

# Rare words and scaling laws in language



Eduardo Goldani Altmann

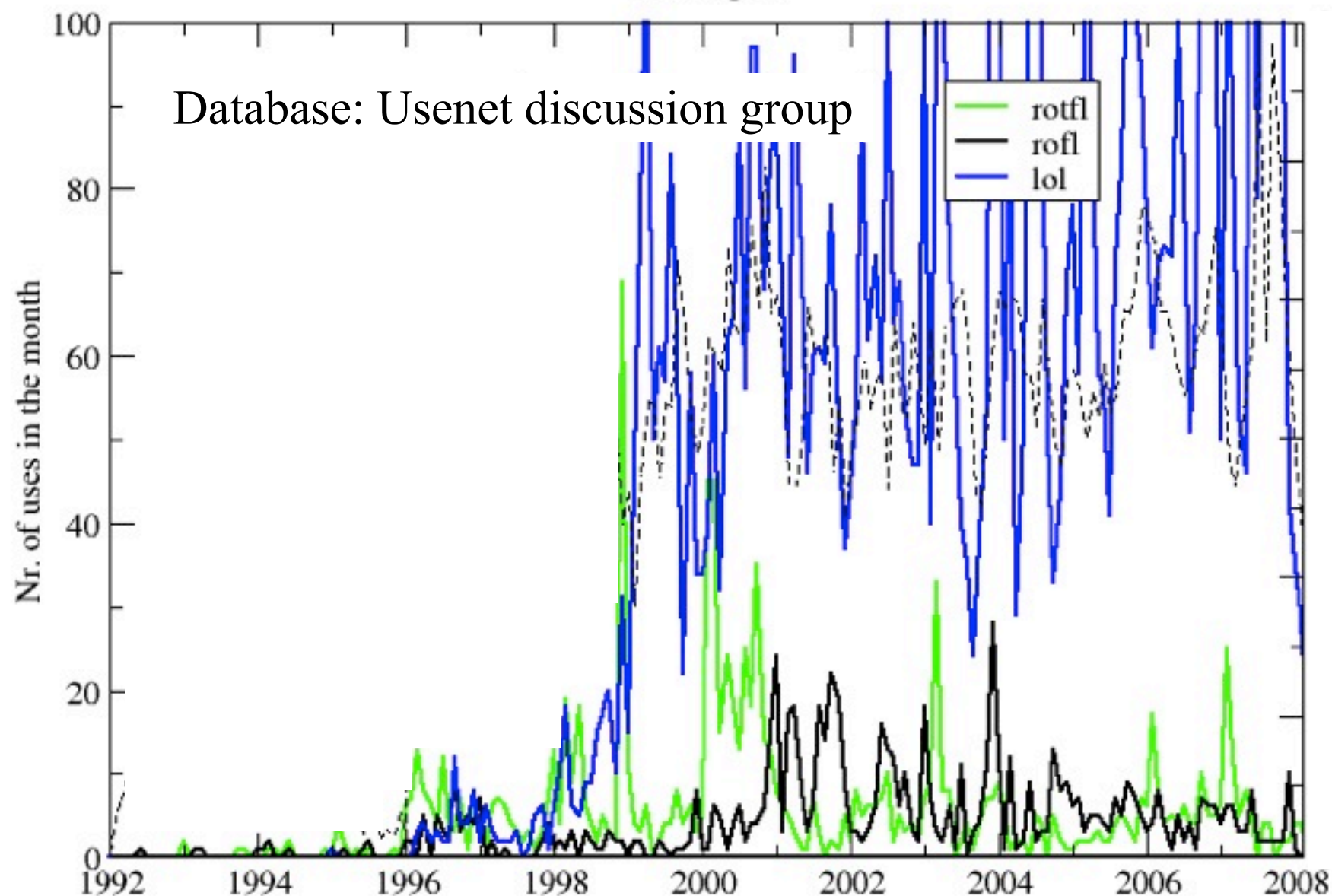
Max Planck Institute  
for the Physics of Complex Systems,  
Dresden, Germany



Paladin Conference, Rome  
September 24, 2013

# Statistical Analysis of Language

- Amount of produced data
  - Text messages: 10 B messages/day worldwide;
  - Twitter: 400 M tweets/day by 200M active users;
  - Wikipedia: 10 M contributors; 1 B words;
  - Google n-grams: 5 M books between [1520, 2000], 100 B words;
- Opportunity for applications (e.g., search engines, data mining) and scientific investigations (language as a lens on human activities and thought).



