



# What echo reduplication reveals about phonological similarity

---

*Sameer ud Dowla Khan, Brown University*

Thursday 20 April 2012

Reed College



# Outline

---

- **Background: echo reduplication**
  - Properties of echo reduplication crosslinguistically
- **Experiment: identity and similarity**
  - Similarity avoidance in Bengali echo reduplication
- **Analysis: measurement of similarity**
  - Measurement of consonant similarity in Bengali
- **General discussion**
  - Summary and further questions



# Background

---

Properties of echo reduplication  
crosslinguistically



# Echo reduplication

---

- **Subtractive reduplication**

*Bengali*

[goli] ‘alley’

[oli goli] ‘alleys, etc.’

- **Fixed-segment ( $S_F$ ) reduplication**

*Bengali*

[kafi] ‘cough’

[kafi t<sub>F</sub>afi] ‘cough, etc.’

*English*

[kɔf] ‘cough’

[kɔf ʃm<sub>F</sub>ɔf] ‘cough<sub>DISMISSIVE</sub>’



# Echo reduplication

---

- **Subtractive reduplication**

*Bengali*

[goli] ‘alley’

[oli goli] ‘alleys, etc.’

- **Fixed-segment ( $S_F$ ) reduplication**

*Bengali*

[kafi] ‘cough’

[kafi t<sub>F</sub>afi] ‘cough, etc.’

*English*

[kɒf] ‘cough’

[kɒf ʃm<sub>F</sub>ɒf] ‘cough<sub>DISMISSIVE</sub>’



## Fixed-segment reduplication

---

- In FSR, **fixed material** ( $S_F$ ) associated with a particular construction is found in the R **instead of a copy of B material**
- The fixed material can be:
  - A consonant (most common)
  - A vowel
  - A CV sequence
  - A stem



# Fixed-segment reduplication

---

- **Consonantal  $S_F$**

*Kashmiri (Koul):*  $S_F = [v_F]$ <sup>1</sup>

[nalkɪ] ‘faucet’

[nalkɪ v<sub>F</sub>alkɪ] ‘faucet, etc.’

- **Vocalic  $S_F$**

*A-Hmao (Mortensen 2005):*  $S_F = [i_F]$

[and<sub>ɸ</sub><sup>h</sup>ǎu] ‘mouth’

[ánd<sub>ɸ</sub><sup>h</sup>i<sub>F</sub> ánd<sub>ɸ</sub><sup>h</sup>ǎu] ‘cheeks, nose, etc.’

<sup>1</sup> IPA: Koul & Wali (2006)



# Fixed-segment reduplication

---

- **[CV] S<sub>F</sub>**

*Tamil (Keane 2001): S<sub>F</sub> = [ki<sub>F</sub>]<sup>1</sup>*

[veɭ:ai] ‘white’

[veɭ:ai ki<sub>F</sub>ɭ:ai] ‘white, etc.’

- **Stem S<sub>F</sub>**

*Russian (Podobryaev 2012): S<sub>F</sub> = [xuj<sub>F</sub>] < ‘penis’*

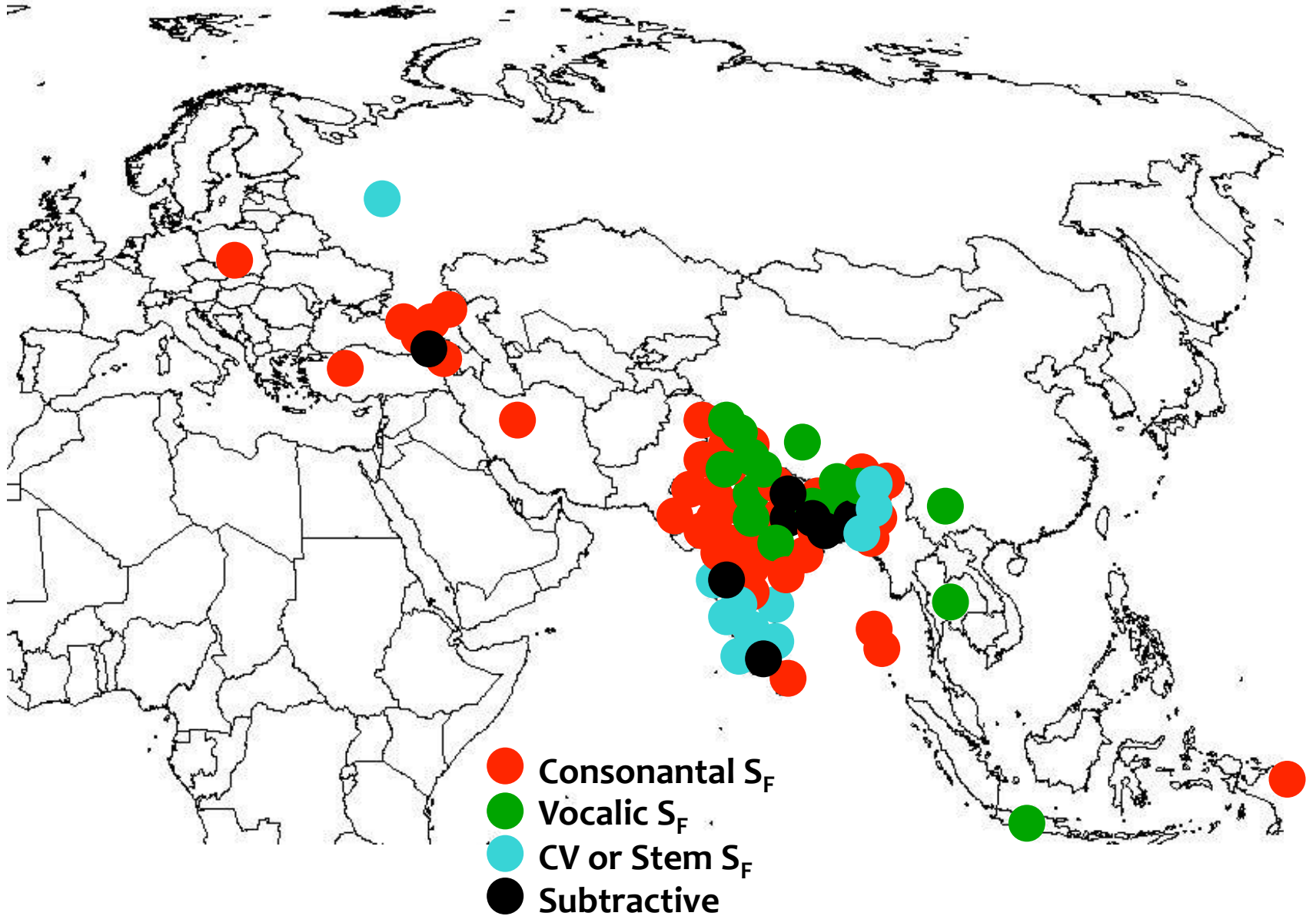
[málʲtɕik] ‘boy’

