# LEXICAL TRANSFORMATIONS IN BLOGSPACE

A CA — RM CULTURAL EVOLUTION

from **The Semantic Drift of Quotations in Blogspace: A Case Study in Short-Term Cultural Evolution**

COGNITIVE SCIENCE
A Multidisciplinary Journal
(2017) 1–32

Sébastien Lerique
(EHESS / Centre Marc Bloch Berlin)

Camille Roth
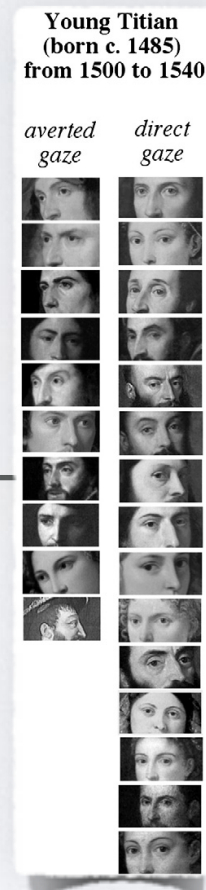(Sciences Po / Centre Marc Bloch Berlin)

# EMPIRICAL STUDY OF CULTURAL EVOLUTION

## IN VIVO

- *using historical data:* e.g.,

    - Morin 2013

    - Miton et al. 2015

# EMPIRICAL STUDY OF CULTURAL EVOLUTION

## IN VIVO

- *using historical data: e.g.,*

  - Morin 2013
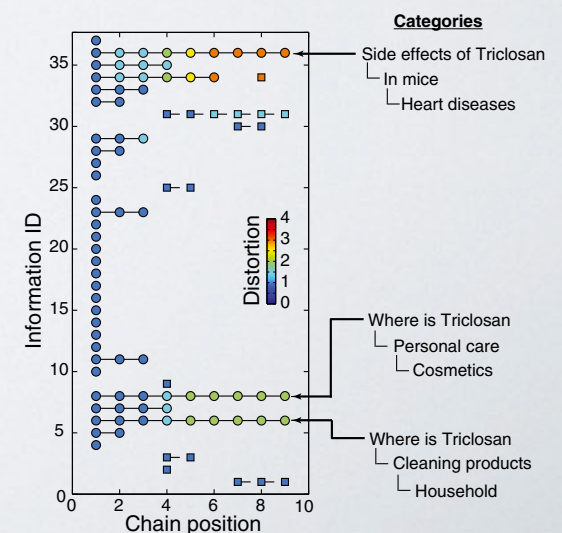
  - Miton et al. 2015

## IN VITRO

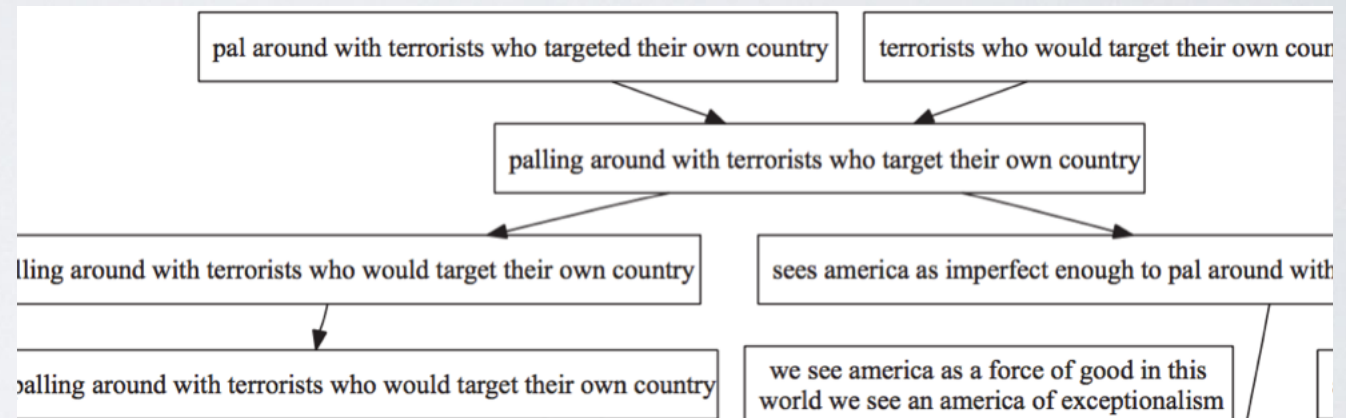- *using transmission chains: e.g.,*

  - Claidière et al. 2014

  - Moussaïd et al. 2015

# IN VIVO ONLINE DATA
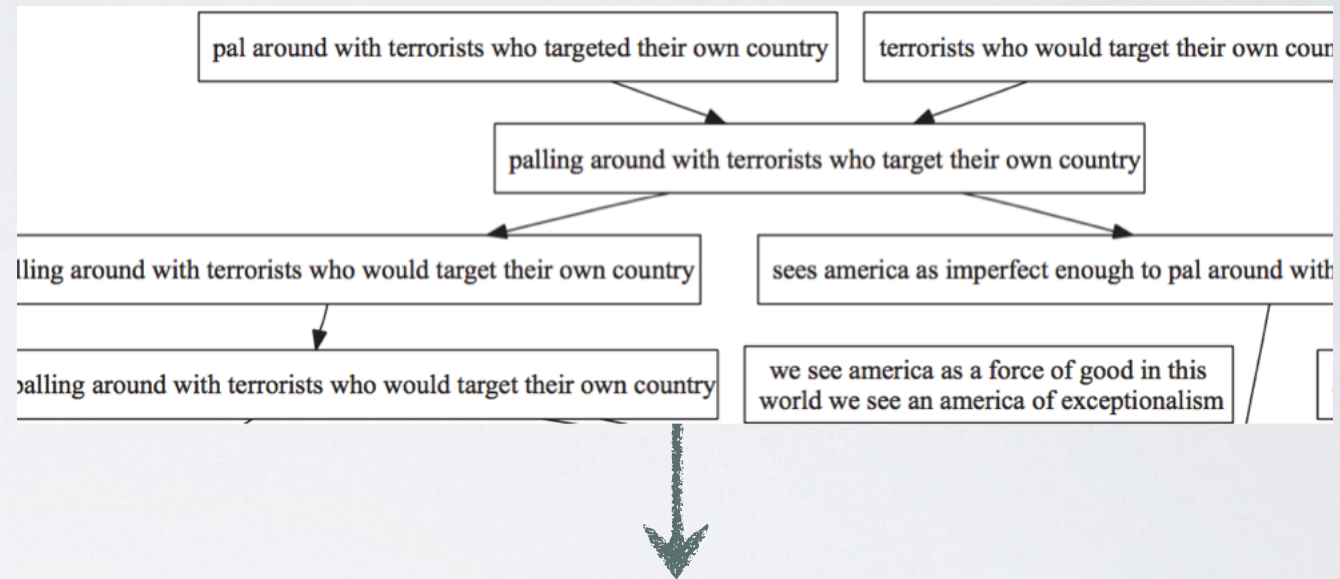
(Leskovec, Backstrom, Kleinberg, 2009)

**Corpus of quotations from a large corpus of (8.5m) blog posts**
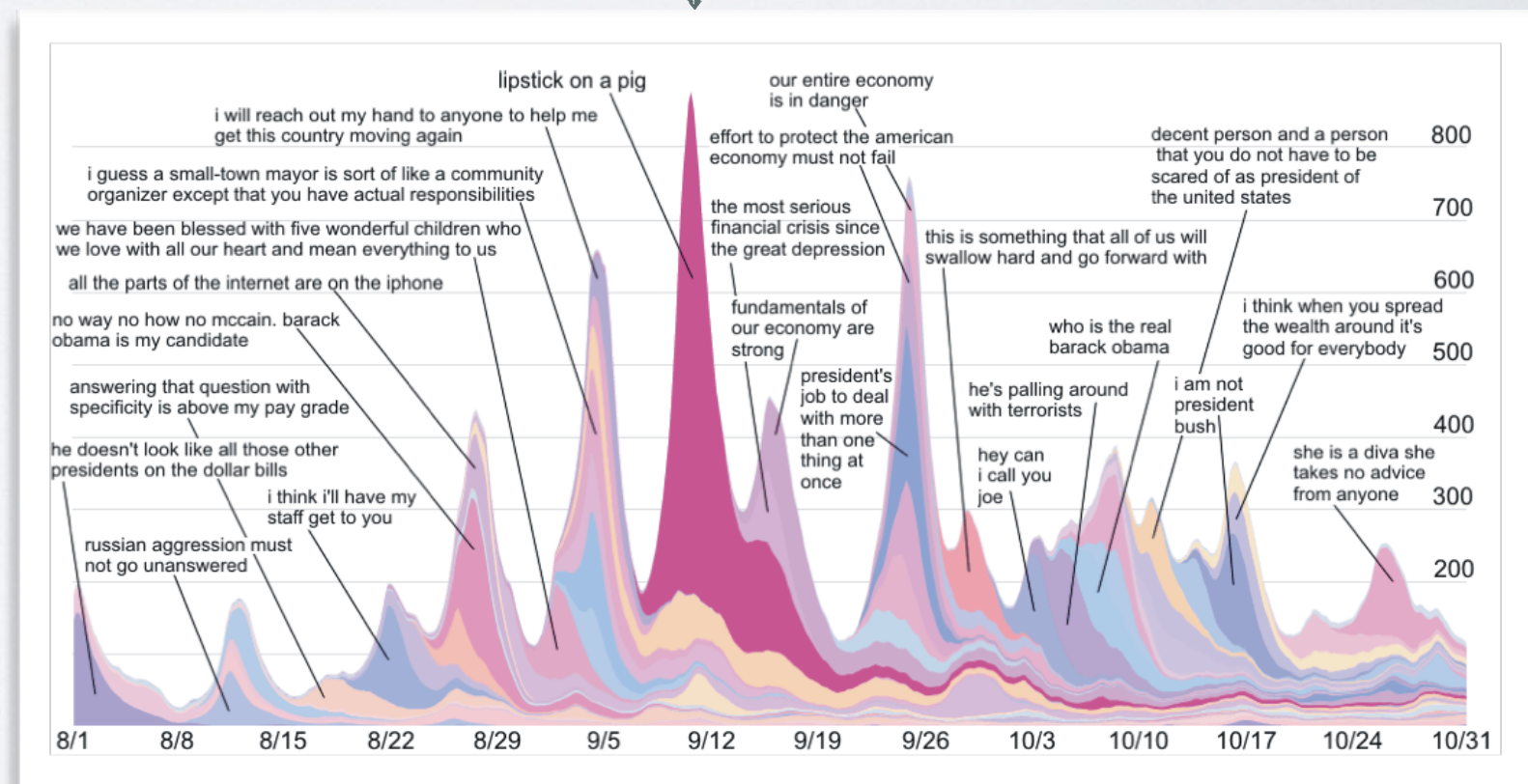
(Aug'08-Apr'09)

# IN VIVO ONLINE DATA

(Leskovec, Backstrom, Kleinberg, 2009)

**Corpus of quotations from a large corpus of (8.5m) blog posts**

(Aug'08-Apr'09)

**Groups (and dynamics) of sentences**

# SENTENCE REFORMULATION

○ *Pakistani President Asif Ali Zardari:*

- ○ "we will not be scared of these cowards"

- ▷ "we will not be **afraid** of these cowards."

○ *US Senator McCain:*

- ○ "I admire Senator Obama and his accomplishments"

- ▷ "I **respect** Senator Obama and his accomplishments."

# SENTENCE REFORMULATION

- *Pakistani President Asif Ali Zardari:*

  - "we will not be scared of these cowards"

  - "we will not be **afraid** of these cowards."