## Plan:

1. Vocabulary Growth

Centuries / millions of books



2. Innovations and Change

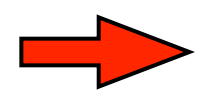


3. Text Analysis

Single book



## Plan:



1. Vocabulary Growth

Centuries / millions of books



2. Innovations and Change



3. Text Analysis

Single book



**Martin Gerlach** and E. G. Altmann,  
*Stochastic model for the vocabulary growth in natural languages*,  
Phys. Rev. X 3, 021006 (2013)

# Motivation



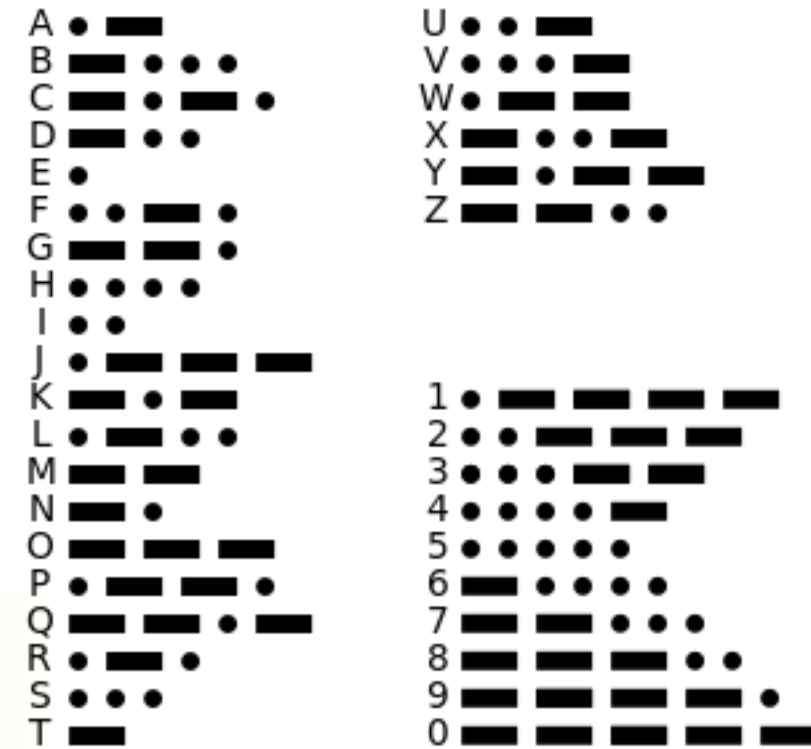
Samuel Morse (1791-1872)  
[Wikipedia]

message

the ... also .. prince ...

e t a i .. q y ..

rank=1 2 3 ...



International Morse Code  
[Wikipedia]

*“The dictionary or vocabulary consists of words alphabetically arranged and regularly numbered... so that each word in the language has its telegraphic number...”*

[Morse’s first telegraph patent as cited by J. Gleick, *The Information*]

How the vocabulary grows in time and with database size?