[málʲtɕik xuj<sub>F</sub>álʲtɕik] ‘boy<sub>DISMISSIVE</sub>’

<sup>1</sup>IPA: Keane (2004)



# Languages reported to have echo reduplication





# Echo reduplication

---

- Typically conveys **generalization**
  - ‘X, etc.’
  - ‘X and associated things’
  - ‘X in general’
  - ‘superset of which X is a member’
- In some lgs, it conveys a **dismissive** tone
  - Russian: [málʲtɕik xuj<sub>F</sub>álʲtɕik] ‘boy<sub>DISMISSIVE</sub>’
  - English: [daktəʁ ʃm<sub>F</sub>aktəʁ] ‘doctor<sub>DISMISSIVE</sub>’



# Obligatory BR-nonidentity

---

- Most salient **phonological property** of echo reduplication is **obligatory BR-nonidentity**
  - $\geq 1$  phonological difference between B and R
- **Presence of  $S_F$**  in R usually enough to **generate BR-nonidentity**

*Kashmiri (Koul)*

[nalki]  $\Rightarrow$  [nalki  $v_F$ alki] ‘faucet, etc.’

- **But what if it isn’ t?**

*Kashmiri (Koul)*

[va:zi]  $\Rightarrow$  ??[va:zi  $v_F$ a:zi]?? ‘cook, etc.’



## Obligatory BR-nonidentity

---

- Lgs avoid such cases of potential BR-identity by either:
  - Using a designated **backup  $S_F$**
  - Choosing from among the **other  $S_F$  options**
  - **Modifying the B** instead of in R
  - Deeming the phrase **ineffable**



# Obligatory BR-nonidentity

---

- Many lgs have a **backup S<sub>F</sub>**, kept on reserve for cases of BR-identity

*Abkhaz (Vaux 1996): S<sub>F</sub> = [m<sub>F</sub>] (⇒ [ʃ<sub>F</sub>])*

*/gáɖzak' / ⇒ [gáɖzak' m<sub>F</sub>áɖzak' ] 'fool, etc.'*

*/ʃək' / ⇒ [ʃək' m<sub>F</sub>ək' ] 'horse, etc.'*

*/maát/ ⇒ \*[maát m<sub>F</sub>aát-] ⇒ [maát ʃ<sub>F</sub>aát-] 'money, etc.'*



## Obligatory BR-nonidentity

---

- Other lgs have **multiple  $S_F$  options**, always choosing one that avoids BR-identity

*Farsi* (Ghaniabadi et al. 2006):  $S_F = [m_F] \sim [p_F]$

/tæɾɒzu/  $\Rightarrow$  [tæɾɒzu **m<sub>F</sub>**æɾɒzu]  $\sim$  [tæɾɒzu **p<sub>F</sub>**æɾɒzu] ‘scale, etc.’

/zæɾif/  $\Rightarrow$  [zæɾif **m<sub>F</sub>**æɾif]  $\sim$  [zæɾif **p<sub>F</sub>**æɾif] ‘slender, etc.’

/mive/  $\Rightarrow$  \*[**m**ive **m<sub>F</sub>**ive]  $\sim$  [**m**ive **p<sub>F</sub>**ive] ‘fruit, etc.’

/pir/  $\Rightarrow$  [**p**ir **m<sub>F</sub>**ir]  $\sim$  \*[**p**ir **p<sub>F</sub>**ir] ‘old, etc.’



# Obligatory BR-nonidentity

---

- Some lgs even go so far as to **modify B** when R with  $S_F$  would be identical to it

*Classical Tibetan* (Beyer 1992):  $S_F = [a_F]$  ( $\Rightarrow$  B  $[o_F]$ )

/ndzog/  $\Rightarrow$  [ndz $a_F$ g ndz $o$ g] ‘jumbled up’

/glen/  $\Rightarrow$  [gl $a_F$ n gl $e$ n] ‘very stupid’

/ŋan/  $\Rightarrow$  \*[ŋ $a_F$ n ŋ $a$ n]  $\Rightarrow$  [ŋ $a$ n ŋ $o_F$ n] ‘miserable’



# Obligatory BR-nonidentity

---

- Lastly, some lgs simply deem cases of echo BR-identity to be **ineffable**

*Turkish (Swift 1963):*  $S_F = [m_F]$

/havtu/<sup>1</sup> ⇒ [havtu **m<sub>F</sub>**avtu] ‘towel, etc.’

/citap/ ⇒ [citap **m<sub>F</sub>**itap] ‘book, etc.’

/masa/ ⇒ \*[masa **m<sub>F</sub>**asa] ‘table, etc.’ ⇒ NO OUTPUT

<sup>1</sup>IPA: Zimmer & Orgun (1999)





## Obligatory BR-nonidentity

---

- Crosslinguistically, **BR-identity in echo reduplication is ungrammatical**
- Trivedi's (1990) survey of FSR in ~100 Indian lgs found **obligatory BR-nonidentity in every lg**
- Seems clear... but I still have one question:  
**How sensitive is BR-nonidentity?**



# Survey

---

- For example, let's consider English

English:  $S_F = [\text{ʃm}_F]$

/daktə/ 'doctor'  $\Rightarrow$  [daktə **ʃm<sub>F</sub>**aktə]

'doctor<sub>DISMISSIVE</sub>'

/skul/ 'school'  $\Rightarrow$  ?

/smuð/ 'smooth'  $\Rightarrow$  ?

/ʃmuz/ 'schmooze'  $\Rightarrow$  ?

/ʃmalts/ 'schmaltz'  $\Rightarrow$  ?

/ʃnaz/ 'schnozz'  $\Rightarrow$  ?