- *US Senator McCain:*

  - "I admire Senator Obama and his accomplishments"

  - "I **respect** Senator Obama and his accomplishments."

- Task similar to word (list) recall

# SENTENCE REFORMULATION

- *Pakistani President Asif Ali Zardari:*
  - "we will not be scared of these cowards"
  - ▷ "we will not be **afraid** of these cowards."

- *US Senator McCain:*
  - "I admire Senator Obama and his accomplishments"
  - ▷ "I **respect** Senator Obama and his accomplishments."

- Task similar to word (list) recall

- Lexical features expected to influence the likelihood of substitution

# SENTENCE REFORMULATION

- *Pakistani President Asif Ali Zardari:*
  - "we will not be scared of these cowards"
  - ▷ "we will not be **afraid** of these cowards."

- *US Senator McCain:*
  - "I admire Senator Obama and his accomplishments"
  - ▷ "I **respect** Senator Obama and his accomplishments."

- Task similar to word (list) recall

- Lexical features expected to influence the likelihood of substitution

  - *for instance:* word frequency, age of acquisition, number of phonemes, phonological neighborhood density, position in a semantic network...

# SENTENCE REFORMULATION

- *Pakistani President Asif Ali Zardari:*

  - "we will not be scared of these cowards"

  - "we will not be **afraid** of these cowards."

- *US Senator McCain:*

  - "I admire Senator Obama and his accomplishments"

  - "I **respect** Senator Obama and his accomplishments."

- Task similar to word (list) recall

- Lexical features expected to influence the likelihood of substitution

  - *for instance:* word frequency, age of acquisition, number of phonemes, phonological neighborhood density, position in a semantic network...

