

# Learning Opaque Generalizations: The Case of Samala (Chumash)

Jeffrey Heinz\*    William Idsardi\*\*

\*University of Delaware

\*\*University of Maryland

LSA 84th Annual Meeting  
Baltimore, MD  
January 9, 2010

# Learning opaque generalizations in phonology

1. How can phonological generalizations be automatically discovered from surface forms when they are obscured by others?
2. Discuss 2 different UG-based proposals which shuffle the data in principled ways to reveal obscured generalization
3. Case Study: Samala (Chumash) (Applegate 1972, 2007), simplified into a phonotactic learning problem
  - Correct misconceptions about the phonology of Samala
  - Study interaction between long-distance and local processes

# Learning opaque generalizations in phonology

1. How can phonological generalizations be automatically discovered from surface forms when they are obscured by others?
2. Discuss 2 different UG-based proposals which shuffle the data in principled ways to reveal obscured generalization
3. Case Study: Samala (Chumash) (Applegate 1972, 2007), simplified into a phonotactic learning problem
  - Correct misconceptions about the phonology of Samala
  - Study interaction between long-distance and local processes

# Learning opaque generalizations in phonology

1. How can phonological generalizations be automatically discovered from surface forms when they are obscured by others?
2. Discuss 2 different UG-based proposals which shuffle the data in principled ways to reveal obscured generalization
3. Case Study: Samala (Chumash) (Applegate 1972, 2007), simplified into a phonotactic learning problem
  - Correct misconceptions about the phonology of Samala
  - Study interaction between long-distance and local processes

# Samala (Inezeño Chumash)



Maria Solares (1842-1923)



John Peabody Harrington (1884-1961)



Dr. Richard Applegate

[www.chumashlanguage.com](http://www.chumashlanguage.com)

# The Corpus

- 4800 words drawn from Applegate 2007, generously provided in electronic form by Applegate (p.c).

## 35 Consonants

	labial	coronal	a.palatal	velar	uvular	glottal
stop	p p <sup>ʔ</sup> p <sup>h</sup>	t t <sup>ʔ</sup> t <sup>h</sup>		k k <sup>ʔ</sup> k <sup>h</sup>	q q <sup>ʔ</sup> q <sup>h</sup>	ʔ
affricates		ts ts <sup>ʔ</sup> ts <sup>h</sup>	tʃ tʃ <sup>ʔ</sup> tʃ <sup>h</sup>			
fricatives		s s <sup>ʔ</sup> s <sup>h</sup>	ʃ ʃ <sup>ʔ</sup> ʃ <sup>h</sup>	x x <sup>ʔ</sup>		h
nasal	m	n n <sup>ʔ</sup>				
lateral		l l <sup>ʔ</sup>				
approx.	w	y				

## 6 Vowels

i	ɨ	u
e		o
a		

(Applegate 1972, 2007)

## Opaque generalizations in Samala

Consider these processes in Samala (Applegate 1972):

1. **Local Assimilation:** [s] becomes [ʃ] before adjacent coronals [t,l,n] only across morpheme boundaries
2. **Sibilant Harmony:** the rightmost sibilant causes sibilants to the left to agree in anteriority

/s-ti-jep-us/ ‘3s tells 3s’

**Local Assimilation**

predicts [ʃtijejus]

**Sibilant Harmony**

predicts [stijejus]



/s-ti-jep-us/ ‘3s tells 3s’

**Local Assimilation**

predicts [ʃtijejus]

which is evidence against  
sibilant harmony!

**Sibilant Harmony**

predicts [stijejus]

/s-ti-jep-us/ ‘3s tells 3s’

**Local Assimilation**

predicts [ʃtijejus]

which is evidence against  
sibilant harmony!

**Sibilant Harmony**

predicts [stijejus]

which is evidence against  
local assimilation!