## Curiosity from literature

---

- In Nevins & Vaux (2003), 95% of speakers avoided [ʃm<sub>F</sub>] in R of [ʃmuz] ‘schmooze’
  - \*[ʃmuz ʃm<sub>F</sub>uz] due to BR-nonidentity
  - Fits with cross-linguistic pattern
- Interestingly, 30% of speakers also avoided [ʃm<sub>F</sub>] in R of [ʃnaz] ‘schnozz’
  - \*[ʃnaz ʃm<sub>F</sub>az]...but why?
  - BR-nonidentity generalized to **BR-dissimilarity**



# Curiosity from literature

---

- What does this mean?
  - Explanation 1:
    - **30%** of speakers have a **BR-dissimilarity constraint**
    - **65%** have a **BR-nonidentity constraint**
  - Explanation 2:
    - **95%** have a **BR-dissimilarity constraint**, of which:
      - **30%** felt  $[\int n]$  and  $[\int m_F]$  are **too similar**
      - **65%** felt  $[\int n]$  and  $[\int m_F]$  are **sufficiently dissimilar**



## Is this English-specific?

---

- Or maybe we're assuming too much from this one data point...
- Is English echo reduplication a weird case?
  - **Not as common** as in other languages
  - [ʃm]-reduplication is somewhat **humorous**
  - [ʃm] and [ʃn] are **highly marked** in English
- Maybe this is just a weird fact of English...



## Motivation for an experiment

---

- To find out if echo reduplication involves BR-nonidentity or BR-dissimilarity...
- We need to study a lg in which:
  - Echo reduplication is a **fully productive, linguistic feature**
  - $S_F$  isn't such a marked sound



# Experiment

---

What echo reduplication reveals  
about phonological similarity



# Experiment: question

---

- Question: **how sensitive is BR-assessment?**
  - **Only sensitive** to exact **BR-identity**
    - Any BR-difference should suffice
  - **Also sensitive** to relative **BR-similarity**
    - Some BR-differences aren't dissimilar enough





# Experiment: language

---

- Test case: **Bengali** echo reduplication
  - Default  $S_F$ : [t<sub>F</sub>]
  - Backup  $S_F$ : [m<sub>F</sub>] [f<sub>F</sub>] [p<sub>F</sub>] [u<sub>F</sub>]...
- Why Bengali?
  - Echo reduplication is a very **common feature**
  - Default [t<sub>F</sub>] is a relatively **unmarked sound**
  - Many **contrastive but phonetically similar** phonemes: [t<sup>h</sup>] [d] [ɽ] [ɽ<sup>h</sup>] [tɕ]...



## Experiment: basic idea

---

- So we know that a word like [b<sup>h</sup>idz:a] ‘having gotten wet’ ⇒ [b<sup>h</sup>idz:a t<sub>F</sub>idz:a]...
- ...and that a word like [tika] ‘vaccine’ ⇒ \*[tika t<sub>F</sub>ika] ⇒ [tika m<sub>F</sub>ika]...
- ...but what about a word like [t<sup>h</sup>ajʃ:a] ‘having stuffed’ ?
  - Will it act like [b<sup>h</sup>idz:a]? [t<sup>h</sup>ajʃ:a t<sub>F</sub>ajʃ:a]
  - Or like [tika]? \*[t<sup>h</sup>ajʃ:a t<sub>F</sub>ajʃ:a] ⇒ [t<sup>h</sup>ajʃ:a m<sub>F</sub>ajʃ:a]



## Experiment: subjects and procedure

---

- **Production experiment** with native speaker adults (n=30)
- Heard audio recording of a word
  - Order was randomized for each speaker
- Asked to **produce echo reduplicated form**
  - Did speaker use default [t<sub>F</sub>]?
    - Or did he/she use a backup S<sub>F</sub>?



# Experiment: stimuli

---

- 60 test words fell under **three conditions**:
  - **Identity**: [t]-initial words
  - **Similarity**: words with [t]-like initials
    - Coronal obstruents: [t<sup>h</sup>] [d] [t̚] [t̚<sup>h</sup>] [t̥̚] [s]~[t̥̚<sup>h</sup>] [ʃ]
  - **Control**: words with **non-[t]-like initials**
    - Coronal sonorants: [n] [l] [ɹ]
    - Non-coronals: [k] [h] [p] [f] [b<sup>h</sup>] [m]

# Experiment: stimuli

Bengali consonant inventory (Khan 2010)

**Identity**

**Similarity**

**Control**

	Labial	Dental	Alveolar	Post-Alv	Velar/Glot
Stop	p b b <sup>h</sup>	t̪ t̪ <sup>h</sup> d̪ d̪ <sup>h</sup>	t t <sup>h</sup> d d <sup>h</sup>		k k <sup>h</sup> g g <sup>h</sup>
Affricate			tʃ tʃ <sup>h</sup> dʒ dʒ <sup>h</sup>		
Fricative	f	s		ʃ	h
Liquid			l ɭ		
Nasal	m		n		(ŋ)

# Experiment: stimuli

Bengali consonant inventory (Khan 2010)

**Identity**

**Similarity**

**Control**

	Labial	Dental	Alveolar	Post-Alv	Velar/Glot
Stop	p b b <sup>h</sup>	t̪ t̪ <sup>h</sup> d̪ d̪ <sup>h</sup>	t t <sup>h</sup> d d <sup>h</sup>		k k <sup>h</sup> g g <sup>h</sup>
Affricate			tʃ tʃ <sup>h</sup> dʒ dʒ <sup>h</sup>		
Fricative	f	s		ʃ	h
Liquid			l ɭ		
Nasal	m		n		(ŋ)



# Experiment: hypothesis 1

---

- Hypothesis 1: BR-assessment is **only identity-sensitive**
  - **Identity** words will **never** use  $[t_F]$
  - **Similarity** words will behave like **Control** words
  - **Control** words will **always** use  $[t_F]$
- **Identity**  $\neq$  **Similarity** = **Control**  
**\*[t...t<sub>F</sub>]**  $\neq$  **[t<sup>h</sup>...t<sub>F</sub>]** = **[b<sup>h</sup>...t<sub>F</sub>]**



## Experiment: hypothesis 2

---

- Hypothesis 2: BR-assessment is **sensitive to phonetic similarity** across phonemes
  - **Identity** words will **never** use  $[t_F]$
  - **Similarity** words will behave like **Identity** words
  - **Control** words will **always** use  $[t_F]$
- **Identity** = **Similarity**  $\neq$  **Control**  
 $*[t\dots t_F] \neq *[t^h\dots t_F] = [b^{\hat{n}}\dots t_F]$





## Experiment: hypothesis 3

---

- Hypothesis 3: BR-assessment is **strongest in cases of identity**, but also **sensitive to phonetic similarity** across phonemes
  - **Identity** words will **never** use  $[t_F]$
  - **Similarity** words will **sometimes** use  $[t_F]$
  - **Control** words will **always** use  $[t_F]$
- **Identity**  $\neq$  **Similarity**  $\neq$  **Control**  
**\*[t...t<sub>F</sub>]**  $\neq$  **?[t<sup>h</sup>...t<sub>F</sub>]** = **[b<sup>h</sup>...t<sub>F</sub>]**