- Address e.g., the "word-frequency paradox" (Mandler et al. 1982)
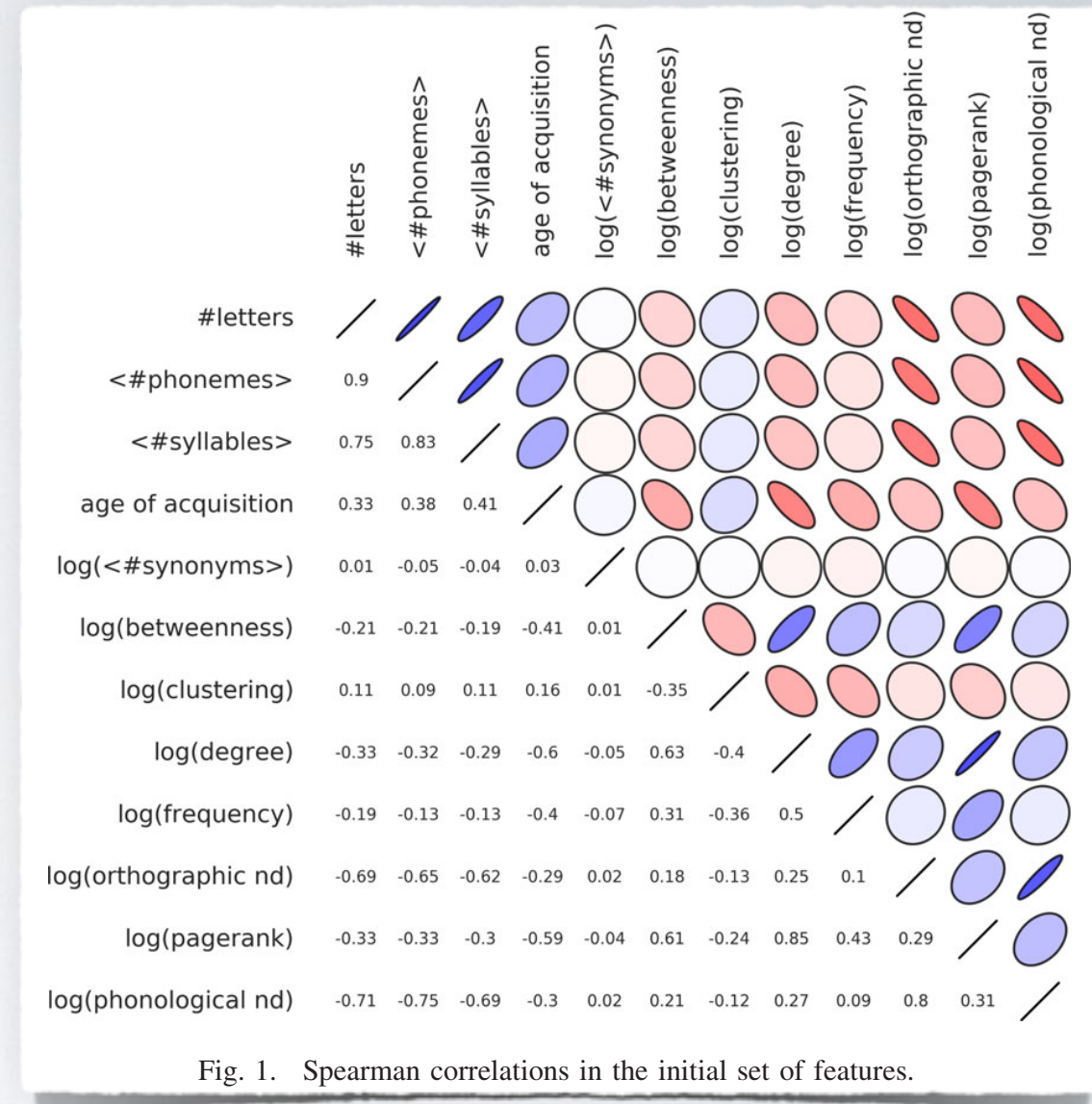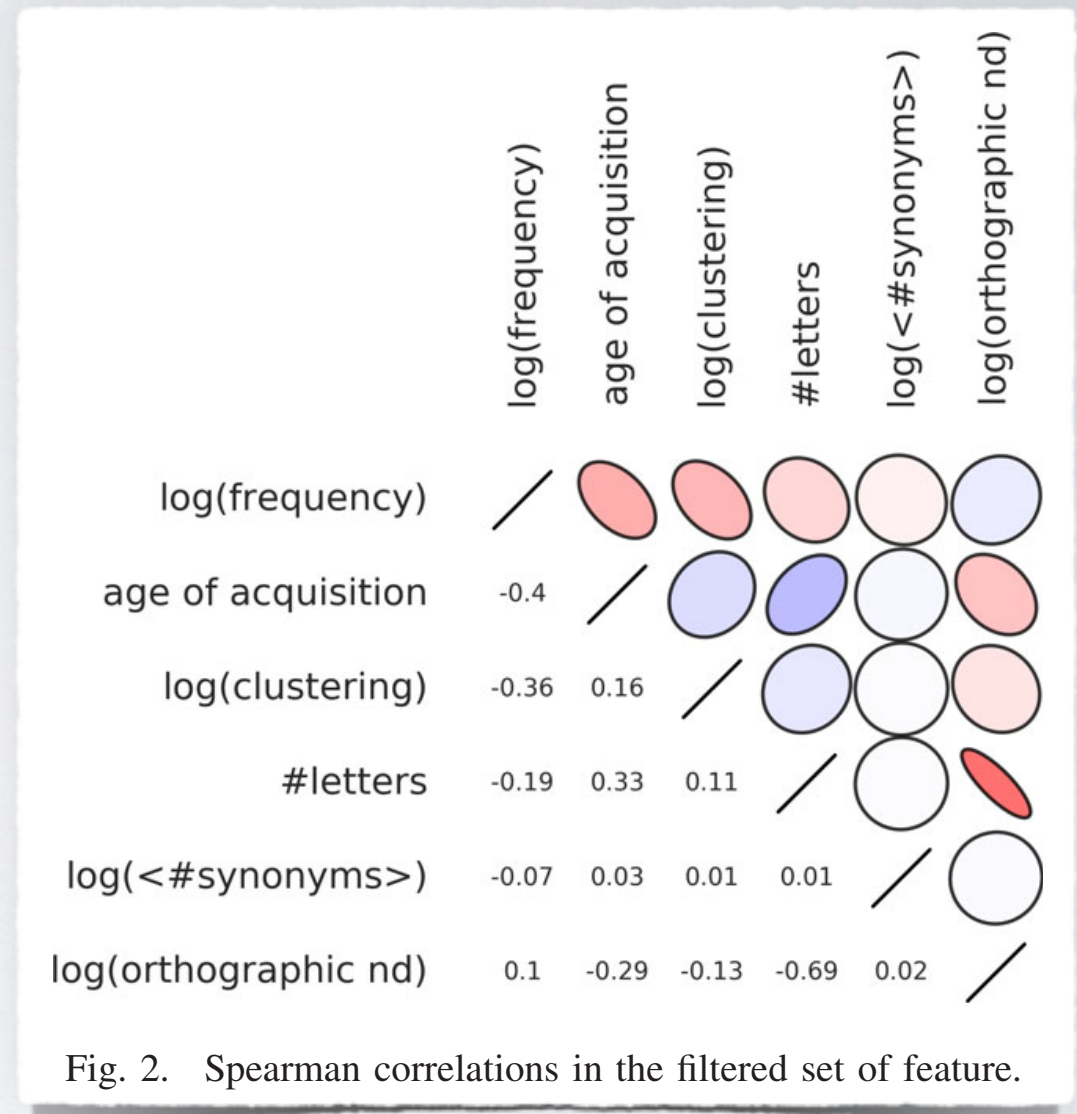
# SENTENCE REFORMULATION

- *Pakistani President Asif Ali Zardari:*
  - "we will not be scared of these cowards"
  - "we will not be **afraid** of these cowards."