## The facts of Samala

### **Local Assimilation**

predicts [ʃtɨjɛpus]

### **Sibilant Harmony**

predicts [stɨjɛpus]

**/s-ti-jep-us/ → stɨjɛpus**

(Applegate 1972, 2007; texts at [www.chumashlanguage.com](http://www.chumashlanguage.com))

Contra much of the secondary phonological literature!

(Poser, 1982, 1993; Hansson, 2001; McCarthy, 2007)

## The misreading

- Applegate (1972:119-120) states that the harmony process has some exceptions, such as when the local process can apply and gives /s-ti-jep-us/ → [ʃtijejus] as an example.
- BUT Applegate meant these were *token* exceptions, not *type* ones. (Applegate p.c.)
- Applegate estimates 95% of the forms like /s-ti-jep-us/ are pronounced like [stijejus] in Harringtons copious notes of Samala (p.c).

## The misreading

- Applegate (1972:119-120) states that the harmony process has some exceptions, such as when the local process can apply and gives /s-ti-jep-us/ → [ʃtijepus] as an example.
- BUT Applegate meant these were *token* exceptions, not *type* ones. (Applegate p.c.)
- Applegate estimates 95% of the forms like /s-ti-jep-us/ are pronounced like [stijepus] in Harringtons copious notes of Samala (p.c).

### Conclusions:

1. The canonical pronunciation is [stijepus].
2. Sibilant Harmony has priority over Local Assimilation.

## Which process has priority is learned

- In Canadian French (Poliquin, 2006), pre-fricative tensing has priority over [ATR] harmony.
- Also, in Shimakonde, two harmony processes interact opaquely (Ettlinger, Bradlow and Wong 2010).
- There is no principle of UG which requires harmony patterns to have greater priority; which generalization obscures the other must be learned.

# The Problem

- Given [stijepus] ‘3s tells 3s’, how do we conclude \*st is active in the language?
- How can generalizations be learned in the face of regular exceptions?

# The Problem

- Given [stijepus] ‘3s tells 3s’, how do we conclude \*st is active in the language?
- How can generalizations be learned in the face of regular exceptions?



## The Problem

- Given [stijepus] ‘3s tells 3s’, how do we conclude \*st is active in the language?
- How can generalizations be learned in the face of regular exceptions?

x	p	t	k	q	$X \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728

**Table:** Counts of s-stop pairs in the corpus (collapsing laryngeal distinctions)

# Translating Samala into a phonotactic learning problem

## Local Assimilation

\*s[+coronal]

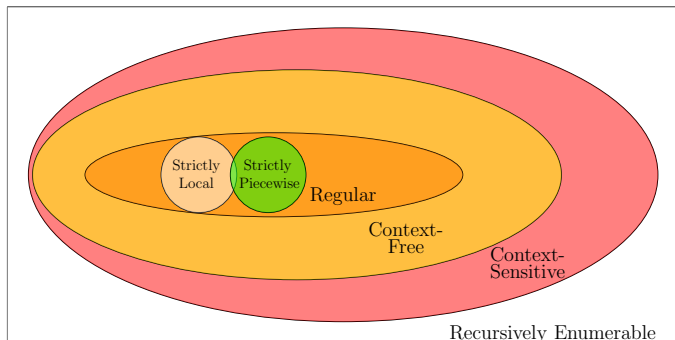
abbreviated \*st

## Sibilant Harmony

\*  $\left[ \begin{array}{c} +\text{strident} \\ \alpha\text{anterior} \end{array} \right] \cdots \left[ \begin{array}{c} +\text{strident} \\ -\alpha\text{anterior} \end{array} \right]$

abbreviated \*s...ʃ

# Learning local and long-distance phonotactic constraints



- Strictly 2-Local (SL) grammars describe constraints like  $*st$
- Strictly 2-Piecewise (SP) grammars describe constraints like  $*s...f$
- SL-k and SP-k constraints are provably efficiently learnable from distribution-free, positive evidence
- SL-k and SP-k distributions are provably efficiently estimable