# Results

---

- The results were surprising:
- While Hypothesis 3 came closest, none of the hypotheses was borne out!
  - **Identity** ≠ **Similarity** = **Control**
  - **Identity** = **Similarity** ≠ **Control**
  - **Identity** ≠ **Similarity** ≠ **Control**
- Lots of variation across speakers and words
- But one thing was clear:



## Results: general pattern

---

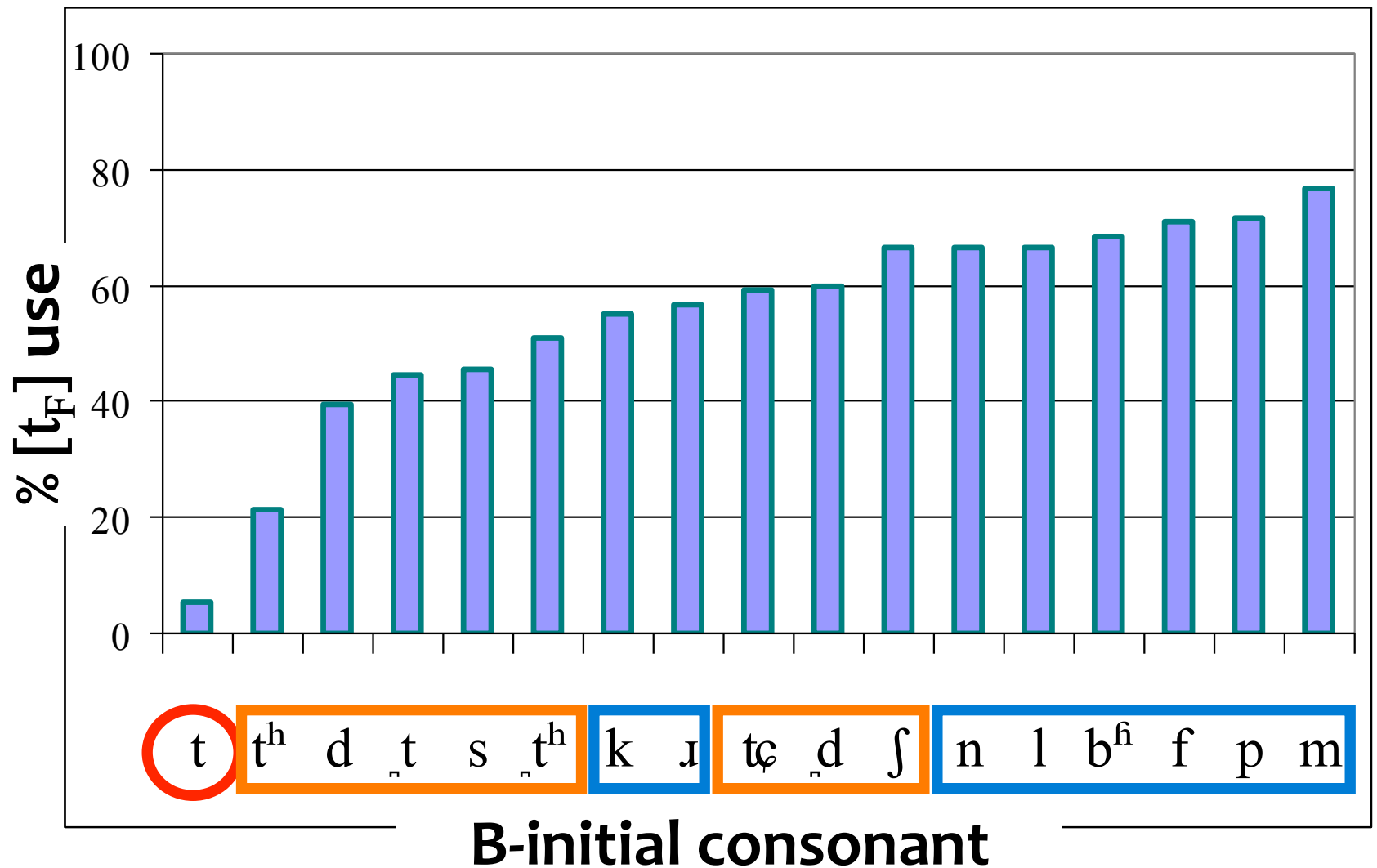
- BR-assessment is **sensitive to similarity**
- But, this sensitivity is **gradient**
- The **more similar** a consonant is to [t], the **less likely** it is to be replaced by [t<sub>F</sub>]

[t] [t<sup>h</sup>] [d] [ṭ] ... [b<sup>h</sup>] [f] [p] [m]

