicated the vowel *length* effect: They preferred voicing alternations for stops that followed long vowels.

In the Dutch lexicon, there are more alternations after high vowels (which are all short) than after non-high short vowels — but speakers did not replicate this trend.

Again, this is natural: Vowel length is correlated with the voicing of a following obstruent in many languages (e.g. English), but vowel height is not.

1.4 A note on methodology

Possible objection: Your dictionary does not represent the knowledge of the people you tested, because it is old/riddled with errors/does not show morphological composition/comes from a different dialect/etc.

Responses:

- (11) The trends that speakers project from their data correlate with the lexicons we have remarkably well. Whatever the shortcomings of our dictionaries are, they are still a very good approximation of the real data in Turkish, Dutch, Hebrew, Tagalog, Hungarian, Spanish, etc.
- (12) We test the trends in the data with *sampling* (Baayen 2008), showing that these trends are strong even in lexicons that only share about 63% of their items.
- (13) It would be ideal to test each speaker twice: Once on their real lexical items, and once with novel items. This way, there is no one idealized lexicon, but rather an actual, separate lexicon for each speaker. I am working on this with Adam Albright and Andrew Nevins.

2 Analysis

2.1 Grammar-based analysis

Work “inside out” (Hayes 1995, 1998, 1999), so the alternations are considered to be irregular intervocalic voicing.

- (14) The UR’s of [sop] and [ɖɔp] are /sop/ and /ɖɔp/
- (15) The UR of the possessive is /u/ (actually just a high vowel)
- (16) /sop + u/ → [sopu] requires IDENT(voice) ≫ *VpV
/ɖɔp + u/ → [ɖɔbu] requires *VpV ≫ IDENT(voice)

Use constraint cloning (Pater 2006, 2009, Coetzee 2008, Becker 2009), which relies on the Recursive Constraint Demotion algorithm (RCD, Tesar & Smolensky 1998, 2000, Tesar 1998, Prince 2002), to detect inconsistent rankings.

- (17) IDENT(voice)_{sop} ≫ *VpV ≫ IDENT(voice)_{ɖɔp}

From this point on, every word that is sensitive to the ranking of IDENT(voice) relative to *VpV will be listed:

(18)

/top + u/	IDENT(voice)	*VpV
a. [☞] top-u		*
b. tob-u	*!	

(19)

/ot + u/	IDENT(voice)	*VpV
a. [☞] ot-u		
b. od-u	*	

- (20) IDENT(voice)_{sop, top, alp, ...} ≫ *VpV ≫ IDENT(voice)_{ɖɔp, harp, ...}

Until the speaker gets:

- (21) IDENT(voice)_{22 items} ≫ *VpV ≫ IDENT(voice)_{8 items}

Novel p-final mono-syllables will have a 8/30 (=27%) chance of alternating with [b].

The result: the lexical statistics are built into the grammar. In other words, the distinction between grammar and lexicon is blurred, so that partially-predictable information is not buried in the lexicon.

2.2 What’s wrong with a UR-based analysis?

The classic generative analysis of Turkish (Lees 1961, Inkelas & Orgun 1995, Inkelas et al. 1997, Petrova et al. 2006, among others):

- (22) The UR’s of [sop] and [ɖɔp] are /sop/ and /ɖɔB/
- (23) The UR of the possessive is /u/ (actually just a high vowel)
- (24) /sop + u/ → [sopu] requires IDENT(voice) ≫ *VpV

sop + u	IDENT(voice)	*VpV
a. [☞] sopu		*
b. sobu	*!	

- (25) /ɖɔB + u/ → [ɖɔbu] is consistent with IDENT(voice) ≫ *VpV

ɖɔB + u	IDENT(voice)	*VpV
a. ɖɔpu	(*)	*!
b. [☞] ɖɔbu	(*)	

The grammar is consistent: IDENT(voice) ≫ *VpV

The problem: The learner has no way to encode the relative numbers of /p/’s and /B/’s in the grammar. Going directly to the lexicon to find them there, unhindered by UG, will find the vowel quality generalizations that speakers don’t have.