# Motivation: invert indexing

Vocabulary size  
=  
memory allocation



The diagram illustrates an inverted index structure. It consists of a 7x7 grid of cells. The columns are labeled 'page 1' through 'page 6' at the top. The rows are labeled with words: 'the', 'it', '...', 'word n', '...', and 'word N' on the left. A vertical double-headed arrow on the left side of the grid spans the height of the word rows, indicating the vocabulary size. Cells are colored green or red to represent word occurrences. 'the' is green in all pages. 'it' is green in pages 1-5 and red in page 6. 'word n' is green in pages 1 and 5, and red in pages 2, 3, 4, and 6. 'word N' is red in pages 1, 3, 4, 5, and 6, and green in page 2. Ellipsis rows are empty.

	page 1	page 2	page 3	page 4	page 5	page 6
"the"	green	green	green	green	green	green
"it"	green	green	green	green	green	red
...						
word n	green	red	red	red	green	red
...						
word N	red	green	red	red	red	red

## Motivation: vocabulary of a language?

*Report on the state of the **German** language* (March 2013)

German Academy for Language and Literature

Union of German Academies of Sciences and Humanities

Year	1905-1914	1948-1957	1995-2004
# distinct words	3,715,000	5,045,000	5,238,000

*Quantitative Analysis of Culture Using Millions of Digitized Books*

Michel et. al., Science (2011) [**English**]

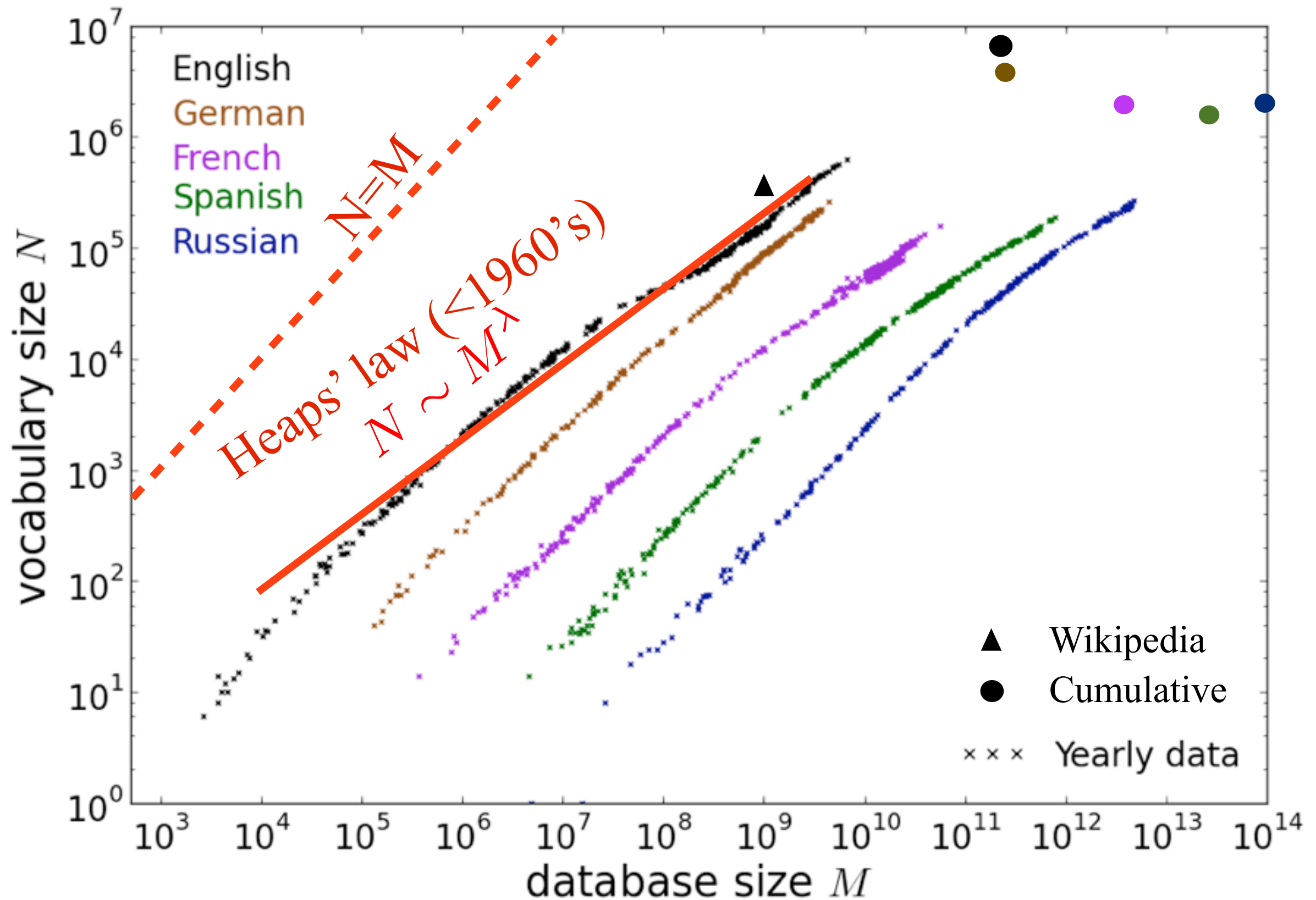
Year	1900	1950	2000
# distinct words	544,000	597,000	1,022,000