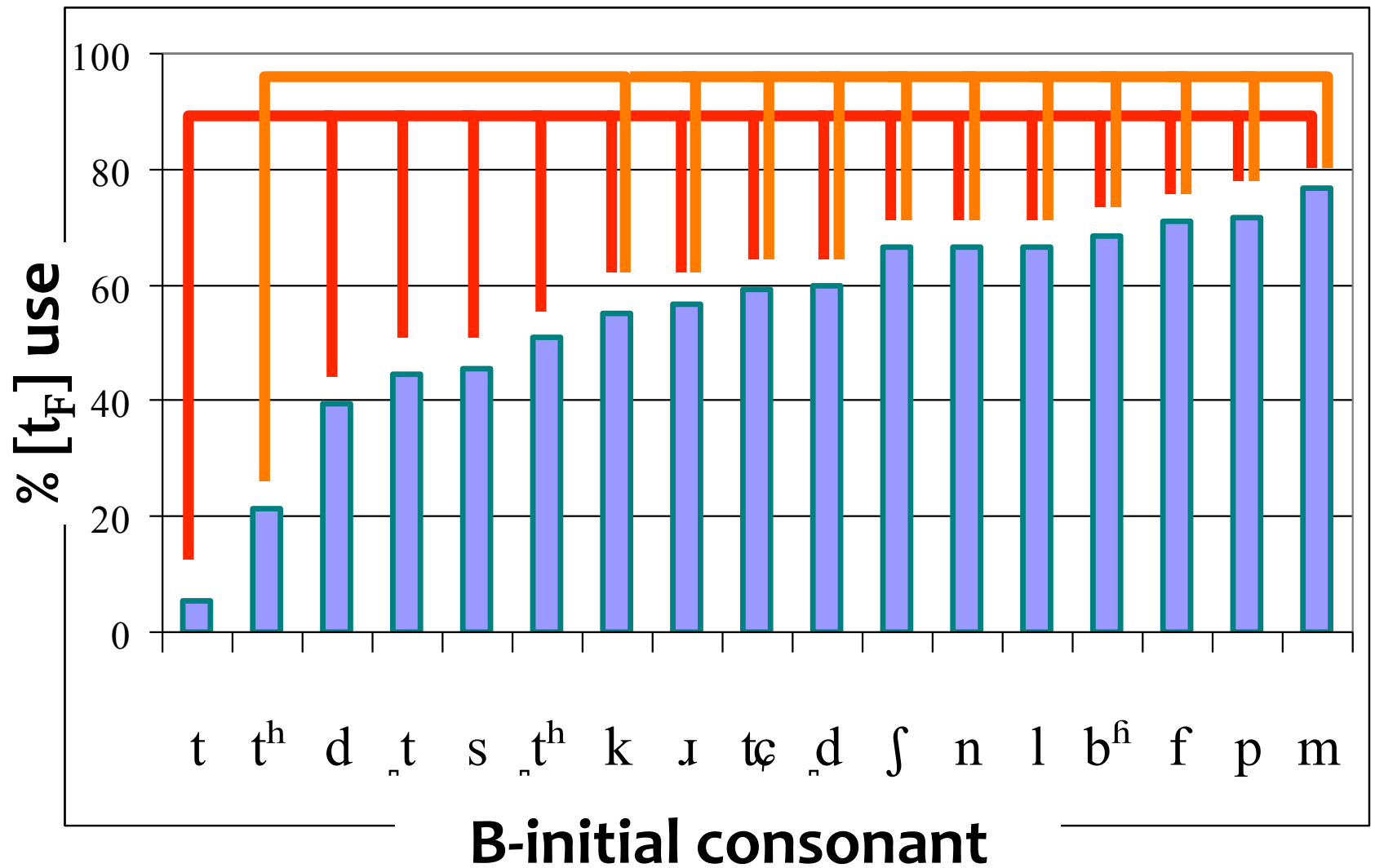
←  
Least likely  
to use [t<sub>F</sub>]

→  
Most likely  
to use [t<sub>F</sub>]

# Results: grouped by initial consonant



## Results: significant differences



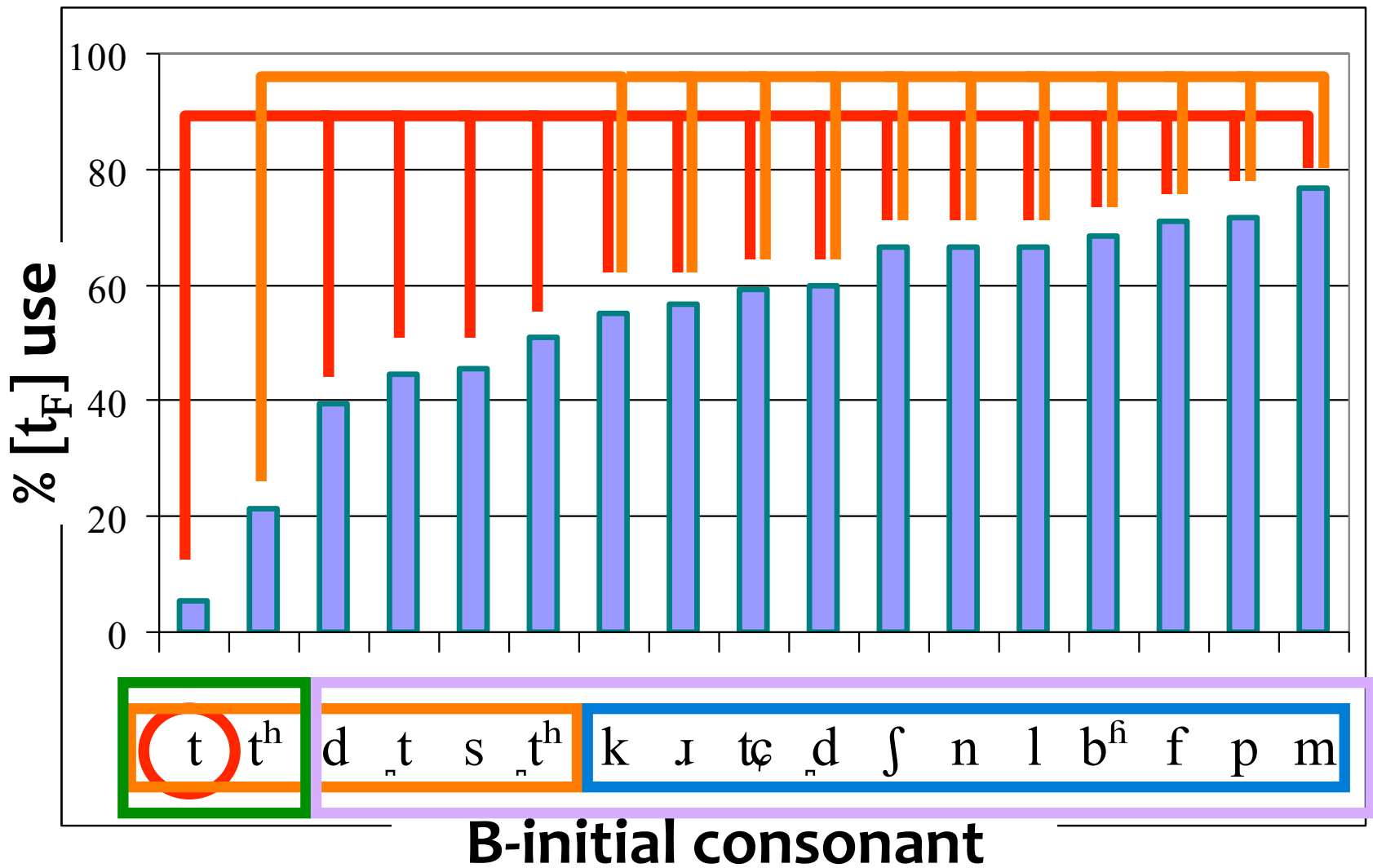


## Results: design issues?

---

- Was there a problem with the setup?
- Should the similarity condition and control condition be redefined?

# Results: clustering





## Results: gradient similarity

---

- No, there is **no clustering of consonants** into two or three categories
- Furthermore, **heavy overlap** across the clusters that are found
- Clearly, **similarity is gradient**





## New question

---

- It's clear some **measurement of similarity** is needed in Bengali echo reduplication
- So then how do speakers **calculate the similarity** of a pair of sounds?