Conclusion: Assume the bases as UR’s, assume that affixes only have segments in them, and try to get everything else by ranking constraints. Clone constraints as necessary.

3 Fallback: When the grammar is not enough

Korean (Albright 2008b):

(26)

Unmarked	Accusative		
naŋ	nasil	‘sickle’	375
naŋ	naŋ ^h il	‘face’	160
naŋ	naŋ ^h il	‘piece’	113
naŋ	nadʒil	‘daytime’	17
naŋ	nadil	‘grain’	1

Assuming /naŋ/ for the roots and /il/ for the accusative can do some work:

(27)

	/naŋ + il/	*VtV	IDENT(voice)	IDENT(asp)
a.	natil	*!		
b.	nadil		*!	
c.	naŋ ^h il			*

(28) /naŋ + il/ → [naŋ^hil], [naŋ^hil]
requires *VtV ≫ IDENT(voice) ≫ IDENT(asp)

(29) /naŋ + il/ → [nadil], [nadʒil]
requires *VtV ≫ IDENT(asp) ≫ IDENT(voice)

(30) IDENT(voice)_{113+160 items} ≫ IDENT(asp) ≫ IDENT(voice)_{1+17 items}

The prediction for a novel form, [paŋ]:

(31) 94% chance of [t^h], [tʰ], 6% chance of [d], [dʒ]

*TI, which wants assibilation before a high vowel (Kim 2001), takes care of [s]:

(32) /naŋ + il/ → [nasil]
requires *TI ≫ IDENT(cont)

(33) /naŋ + il/ → [naŋ^hil], [naŋ^hil], [nadil], [nadʒil]
requires IDENT(cont) ≫ *TI

(34) IDENT(cont)_{113+160+1+17 items} ≫ *TI ≫ IDENT(cont)_{375 items}

The prediction for a novel form, [paŋ]:

(35) 56% chance of [s], 44% chance of [t^h], [tʰ], [d], [dʒ]

But are there plausible constraints that will map /naŋ + il/ to [nadʒil] or [naŋ^hil]? It seems awfully hard to palatalize without a front vowel around.

With [naŋ^hil] as the intended winner, [naŋ^hil] is most faithful to it, but still incurs an IDENT(ant) violation → add the missing feature as floating in the UR of the accusative affix: /[-ant] il/.

(36) /naŋ + [-ant] il/ → [naŋ^hil], [nadʒil]
requires MAX(float) ≫ IDENT(ant)

(37) /naŋ + [-ant] il/ → [naŋ^hil], [nadil]
requires IDENT(ant) ≫ MAX(float)

(38) /naŋ + [-ant] il/ → [nasil]
requires *f ≫ IDENT(ant), MAX(float)

(39) *f ≫ IDENT(ant)_{113+1 items} ≫ MAX(float) ≫ IDENT(ant)_{160+17 items}

The prediction for a novel form, [paŋ]:

(40) 61% chance of [tʰ],[dʒ], 39% chance of [t^h], [d]

Summary of the predictions:

(41)

	IDENT(cont) vs. *TI	IDENT(voice) vs. IDENT(asp)	IDENT(ant) vs. MAX(float)	
[s]	56%			= 56%
[tʰ]		94%	61%	= 25%
[t ^h]	44%		39%	= 16%
[dʒ]		6%	61%	= 2%
[d]			39%	= 1%

The high probability of [s] and [ʃ^h] conforms with the report in Albright (2008b) about the treatment of novel forms, loanwords, and many native items.

My analysis expresses the language-specific frequencies of mappings in terms of rankings of universal constraints.

4 Last resort: Suppletion and diacritics

It's certainly not the case that every paradigmatic relation can be derived with phonological mechanisms, e.g. English go ~ went.

English *ɔt*-takers: *teach, catch, think, bring, seek, fight, buy* — how many of those can map to their past tense using phonological mechanisms?

The rhymes of [brɪŋ] and [baɪ] don't share any features with [ɔt] beyond [consonantal]. If we assume a floating pair of segments, /ɔt/, they can dock correctly and replace the root segments.