**Problem**: role of database size?



# Vocabulary growth with database size

Limit vocabulary?



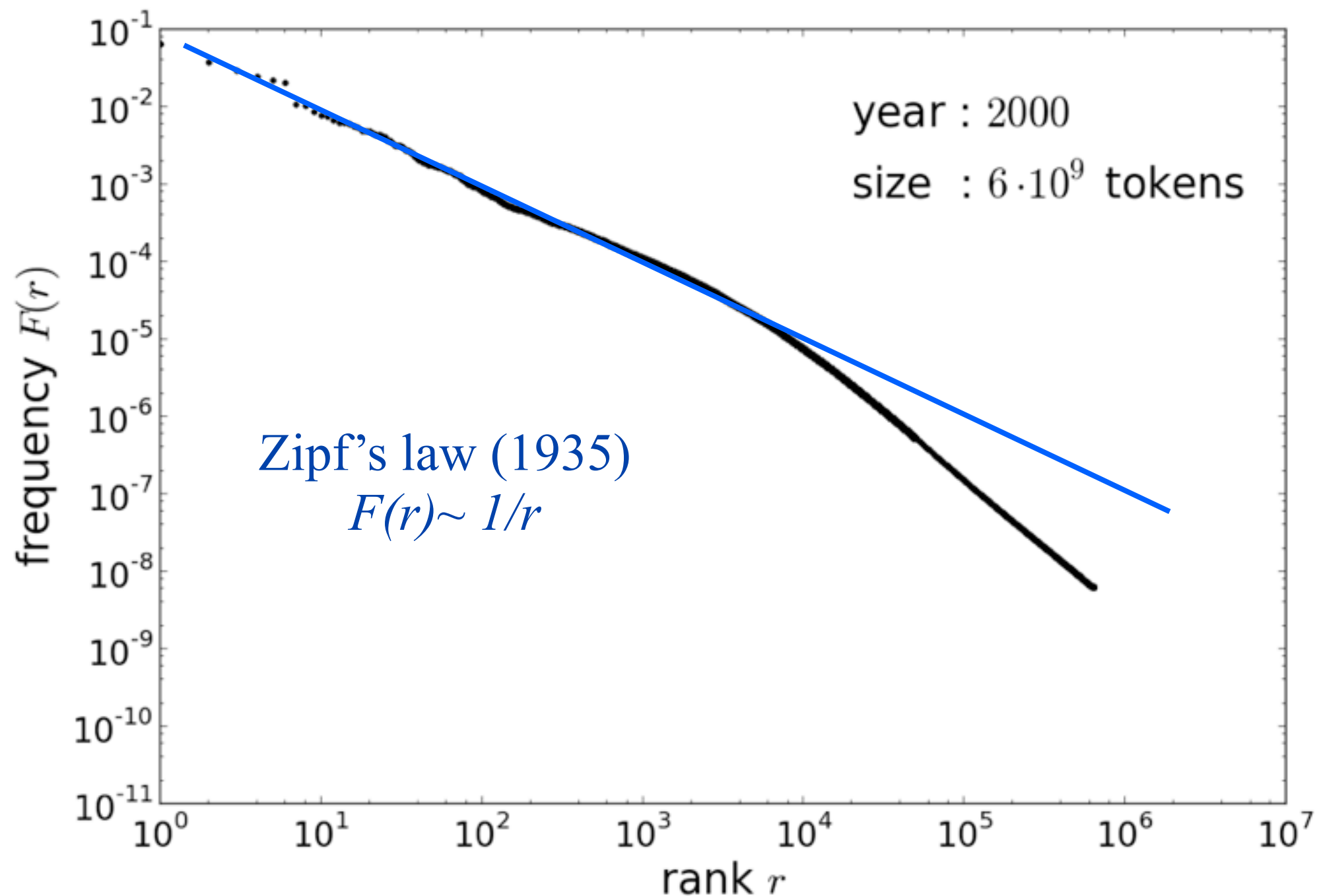


# Descriptive model

**Simple mode**: usage of each word follows a Poisson process with fixed frequency

$$\langle N(M) \rangle = \sum_r 1 - e^{-F(r)M}$$

where  $F(r)$  is the frequency of the  $r$ -th most frequent word ( $r = \text{rank}$ ).