# Analysis

---

Measurement of consonant  
similarity in Bengali



## Models of similarity

---

- Phonological similarity has been **measured in different ways** in the literature
- Most metrics incorporate:
  - **Shared natural classes**
  - Correlation with **lexical cooccurrence**
- Can either of these model Bengali speakers' notions of similarity?



# Shared natural classes: introduction

---

- One method<sup>1</sup> of calculating similarity is by comparing the **number of natural classes** of which two sounds are members
  - [t] and [n] share [+cor] but not [voi] or [son]
  - [t] and [p] share [-son] and [-voi] but not [cor]
  - [t] and [v] share [-son] but not [voi] or [cor]
- The **more similar** two consonants are, the **more natural classes** they will share

<sup>1</sup> Frisch, Pierrehumbert, & Broe (2004)



## Shared natural classes: introduction

---

- This measure takes **lg-specific details** into account, due to different inventories
  - Aspiration is contrastive in Bengali, not English
  - [b] and [d] share [-asp] in Bengali, not English
- Can we apply the SNC metric to the echo reduplication results in Bengali?



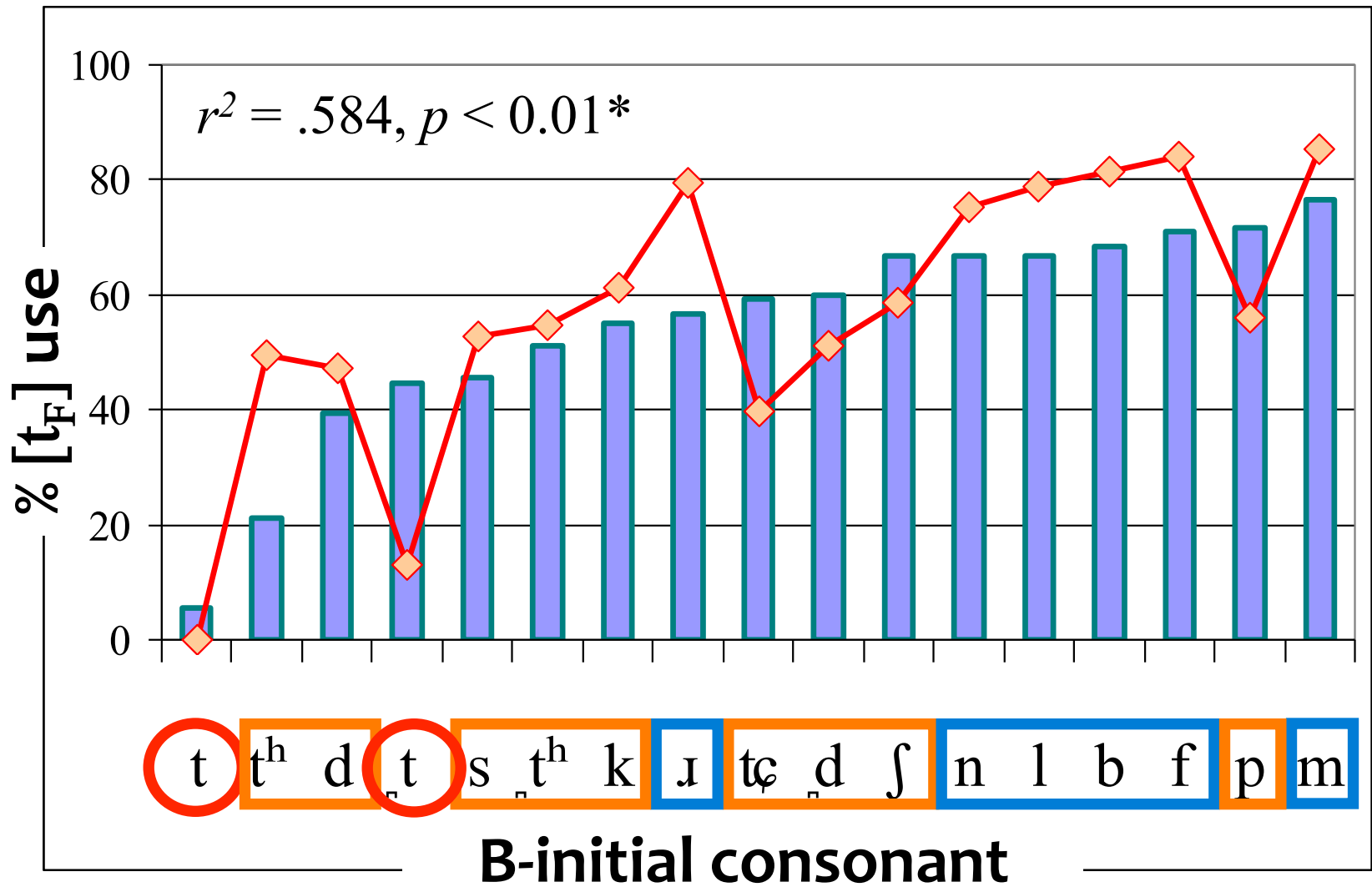
## Shared natural classes metric

---

- In a model of similarity based on **shared natural classes**...
- ...the similarity of a consonant  $C_1$  to [t] can be calculated as follows

$$sim(C_1, t) = \frac{\# \text{ shared natural classes}}{\# \text{ shared natural classes} + \# \text{ non-shared natural classes}}$$

# SNC predictions (line) vs. observed (bars)





## SNC metric: discussion

---

- SNC metric does fairly well ( $r^2 = .584$ )
- However, where it doesn't do well is the most crucial area: **coronal obstruents**
  - How can we adjust this model to reflect that [t] is more similar to [t<sup>h</sup>] than to [ṭ]?
  - Is there a way to designate **certain features as being more important** than others?





# Feature weighting

---

- What if we incorporated **different weights for different features**, reflecting their importance in similarity measurement?
  - Weighting [distributed] over [aspiration] will make ([t], [t]) more different than ([t], [t<sup>h</sup>])
- What would such a metric look like?



## Similarity equation

---

- In a model of similarity based on **shared weighted features**...
- ...the similarity of a consonant  $C_1$  to [t] can be calculated as follows

$$sim(C_1, t) = \exp\left(-\sum_{i=1}^{\#features} w_i (1 - \delta_i(C_1, t))\right)$$

$w_i$  = weight of the feature  $f_i$

$\delta_i(C_1, t) = 1$  (feature value shared) or 0 (not shared)



## Where do these weights come from?

---

- So how do we determine  $w_i$ ?
- Maybe **lexical statistics**?