(42)

	baɪ + {d, ɔt}	MAX(float)	MAX(root)
a.	ɪ ³⁰ bɔt		**
b.	bat	*	*
c.	baɪ	**	
d.	bard		

Cloning MAX(float) or MAX(root) will give a small probability to *ɔt*-taking, but will say nothing about the possible shapes of *ɔt*-takers.

The fact that the regular [bard] harmonically bounds the intended winner is also a hint that something non-phonological is going on, prompting the speaker to assume suppletion or some phonology-free diacritic.

Either cloning MAX(float) or using diacritics is equally bad for finding out what kind of roots are *ɔt*-takers, and indeed speakers have no clue about *ɔt*-taking.

5 Conclusions

Speakers learn statistical trends in their lexicon, and they do so in terms of UG.

Now we have two ways of studying UG: Study regular phonology typologically, and study irregular morpho-phonology in individual languages.

To make sure that the grammar gets to see lexical statistics, don't bury them in the lexicon, and work "inside-out":

- Assume the paradigm's base as the UR, derive the other forms from it.
- Assume that affixes only have segments in them, and try to get the rest from constraint interactions. Clone constraints as necessary.
- If no grammar can be found, assume that missing structure is floating in the UR's of affixes, and try to get the rest from the grammar.
- If everything else fails, assume suppletion and/or diacritics.

This approach learns lexical trends and projects them onto novel words using an Optimality Theoretic grammar.

References

- Albright, Adam (2002). *The Identification of Bases in Morphological Paradigms*. Ph.D. dissertation, UCLA.
- Albright, Adam (2008a). A Restricted Model of UR Discovery: Evidence from Lakhota. Ms. MIT.
- Albright, Adam (2008b). Explaining universal tendencies and language particulars in analogical change. In Jeff Good (ed.) *Language Universals and Language Change*, Oxford University Press. 36 pp.
- Apoussidou, Diana (2007). *The Learnability of Metrical Phonology*. Ph.D. dissertation, University of Amsterdam.
- Avery, Peter (1996). *The Representation of Voicing Contrasts*. Ph.D. dissertation, University of Toronto.
- Baayen, R. Harald (2008). *Analyzing Linguistic Data: A practical introduction to statistics*. Cambridge University Press.
- Becker, Michael (2009). *Phonological Trends in the Lexicon: The Role of Constraints*. Ph.D. dissertation, UMass Amherst.
- Becker, Michael, Nihan Ketrez & Andrew Nevins (2008). The surfeit of the stimulus: Analytic biases filter lexical statistics in Turkish devoicing neutralization. ROA-1001.
- Beckman, Jill (1997). Positional faithfulness, positional neutralisation and Shona vowel harmony. *Phonology* 14. 1–46.