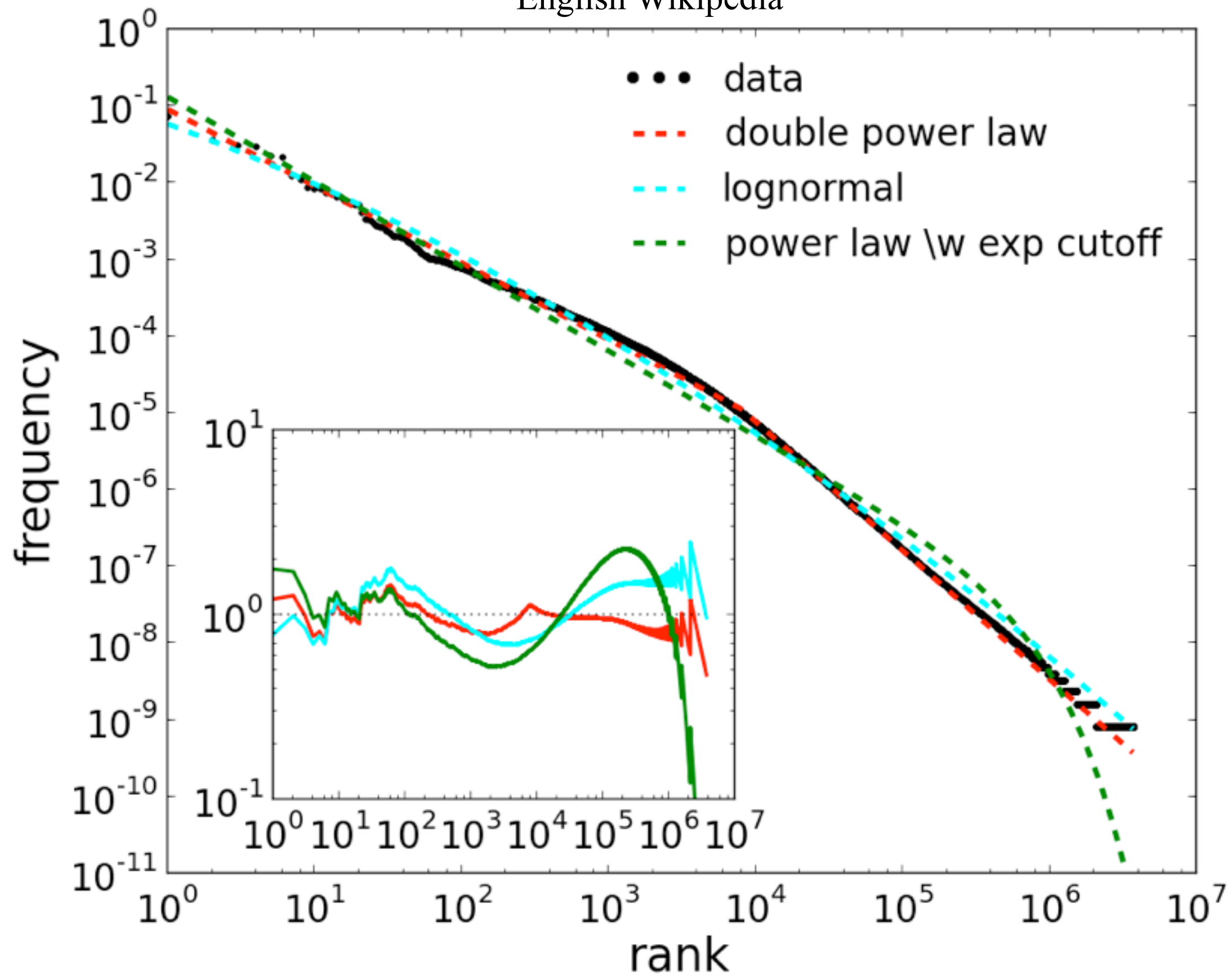


# Generalized Zipf's law

i	distribution	$F(r; \Omega)$	set of parameters $\Omega$
1	Power-Law	$Cr^{-\gamma}$	$\gamma$
2	Shifted Power-Law	$C(r + b)^{-\gamma}$	$\gamma, b$
3	Power-Law with Exponential cutoff (beginning)	$C \exp(-b/r) r^{-\gamma}$	$\gamma, b$
4	Power-Law with Exponential cutoff (tail)	$C \exp(-br) r^{-\gamma}$	$\gamma, b$
5	Log-normal	$Cr^{-1} \exp(-\frac{1}{2} (\ln r - \mu)^2 / \sigma^2)$	$\mu, \sigma$
6	Weibull	$Cr^{\gamma-1} \exp(-br^{-\gamma})$	$\gamma, b$
7	Double Power-Law	$C \begin{cases} r^{-1}, & r \leq b \\ b^{\gamma-1} r^{-\gamma} & r > b, \end{cases}$	$\gamma, b$