(McNaughton and Papert 1971, Rogers and Pullum 2007, Heinz 2007, Rogers et. al to appear, Garcia et. al 1991, Jurafsky and Martin 2008, Heinz and Rogers in prep, Vidal et. al 2005a,b)

## Strictly Local and Strictly Piecewise

Strictly 2-Local (e.g. $*st$ )	Strictly 2-Piecewise (e.g. $*s\dots f$ )
Contiguous subsequences	Subsequences (discontiguous OK)
Immediate Predecessor	Predecessor
Concatenation ( $\cdot$ )	Less than ( $<$ )
0 = have not just seen an [a]	0 = have never seen an [a]
1 = have just seen an [a]	1 = have seen an [a] earlier

(McNaughton and Papert 1971, Simon 1975, Rogers and Pullum 2007, Rogers et. al. 2009, Heinz and Rogers in prep)

## The Estimation of $SL_2$ Distributions (bigram model)

x	p	t	k	q	$x \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728

**Table:** Counts of s-stop pairs in the corpus (collapsing laryngeal distinctions)

(Garcia et. al 1991, Jurafsky and Martin 2008)

## The Estimation of $SL_2$ Distributions (bigram model)

x	p	t	k	q	$x \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728
Counts(x)	1333	1679	1373	1130	28029

**Table:** Counts of s-stop pairs in the corpus (collapsing laryngeal distinctions)

(Garcia et. al 1991, Jurafsky and Martin 2008)

## The Estimation of $SL_2$ Distributions (bigram model)

x	p	t	k	q	$x \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728
Counts(x)	1333	1679	1373	1130	28029

**Table:** Counts of s-stop pairs in the corpus (collapsing laryngeal distinctions)

Chi-squared test not significant,  $p=0.264$

(Garcia et. al 1991, Jurafsky and Martin 2008)

## The Estimation of SP<sub>2</sub> Distributions

$P(x   b <)$		x			
		s	$\widehat{ts}$	ʃ	$\widehat{tʃ}$
b	s	0.0325	0.0051	0.0013	0.0002
	$\widehat{ts}$	0.0212	0.0114	0.0008	0.
	ʃ	0.0011	0.	0.067	0.0359
	$\widehat{tʃ}$	0.0006	0.	0.0458	0.0314

**Table:** SP<sub>2</sub> probabilities of sibilant occurring sometime after another one (collapsing laryngeal distinctions)

(Rogers et. al to appear, Heinz and Rogers in prep)



## Proposal #1

Remove data points confounded by the obscuring generalization and re-estimate

- Since Sibilant Harmony has priority over Local Assimilation, we'd like to remove words with sibilant harmony since they lead us to overestimate *st*.
1. Identify the obscuring generalization through correlation
  2. Remove all data points which conform to the obscuring generalization
  3. Re-estimate

## Proposal #1 (detail)

s t i j e p u s	u s t a s i n
s u m u ?	p a l u w o y o tʃ
a s p a x a n u s	n i p o w
? o x p o n u ʃ	s e
ts̃ a y a ?	m a n i s u s
q a l i w i l p i	s u s t a k u y u s
e x q e n	n i p a t
u s w a w a n u s	i ʃ o y

Table: Example words illustrating proposal #1

## Proposal #1 (detail)

	st	s...s		st	s...s
stijepus	1	1	ustasin	1	1
sumu?	0	0	paluwoyo $\widehat{tj}$	0	0
aspaxanus	0	1	nipow	0	0
?oxponuf	0	0	se	0	0
$\widehat{ts}$ aya?	0	0	manisus	0	1
qaliwilpi	0	0	sustakuyus	1	1
exqen	0	0	nipat	0	0
uswawanus	0	1	ifoy	0	0

**Table:** Example subset of words illustrating proposal #1.  
**Check for correlation.**

## Proposal #1 (detail)

	st	s...s		st	s...s
stijepus	1	1	ustasin	1	1
sumu?	0	0	paluwoyotj	0	0
aspaxanus	0	1	nipow	0	0
?oxponuf	0	0	se	0	0
tsaya?	0	0	manisus	0	1
qaliwilpi	0	0	sustakuyus	1	1
exqen	0	0	nipat	0	0
uswawanus	0	1	ifoy	0	0