- Beckman, Jill (1998). *Positional Faithfulness*. Ph.D. dissertation, UMass, Amherst.
- Berko, Jean (1958). The child's learning of English morphology. *Word* 14. 150–177.
- Boersma, Paul (2001). Phonology-semantics interaction in ot, and its acquisition. In Robert Kirchner, Wolf Wikeley & Joe Pater (eds.) *Papers in Experimental and Theoretical Linguistics*, University of Alberta, vol. 6. 24–35.
- Casali, Roderic (1998). *Resolving Hiatus*. Garland, New York.
- Coetzee, Andries W. (2008). Grammaticality and ungrammaticality in phonology. *Language* 84. 218–257.
- Ernestus, Miriam & R. Harald Baayen (2003). Predicting the Unpredictable: Interpreting Neutralized Segments in Dutch. *Language* 79. 5–38.
- Hayes, Bruce (1995). On what to teach the undergraduates: Some changing orthodoxies in phonological theory. In Ik-Hwan Lee (ed.) *Linguistics in the Morning Calm* 3, Seoul: Hanshin. 59–77.
- Hayes, Bruce (1998). On the richness of paradigms, and the insufficiency of underlying representations in accounting for them. Handout for talk at Stanford.
- Hayes, Bruce (1999). Phonological Restructuring in Yidip and its Theoretical Consequences. In Ben Hermans & Marc van Oostendorp (eds.) *The derivational residue in phonology*, Amsterdam: Benjamins. 175–205.
- Hayes, Bruce, Kie Zuraw, Péter Siptár & Zsuzsa Londe (to appear). Natural and unnatural constraints in hungarian vowel harmony. *Language*.
- Inkelas, Sharon, Aylin Kuntay, John Lowe, Orhan Orgun & Ronald Sprouse (2000). Turkish Electronic Living Lexicon (TELL). Website, <http://socrates.berkeley.edu:7037/>.
- Inkelas, Sharon & Cemil Orhan Orgun (1995). Level ordering and economy in the lexical phonology of Turkish. *Language* 71. 763–793.
- Inkelas, Sharon, Cemil Orhan Orgun & Cheryl Zoll (1997). The implications of lexical exceptions for the nature of the grammar. In Iggy Roca (ed.) *Derivations and Constraints in Phonology*, Oxford: Clarendon. 393–418.
- Jarosz, Gaja (2006). *Rich Lexicons and Restrictive Grammars - Maximum Likelihood Learning in Optimality Theory*. Ph.D. dissertation, Johns Hopkins University.
- Kaisse, Ellen (1986). Locating Turkish devoicing. In M. Dalrymple et al. (ed.) *WCCFL* 5, Stanford Linguistics Association. 119–128.
- Kallestinova, Elena (2004). Voice and aspiration of stops in Turkish. In Grzegorz Dogil (ed.) *Folia Linguistica* 38: *Special Issue on Voice*, Berlin: Mouton de Gruyter. 117–143.
- Kenstowicz, Michael & Charles Kisseberth (1979). *Generative Phonology: Description and Theory*. Academic Press, New York.
- Kim, Hyunsoon (2001). A phonetically based account of phonological stop assibilation. *Phonology* 18. 81–108.
- Lees, Robert (1961). *The Phonology of Modern Standard Turkish*. Bloomington: Indiana University Press.
- Merchant, Nazarré (2008). *Discovering Underlying Forms: Contrast Pairs and Ranking*. Ph.D. dissertation, Rutgers University.
- Pater, Joe (2006). The locus of exceptionality: Morpheme-specific phonology as constraint indexation. In Leah Bateman & Adam Werle (eds.) *UMOP: Papers in Optimality Theory III*, Amherst, MA: GLSA. 1–36.
- Pater, Joe (2009). Morpheme-specific phonology: Constraint indexation and inconsistency resolution. In Steve Parker (ed.) *Phonological Argumentation: Essays on Evidence and Motivation*, Equinox. 1–33.
- Petrova, Olga, Rosemary Plapp, Catherine Ringen & Szilárd Szentgyörgyi (2006). Voice and aspiration: Evidence from Russian, Hungarian, German, Swedish, and Turkish. *The Linguistic Review* 23. 1–35.
- Prince, Alan (2002). Arguing optimality. ROA 562-1102.
- Tesar, Bruce (1998). Using the mutual inconsistency of structural descriptions to overcome ambiguity in language learning. In P. Tamanji & K. Kusumoto (eds.) *Proceedings of NELS* 28. Amherst, MA: GLSA, 469–483.
- Tesar, Bruce (2006). Learning from paradigmatic information. In *Proceedings of NELS* 36.
- Tesar, Bruce & Alan Prince (2006). Using phonotactics to learn phonological alternations. In *CLS* 39. Available as ROA-620.
- Tesar, Bruce & Paul Smolensky (1998). Learnability in Optimality Theory. *Linguistic Inquiry* 29. 229–268.
- Tesar, Bruce & Paul Smolensky (2000). *Learnability in Optimality Theory*. Cambridge, MA: MIT Press.
- Zimmer, Karl & Barbara Abbott (1978). The k/∅ alternation in Turkish: Some experimental evidence for its productivity. *Journal of Psycholinguistic Research* 7. 35–46.