- *US Senator McCain:*
  - "I admire Senator Obama and his accomplishments"
  - "I **respect** Senator Obama and his accomplishments."

- Task similar to word (list) recall

- Lexical features expected to influence the likelihood of substitution

  - *for instance:* word frequency, age of acquisition, number of phonemes, phonological neighborhood density, position in a semantic network...

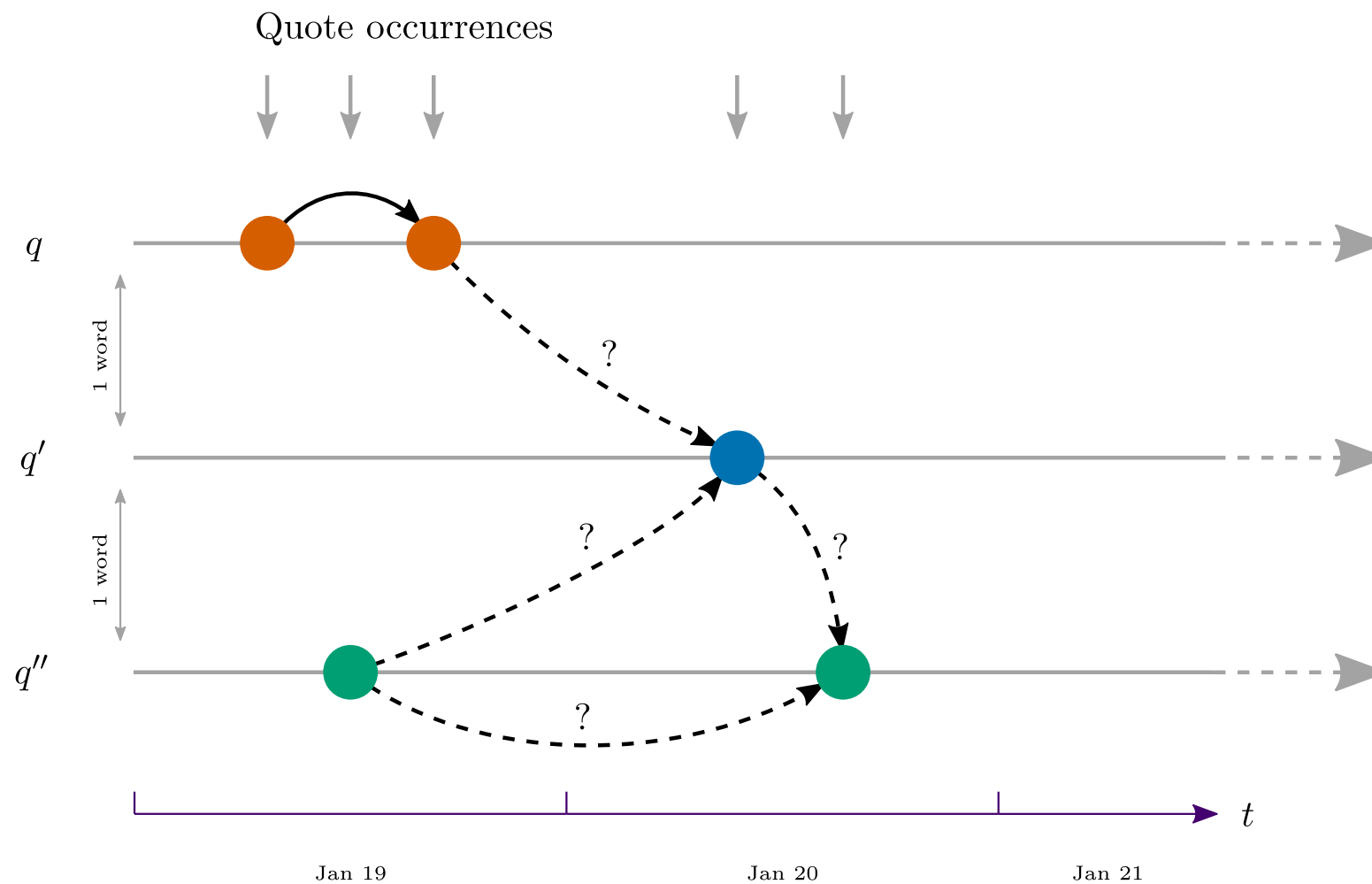- Address e.g., the "word-frequency paradox" (Mandler et al. 1982)



Fig. 1. Spearman correlations in the initial set of features.

# SENTENCE REFORMULATION

- *Pakistani President Asif Ali Zardari:*
  - "we will not be scared of these cowards"
  - "we will not be **afraid** of these cowards."

- *US Senator McCain:*
  - "I admire Senator Obama and his accomplishments"
  - "I **respect** Senator Obama and his accomplishments."

- Task similar to word (list) recall

- Lexical features expected to influence the likelihood of substitution

  - *for instance:* word frequency, age of acquisition, number of phonemes, phonological neighborhood density, position in a semantic network...

- Address e.g., the "word-frequency paradox" (Mandler et al. 1982)

|  | log(frequency) | age of acquisition | log(clustering) | #letters | log(<#synonyms>) | log(orthographic nd) |
|---|---|---|---|---|---|---|
| log(frequency) |  |  |  |  |  |  |
| age of acquisition | -0.4 |  |  |  |  |  |
| log(clustering) | -0.36 | 0.16 |  |  |  |  |
| #letters | -0.19 | 0.33 | 0.11 |  |  |  |
| log(<#synonyms>) | -0.07 | 0.03 | 0.01 | 0.01 |  |  |
| log(orthographic nd) | 0.1 | -0.29 | -0.13 | -0.69 | 0.02 |  |

Fig. 2.   Spearman correlations in the filtered set of feature.