**Table:** Example subset of words illustrating proposal #1.  
**Check for correlation.**

## Proposal #1 (detail)

	st	s...s		st	s...s
<del>stijepus</del>	1	1	<del>ustasin</del>	1	1
sumuʔ	0	0	paluwoyotʃ̂	0	0
<del>aspaxanus</del>	0	1	nipow	0	0
ʔoxponuf	0	0	se	0	0
ʔsayaʔ	0	0	<del>manisus</del>	0	1
qaliwilpi	0	0	<del>sustakuyus</del>	1	1
exqen	0	0	nipat	0	0
<del>uswawanus</del>	0	1	ifoy	0	0

**Table:** Example subset of words illustrating proposal #1.

**Remove s...s words.**

## Proposal #1 (detail)

	st	s...s		st	s...s
<del>stijepus</del>	1	1	<del>ustasin</del>	1	1
sumu?	0	0	paluwoyotj	0	0
<del>aspaxanus</del>	0	1	nipow	0	0
?oxponuf	0	0	se	0	0
tsaya?	0	0	<del>manisus</del>	0	1
qaliwilpi	0	0	<del>sustakuyus</del>	1	1
exqen	0	0	nipat	0	0
<del>uswawanus</del>	0	1	ifoy	0	0

**Table:** Example subset of words illustrating proposal #1.  
**Estimate SL2 again.**

## Results

- Only 14 of the 29 *st* words are in *s...s* words!
- The other 15 are within morphemes.

x	p	t	k	q	$x \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728

**Table:** Counts of s-stop pairs in the corpus (collapsing laryngeal distinctions).

## Results

- Only 14 of the 29 *st* words are in *s...s* words!
- The other 15 are within morphemes.

x	p	t	k	q	$x \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	24	15	28	16	511

**Table:** Counts of s-stop pairs in the corpus (collapsing laryngeal distinctions). **Results after removing s...s words.**



## Results

- Only 14 of the 29 *st* words are in *s...s* words!
- The other 15 are within morphemes.

x	p	t	k	q	$x \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	24	0	28	16	511

**Table:** Counts of s-stop pairs in the corpus (collapsing laryngeal distinctions). **Desired Results!**

## Summary of Proposal #1

- Check for an interaction between two different initial estimations of probability distributions and then revise.
- This procedure fails here because of another confound: morphological context.
- If we had a way of detecting this (e.g. Goldsmith 2001), it too could be subject to the above procedure.

## Proposal #2

Search for SL2 constraints via comparison to similar sounds

- Prior knowledge of where to search can provide direct evidence not only of \*st, but also of the repair (s→ʃ).
- To illustrate, compare *sx* and *ʃx* counts with a chi-squared analysis.

## Searching for \*st despite the confound

x	p	t	k	q	$X \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728
Counts(fx)	33	134	48	18	762

**Table:** Counts of s-stop and f-stop pairs in the corpus (collapsing laryngeal distinctions).

## Searching for \*st despite the confound

x	p	t	k	q	$X \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728
Counts(jx)	33	134	48	18	762
Counts(x)	1333	1679	1373	1130	28029

**Table:** Counts of s-stop and j-stop pairs in the corpus (collapsing laryngeal distinctions).

## Searching for \*st despite the confound

x	p	t	k	q	$X \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728
Counts(jx)	33	134	48	18	762
Counts(x)	1333	1679	1373	1130	28029

**Table:** Counts of s-stop and j-stop pairs in the corpus (collapsing laryngeal distinctions).

## Searching for \*st despite the confound

x	p	t	k	q	$X \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728
Counts(jx)	33	134	48	18	762
Counts(x)	1333	1679	1373	1130	28029