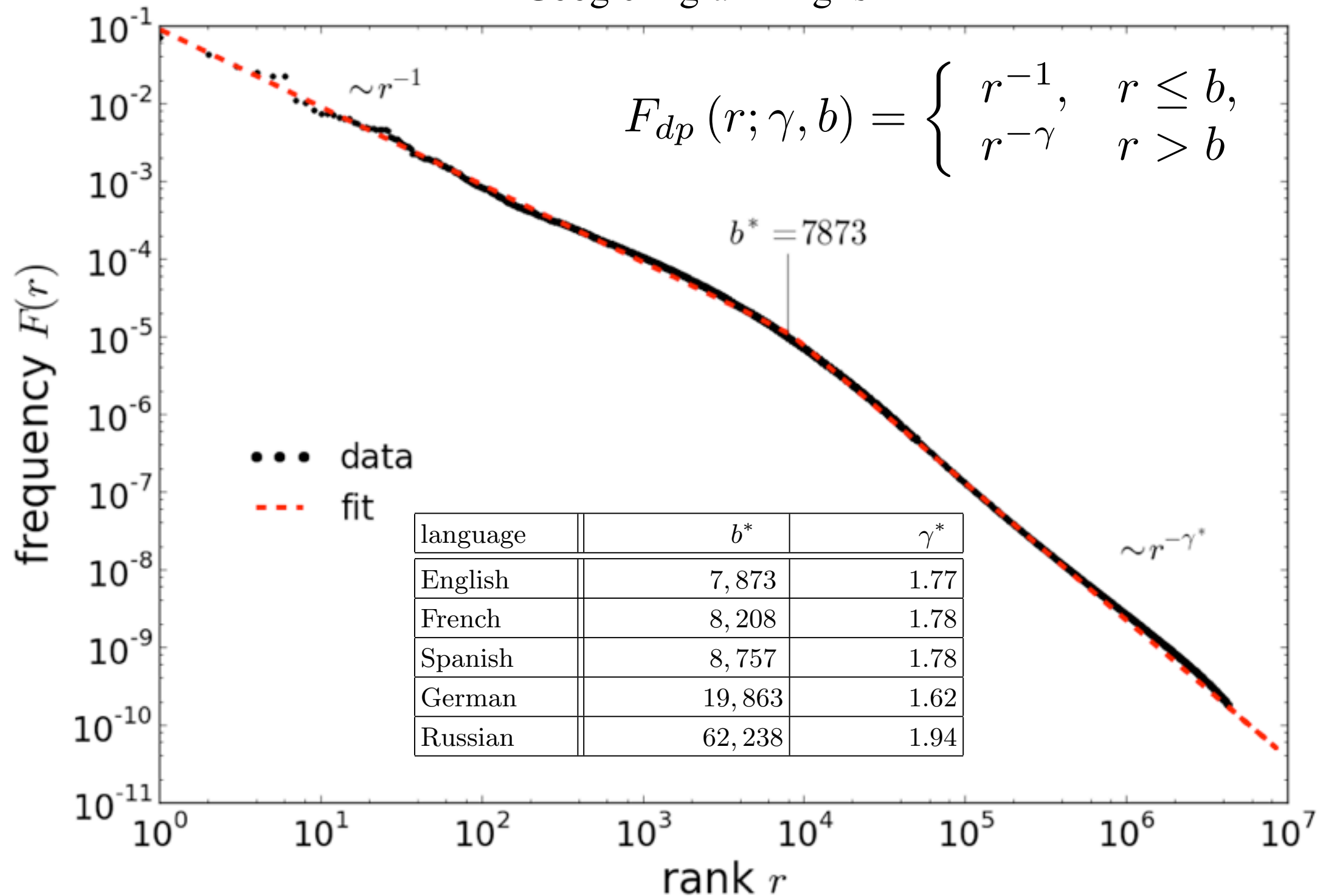
# Generalized Zipf's law

English Wikipedia



# Generalized Zipf's law

Google n-gram English



# Descriptive model

**Simple mode**: usage of each word follows a Poisson process with fixed frequency

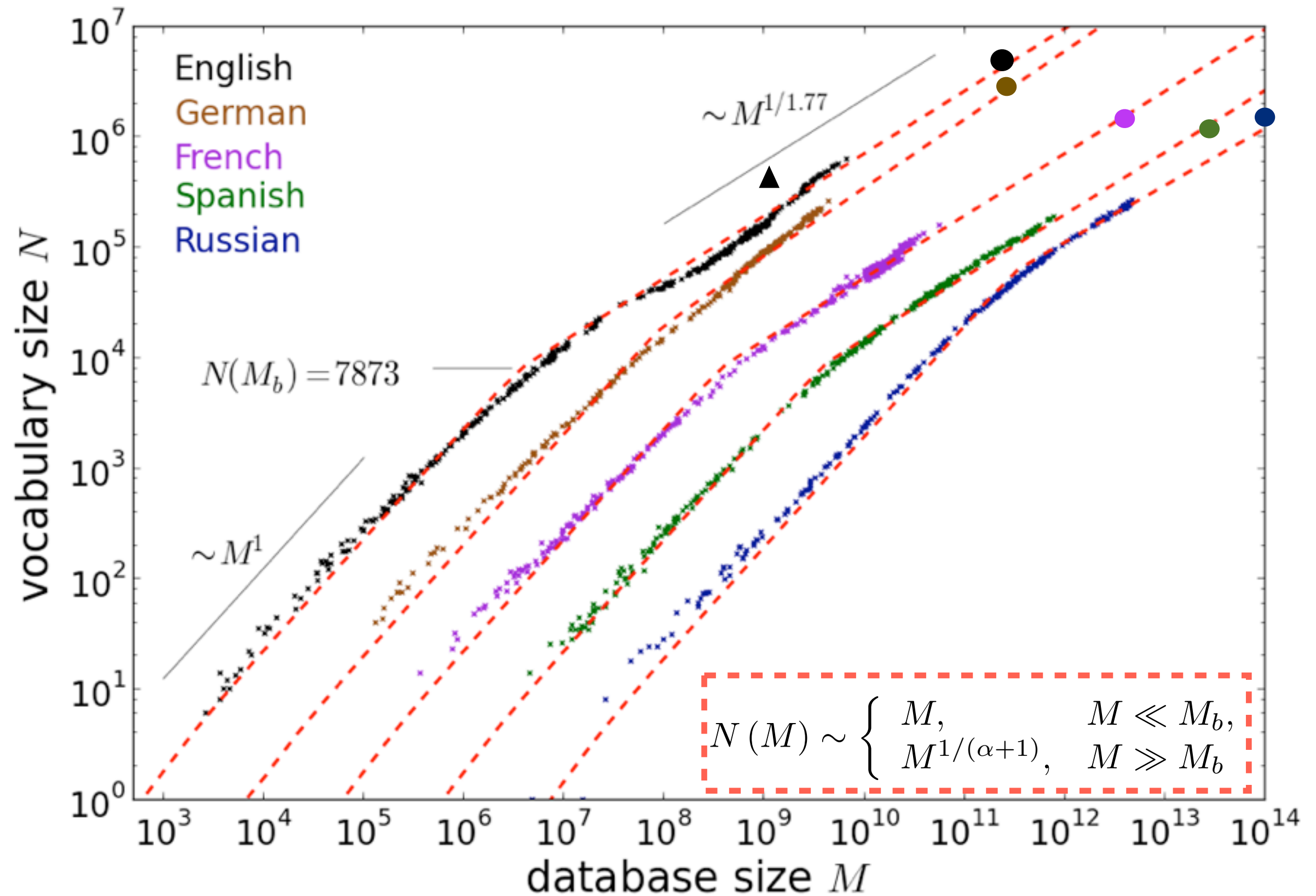
$$\langle N(M) \rangle = \sum_r 1 - e^{-F(r)M}$$

where  $F(r)$  is the frequency of the  $r$ -th most frequent word ( $r = \text{rank}$ ).