# SUBSTITUTION MODEL



Fig. 3. Possible paths from occurrence to occurrence: $q$, $q'$, and $q''$ are three quotation variants belonging to the same cluster. $q$ and $q''$ differ by two words, but $q'$ differs from both $q$ and $q''$ by one word. The second occurrence of $q$ can safely be considered a faithful copy of the first, but the occurrences of $q'$ and $q''$ are uncertain: While the first occurrence of $q'$ is most likely a substitution for $q$, it could also stem from $q''$; conversely, the second occurrence of $q''$ could also be a substitution for $q'$ instead of being a faithful copy of its first occurrence.

# SUSCEPTIBILITY $\quad \sigma_g = \dfrac{s_g}{s_g^0}$



Fig. 5. Part-of-Speech-related results: Categories are simplified from the TreeTagger tag set: *C* means *Closed class-like* (see main text for details), *J* means adjective, *N* noun, *R* adverb, and *V* means verb. The top panel shows the actual $s_{POS}$ and $s_{POS}^0$ counts. The bottom panel shows the substitution susceptibility $\sigma_{POS}$, which is the ratio between the two previous counts. Confidence intervals are computed with the Goodman (1965) method for multinomial proportions.
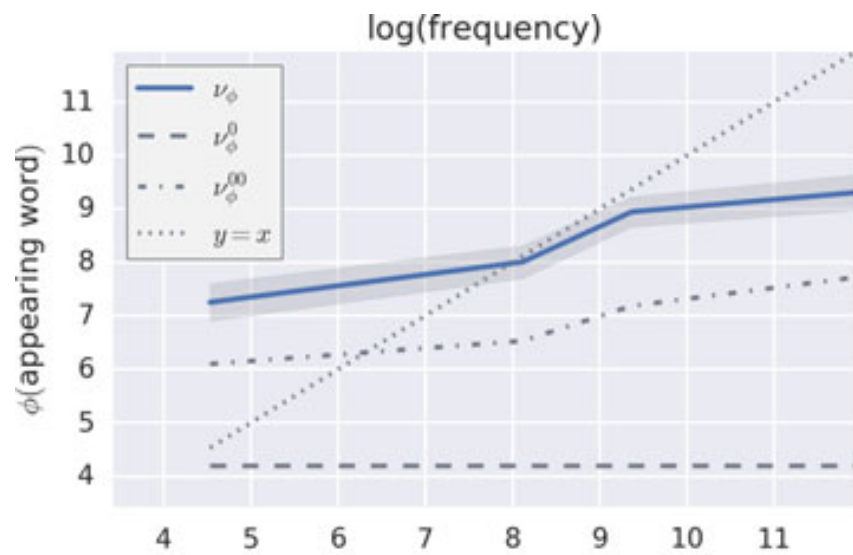
# SUSCEPTIBILITY $\quad \sigma_g = \dfrac{s_g}{s_g^0}$



On the whole, the trends observed are consistent with known effects of word frequency, age of acquisition, and number of letters, indicating that the triggering of a substitution could behave quite similarly to word recall in standard tasks.
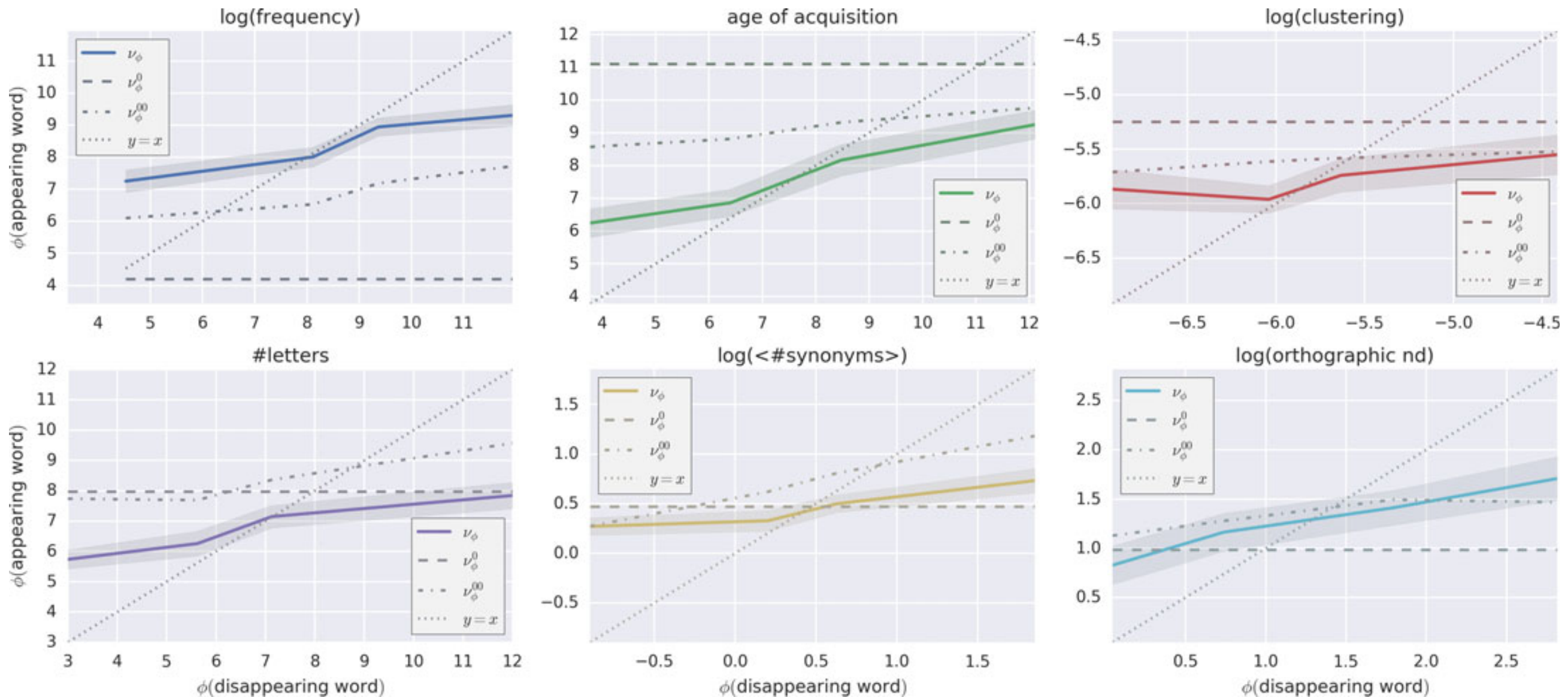
# FEATURE VARIATION