## Lexical cooccurrence: introduction

---

- Many studies<sup>1</sup> claim that the **lexicon of a lg reflects notions of similarity**
- The **more similar** two consonants are, the **less often they will cooccur** within roots
  - Words like [fʌɔʒ] and [pɛg] are common
  - Words like \*[ʃʌɔʒ] and \*[pɛb] are underattested
- Can we apply this to the Bengali data?

<sup>1</sup> McCarthy (1994) and many others



## Lexical cooccurrence: metric

---

- We can turn this around:
- Sound pairs that are **underattested** within roots must be perceived as **more similar**
- Thus, **lexical statistics** can be converted into a **similarity score**



## Lexical cooccurrence: metric

---

- In a model of similarity based on **lexical cooccurrence statistics**...
- ...the **similarity of a consonant  $C_1$  to [t]** can be calculated as follows

$$sim(C_1, t) = \frac{obs[C_1VCV]}{\text{all roots}} \times \frac{obs[CVtV]}{\text{all roots}}$$

---

$$\frac{obs[C_1VtV]}{\text{all roots}}$$

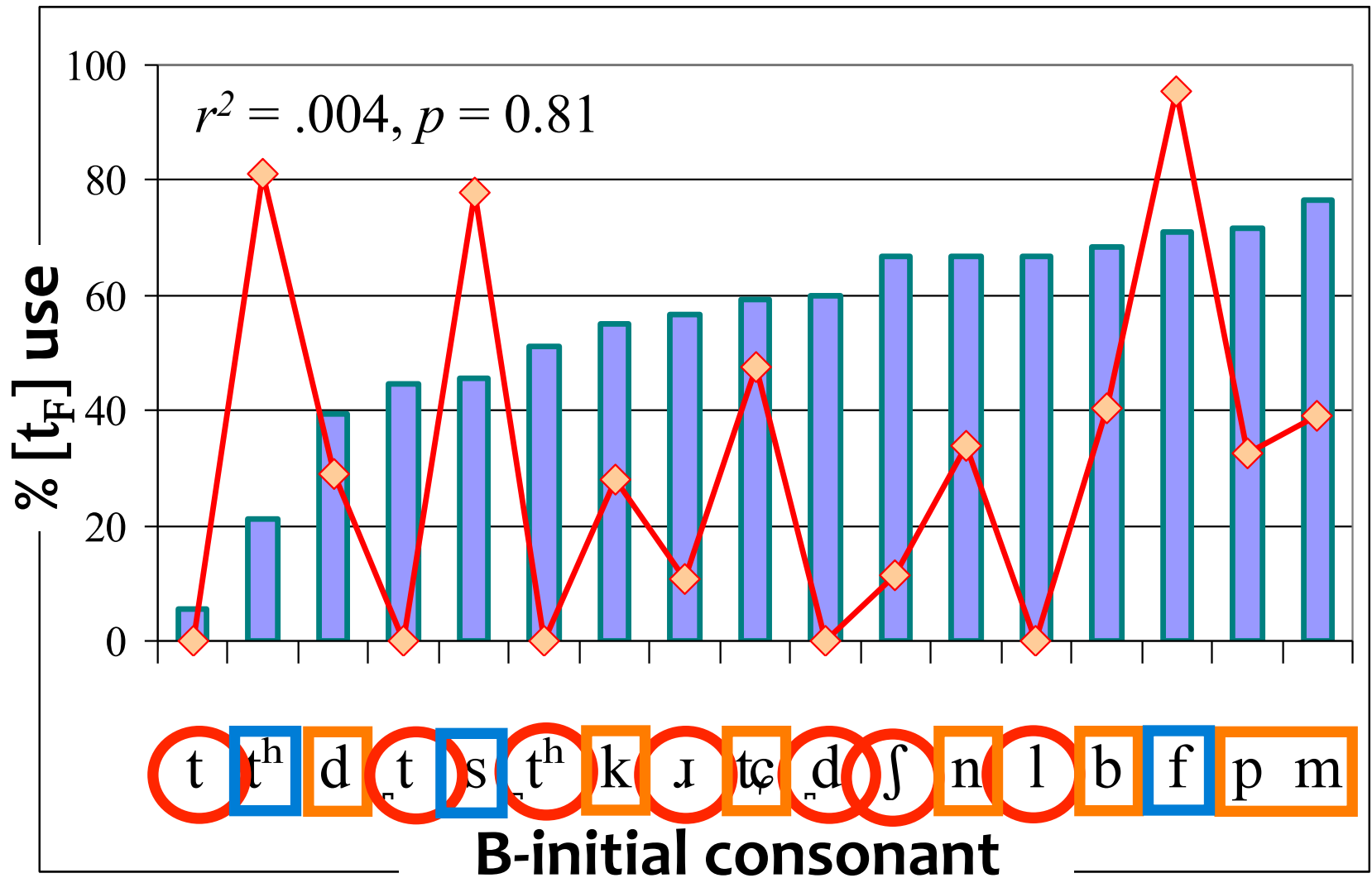


## Lexical cooccurrence: data

---

- Used a **Bengali corpus** (Mallik et al. 1998) to examine roots where **[t]** cooccurs with a **consonant (C)**
- Plugged in the numbers to get a **similarity score** for each C paired with [t]
- Compared those similarity scores to the  $[t_F]$ -use patterns from my experiment

# Lexical cooccurrence: results







## Lexical cooccurrence: discussion

---

- The lexical cooccurrence model of similarity **fails to predict** the observed  $[t_F]$ -avoidance patterns ( $r^2 = .004$ )
- It appears that the **Bengali lexicon does not reflect the notions of similarity** at work in the productive grammar
- Thus, we cannot use lexical statistics to adjust our natural classes model



## Can weights be used at all?

---

- So how else can we determine what the weights should be?
- Can weights help us at all?
- Let's see if we can use the **variation in the data itself** to determine the weights...
- ...and then worry about **where the weights are coming from** at some other time



## Probability equation

---

- Probability of  $[t_F]$ -use in the echo R of a  $C_1$ -initial B can be calculated as follows