**Table:** Counts of s-stop and j-stop pairs in the corpus (collapsing laryngeal distinctions).

## Chi-squared Test

x	p	t	k	q	$X \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728
Counts(fx)	33	134	48	18	762
Counts(x)	1333	1679	1373	1130	28029

**Table:** Counts of s-stop and f-stop pairs in the corpus (collapsing laryngeal distinctions).

x	p	t	k	q	$x \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	0.106	-5.292	-0.318	0.616	1.706
Counts(fx)	-0.097	4.871	0.293	-0.567	-1.571

**Table:** Residuals from  $\chi^2$  test on counts of s-stop and f-stop pairs in the corpus (collapsing laryngeal distinctions).  $\chi^2 = 58.0274$ ,  $df = 4$ ,  $p\text{-value} = 7.53e-12$ .



## Chi-squared Test

x	p	t	k	q	$X \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728
Counts(fx)	33	134	48	18	762
Counts(x)	1333	1679	1373	1130	28029

**Table:** Counts of s-stop and f-stop pairs in the corpus (collapsing laryngeal distinctions).

x	p	t	k	q	$x \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	0.106	-5.292*	-0.318	0.616	1.706
Counts(fx)	-0.097	4.871*	0.293	-0.567	-1.571

**Table:** Residuals from  $\chi^2$  test on counts of s-stop and f-stop pairs in the corpus (collapsing laryngeal distinctions).  $\chi^2 = 58.0274$ ,  $df = 4$ , p-value =  $7.53e-12$ . Highlighted cells  $p < 0.05$  (critical value=3.84)

## Unigram counts of C2 are misleading

x	p	t	k	q	$X \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728
Counts(fx)	33	134	48	18	762
Counts(x)	1333	1679	1373	1130	28029

**Table:** Counts of s-stop and f-stop pairs in the corpus (collapsing laryngeal distinctions).

x	p	t	k	q
Count(fx)	-3.006	7.058*	-1.265	-4.183*
Count(x)	0.618	-1.451	0.260	0.860

**Table:** Residuals from  $\chi^2$  test on counts of fx pairs with counts of x in the corpus (collapsing laryngeal distinctions).  $\chi^2 = 1.2497$ ,  $df = 3$ , p-value  $< 2.2e-16$ . Highlighted cells  $p < 0.05$  (critical value=3.84)

## Unigram counts of C2 are misleading

x	p	t	k	q	$X \notin \Sigma - \{p, t, k, q\}$
Counts(sx)	29	29	37	20	728
Counts(fx)	33	134	48	18	762
Counts(x)	1333	1679	1373	1130	28029

**Table:** Counts of s-stop and f-stop pairs in the corpus (collapsing laryngeal distinctions).

Would we conclude that  $q \rightarrow t/f$  \_\_\_\_\_?

x	p	t	k	q
Count(fx)	-3.006	7.058*	-1.265	-4.183*
Count(x)	0.618	-1.451	0.260	0.860

**Table:** Residuals from  $\chi^2$  test on counts of fx pairs with counts of x in the corpus (collapsing laryngeal distinctions).  $\chi^2 = 1.2497$ ,  $df = 3$ , p-value  $< 2.2e-16$ . Highlighted cells  $p < 0.05$  (critical value=3.84)

## Proposal #2 Summary

- Prior knowledge guides the right comparisons to make correct inferences despite confounded data
- Generally, the idea is to compare  $ax$  sequences (SL or SP) with  $bx$  sequences where  $a$  and  $b$  are similar.

## Conclusion

1. We corrected a misreading in earlier literature

**/s-ti-jep-us/ → [stijepus], not \*[ftijepus]**

2. We identified a new well-defined learning problem and explored two different approaches
3. Correct statistical inference is possible, but only with the right model, i.e. structured probabilistic models  
(Yang 2000, Goldwater 2006, Hayes and Wilson 2008, and many others)

# Acknowledgements

Thanks to

- Dr. Richard Applegate
- 2008-2009 U. of Delaware Research Fund Grant
- National Institutes of Health #7R01DC005660-07