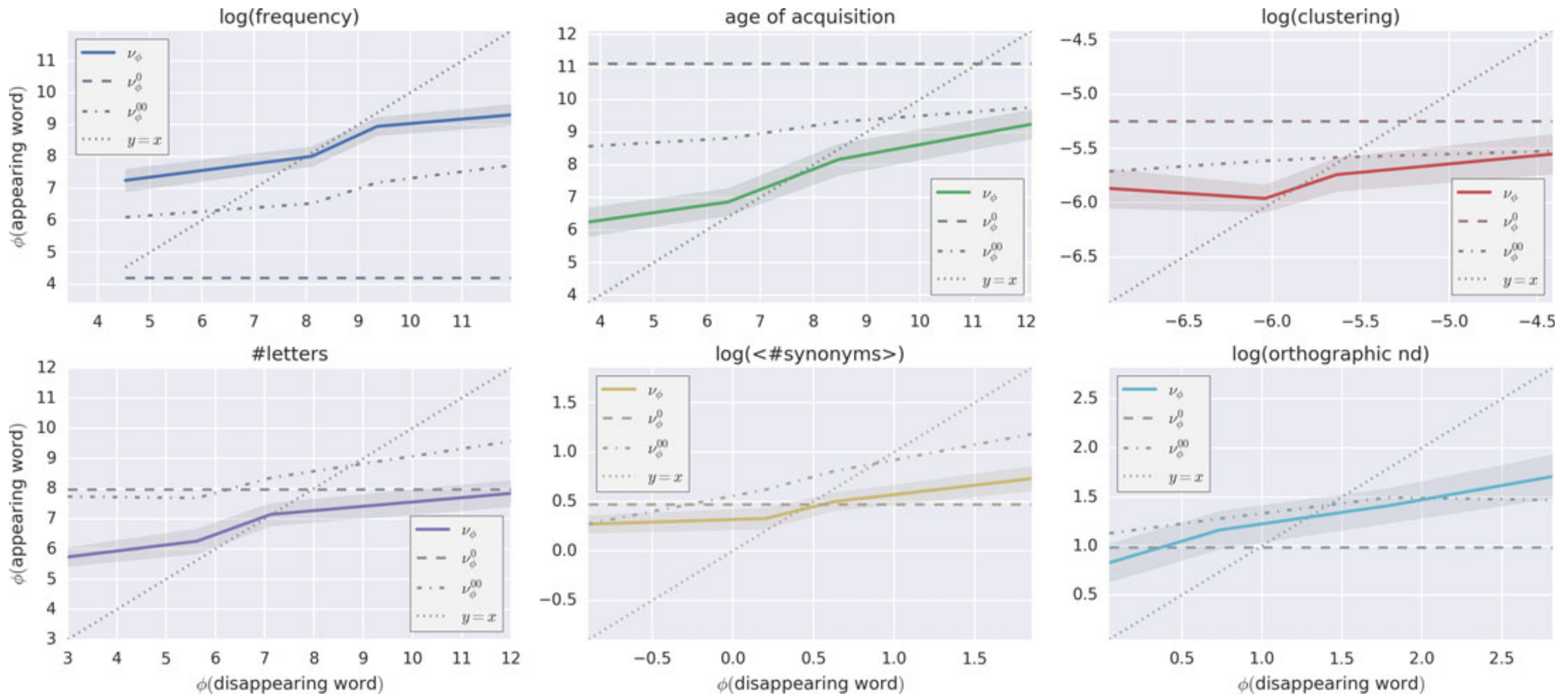$$\nu_\phi(f) = \langle \phi(w') \rangle_{\{w \to w' | \phi(w) = f\}}$$

# FEATURE VARIATION

$$\nu_\phi(f) = \langle \phi(w') \rangle_{\{w \to w' | \phi(w) = f\}}$$



First, there is a single intersection of $\nu_\phi$ with $y=x$ and the slope of $\nu_\phi$ remains smaller than 1:
**the substitution process exhibits a single attractor**

# FEATURE VARIATION

$$\nu_\phi(f) = \langle \phi(w') \rangle_{\{w \to w' | \phi(w) = f\}}$$



First, there is a single intersection of $\nu_\phi$ with $y=x$ and the slope of $\nu_\phi$ remains smaller than 1:
**the substitution process exhibits a single attractor**

Second, the comparison with $\nu_\phi^0$ and $\nu_\phi^{00}$ shows that there are two classes of attractors, depending on whether:

1. there is a triple intersection (of $y = x$, $\nu_\phi$, and $\nu_\phi^0$ or $\nu_\phi^{00}$);
2. or $\nu_\phi$ always remains above or below $\nu_\phi^0$ and $\nu_\phi^{00}$.

# COMBINED EFFECTS

To make sure our observations are not the product of correlations or interactions, we model the variations of the six features as a linear function of the start word's feature values:
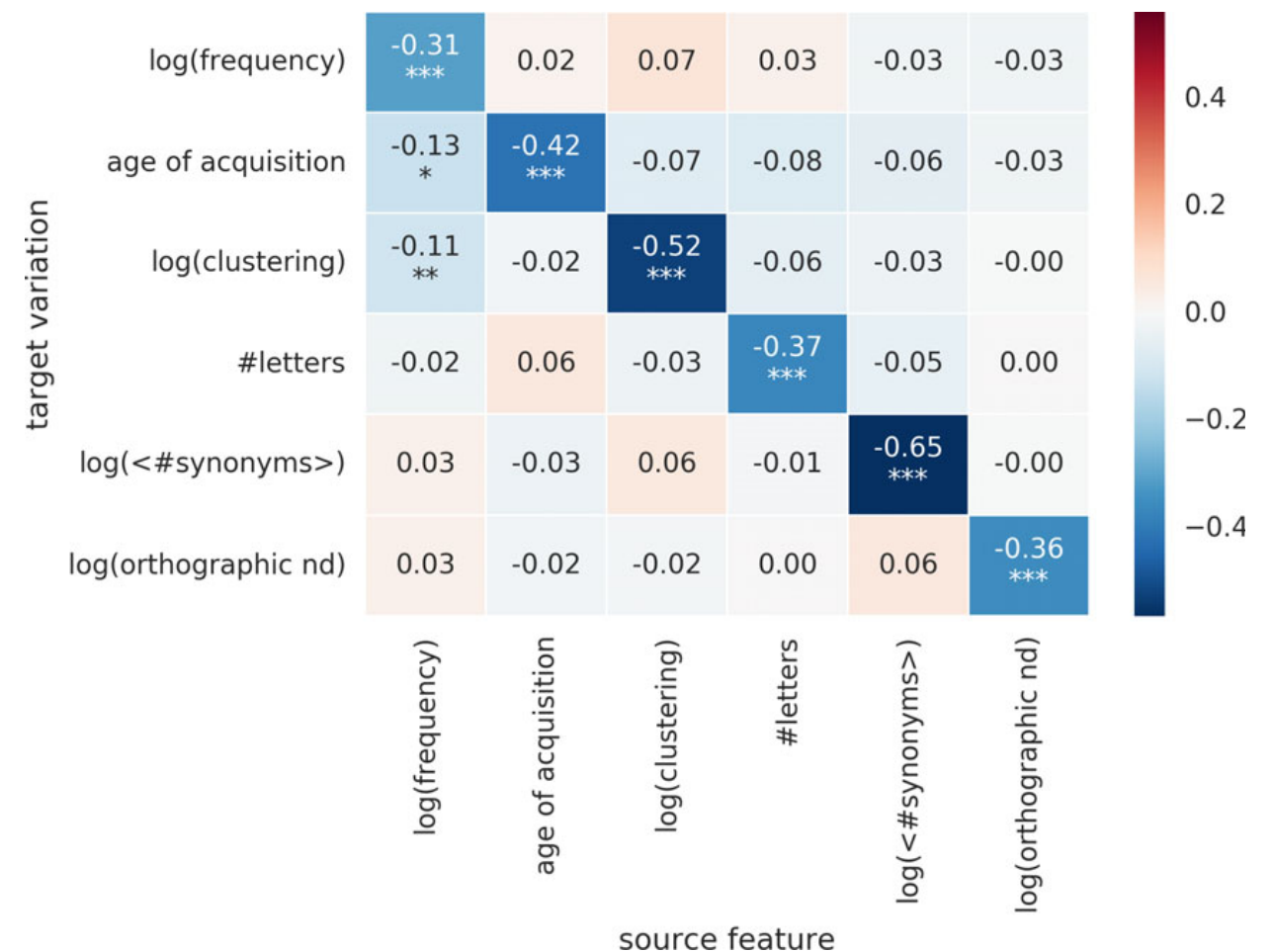
$$\phi(w') - \phi(w) = A + B \cdot \phi(w)$$