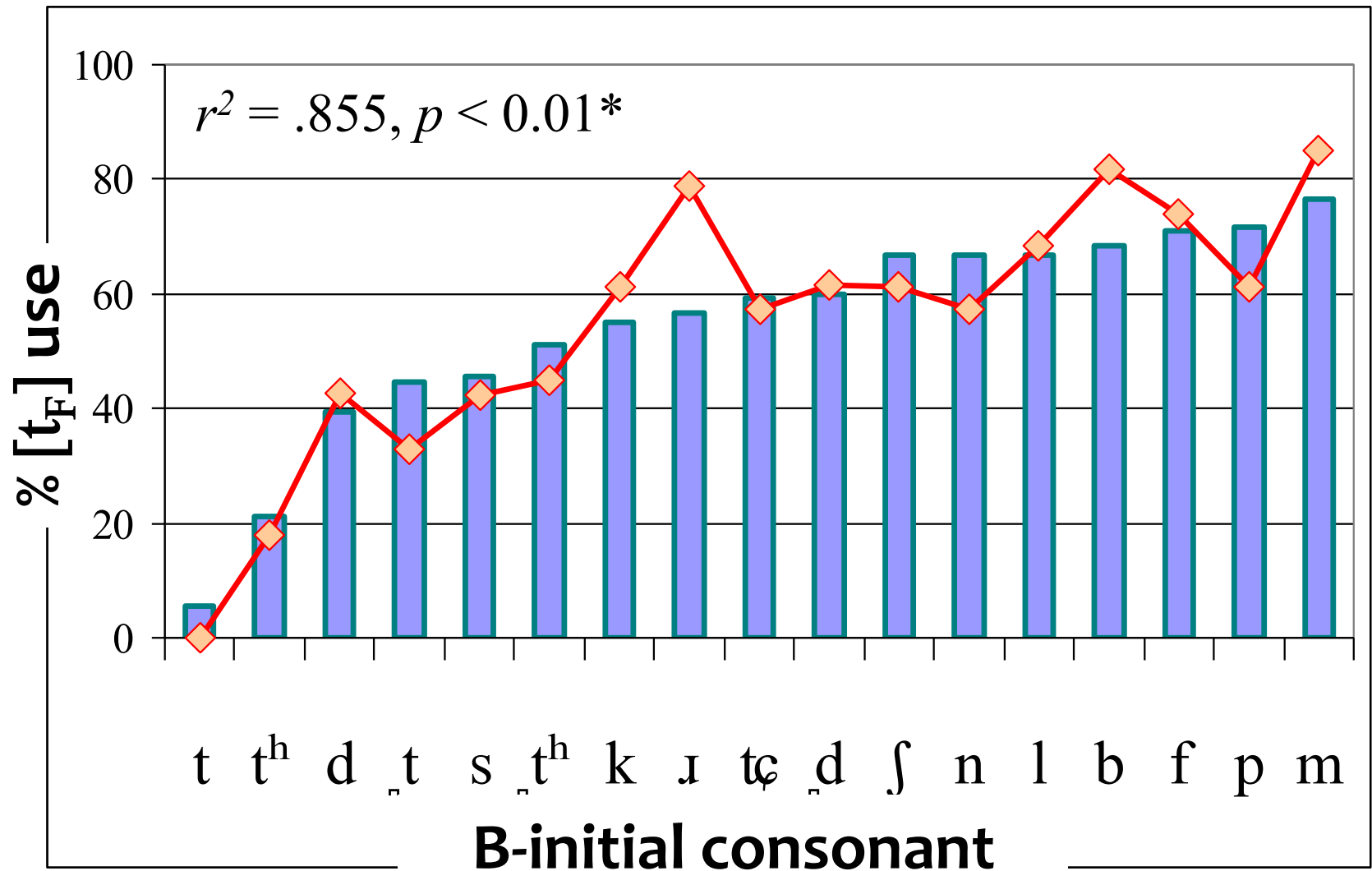
$$P = ((m!) \div (n!(m-n)!)) (1 - \text{sim}(C_1, t))^n (\text{sim}(C_1, t))^{m-n}$$

Probability that  $C_1$ -initial base will be reduplicated with  $[t_F]$   $n$  times out of a total of  $m$  trials

$m$  = number of reduplications for  $C_1$ -initial word

$n$  = number of reduplications with  $[t_F]$  for  $C_1$ -initial word

# Feature weighting (line) vs. observed (bars)





## Feature weighting: discussion

---

- A model of similarity that takes **feature weights** into account can closely model the data ( $r^2 = .855$ )
- Of course, in our case, we used the data to determine the weights

*I'll talk about some ideas of where this could independently come from in a minute...*



## General discussion

---

Summary and further questions



## Summary

---

- Crosslinguistically, B and R in echo reduplication must be sufficiently different
- In most lgs/studies, this is taken to be a **categorical nonidentity constraint**
  - “B and R must be **non-identical**”
  - Assumes that sound pairs can be categorized as either “identical” or “non-identical”
  - Even one BR feature mismatch should suffice



## Summary

---

- Data from English show that the constraint is actually **sensitive to phonetic similarity**, not just identity
  - “B and R must be **dissimilar**”
  - Still assumes categorical grouping of sound pairs: “identical”, “similar”, and “dissimilar”





## Summary

---

- Experimental data from Bengali confirm that **BR-assessment is sensitive to phonetic similarity**, not just identity
- The data also show that in fact, **similarity is gradient**
  - **Cannot group sound pairs categorically** as “identical”, “similar”, and “dissimilar”



## Summary

---

- Speakers compute the **similarity score of two sounds** using a metric that takes different **feature weights** into account
- We can derive the feature weights from the pattern itself...
- ...but where do speakers actually get these weights from independently?



# Summary

---

- Weights do not come from the lexicon
  - Similarity in echo reduplication is **not correlated with lexical patterns**



## Further questions

---

- Alternatively, feature weights could come **directly from the phoneme inventory**
  - The features that are weighted heavily are:
    - [voice]: 0.554
    - [distributed] (=dental vs. alveolar): 0.400
    - [strident]: 0.249
    - [spread glottis] (=aspiration): 0.198
  - All others are weighted 0.1



## Further questions

---

- These are also the features that help make the Bengali phoneme inventory so **coronal-heavy**
- In fact, of all features, these make the **most phonemic contrasts on their own** in the lg
- Thus, there might be **independent evidence of their “weight”**
- Need to do more work to confirm this



## Further questions

---

- How much of this is **language-specific**?
  - Coronal-heavy inventory?
- How much is **universal**?
  - Feature inventory?
- Future studies on **similarity-sensitive phenomena in other lgs** will take on these questions to build this model of similarity



# Thank you!

---

Special thanks to Kie Ross Zuraw, Colin Wilson, Bruce Hayes, Farida Amin Khan, the 30 participants in my experiment, and everyone in attendance here.