where $\phi$ is the vector of all six features of a word, $A$ is an intercept vector, and $B$ is a $6 \times 6$ coefficients matrix. This regression achieves an overall $R^2$ of .33.
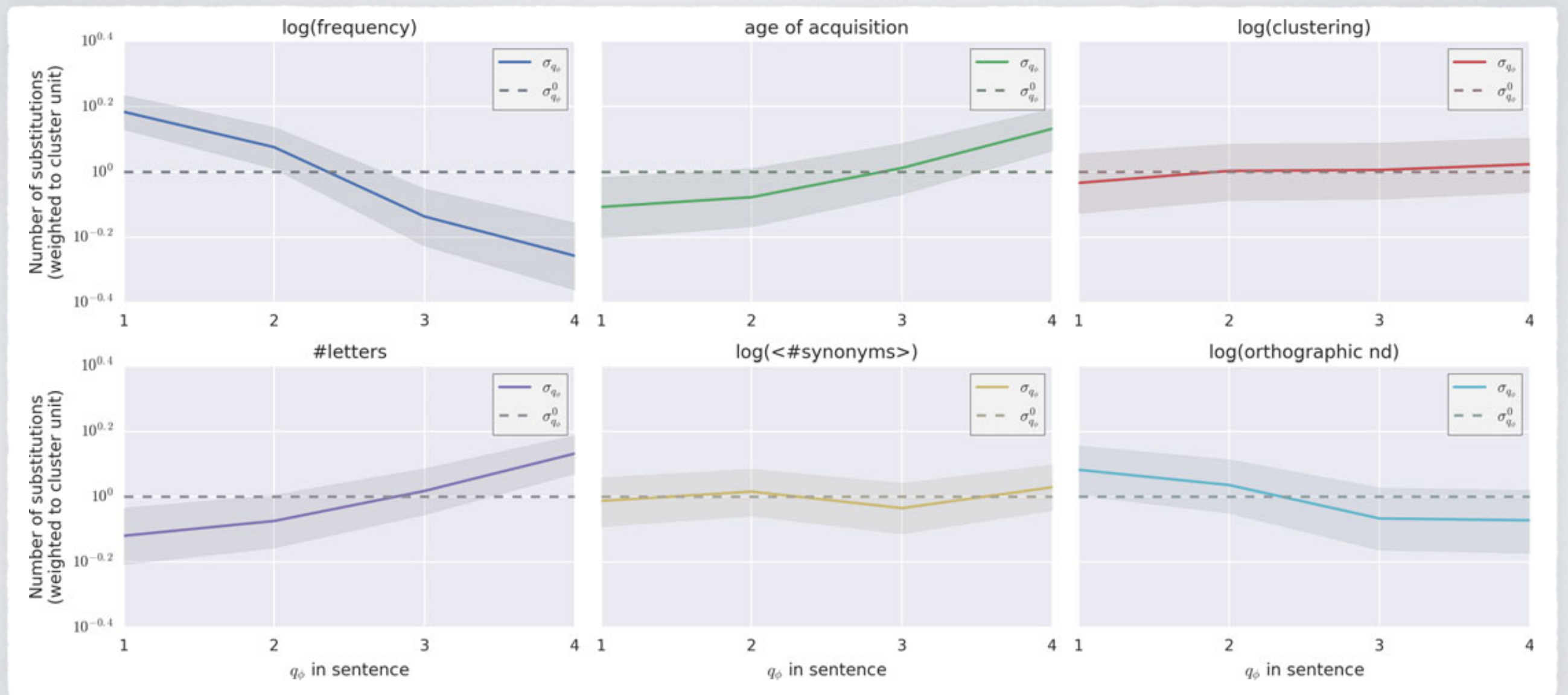
- *Burmese poet Saw Wai (Nov 2008):*
  - "Senior general Than Shwe is foolish with power"
  - "Senior general Than Shwe is **crazy** with power"

"foolish": *8.94 y.o., 675 times, cc of .0082*
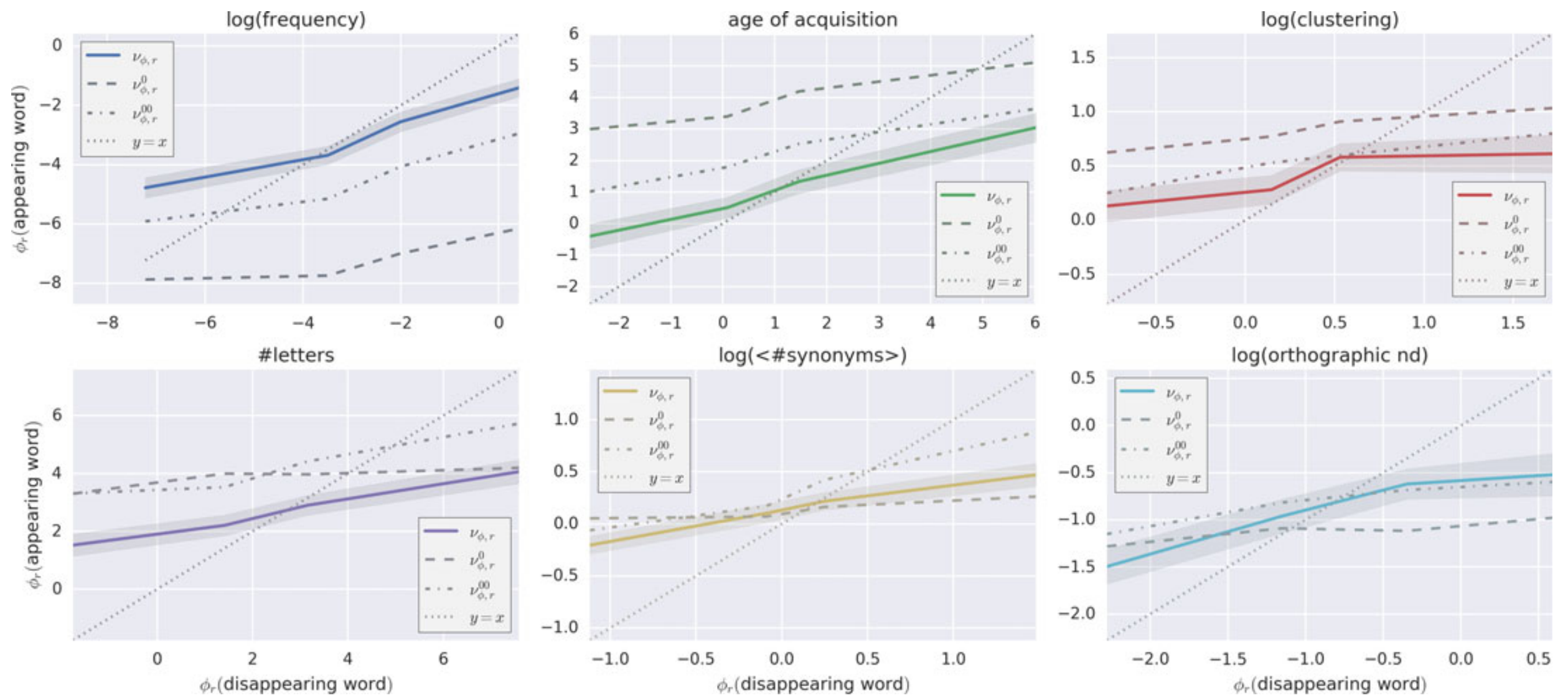>"crazy": *5.22 y.o., 4100 times, cc of .0017*

# TAKING SENTENCE CONTEXT INTO ACCOUNT



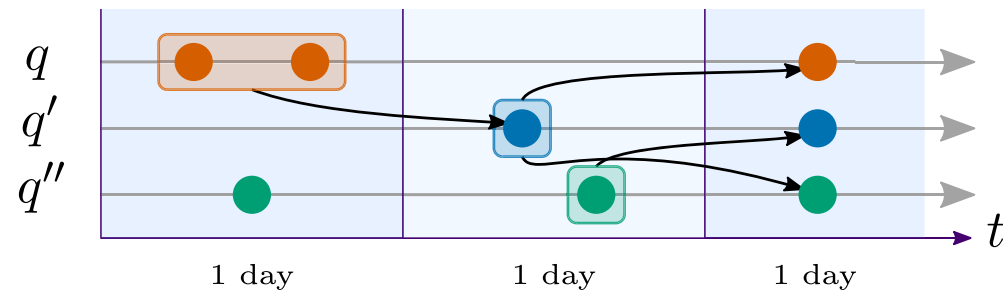susceptibility based on the position of the word in the sentence (quartiles)

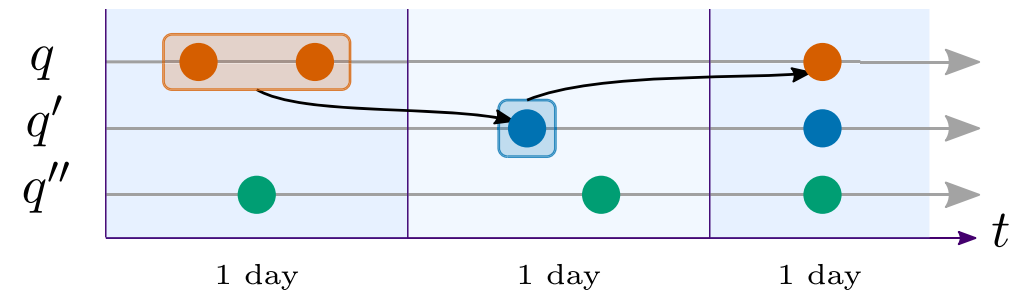# TAKING SENTENCE CONTEXT INTO ACCOUNT



feature variation w.r.t. median feature value in the sentence

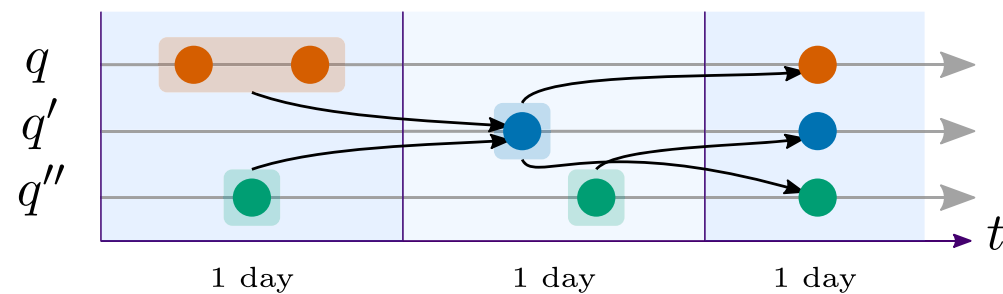The speaker says "Thanks" –> "<u>Danke</u>"
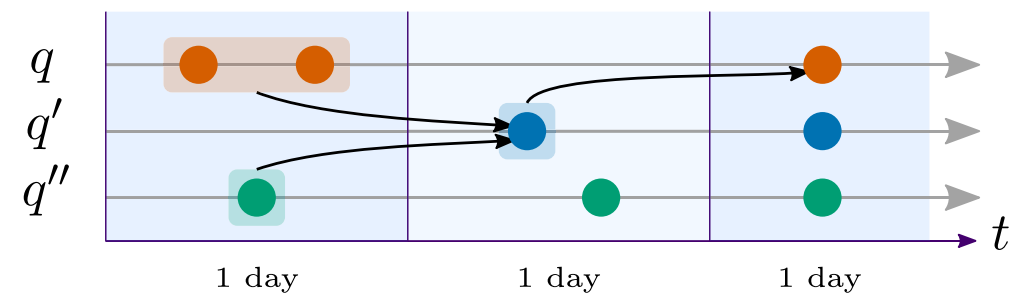
# SUBSTITUTION MODEL VARIANTS



(a) Source must be majority in preceding bin, destination can be anything

(b) Source must be majority in preceding bin, destination must not appear in preceding bin

(c) Source can be anything, destination can be anything

(d) Source can be anything, destination must not appear in preceding bin

(i) bin position   (ii) bin length   (iii) candidate sources   (iv) candidate destinations