$$F_{dp}(r; \gamma, b) = \begin{cases} r^{-1}, & r \leq b, \\ r^{-\gamma} & r > b \end{cases}$$

$$N_{dp}(N_c) = \begin{cases} M, & M \ll M_b, \\ M^{1/\gamma}, & M \gg M_b \end{cases}$$

# Descriptive model



# Generative model (Yule-Simon type)

$M$ -th word in the database

Is it a new word?

$p_{new}$

Yes No

$1 - p_{new}$

Is it a core-word?

$p_c$

Yes No

$1 - p_c$

Choose a previous word, proportional to freq.

$N_c \mapsto N_c + 1$

$N_{\bar{c}} \mapsto N_{\bar{c}} + 1$

Assumptions:

1. Core vocabulary is finite:

$$N_c \leq N_c^{max} \Rightarrow p_c \rightarrow 0$$

2.  $p_{new}$  decays with  $N$ :

$$p_{new} \mapsto p_{new} (1 - \alpha/N)$$

$$N(M) \sim \begin{cases} M, & M \ll M_b, \\ M^{1/(\alpha+1)}, & M \gg M_b \end{cases}$$

with  $N(M_b) = N_c^{max}$



## Plan:

1. Vocabulary Growth

Centuries / millions of books



➔ 2. Innovations and Change



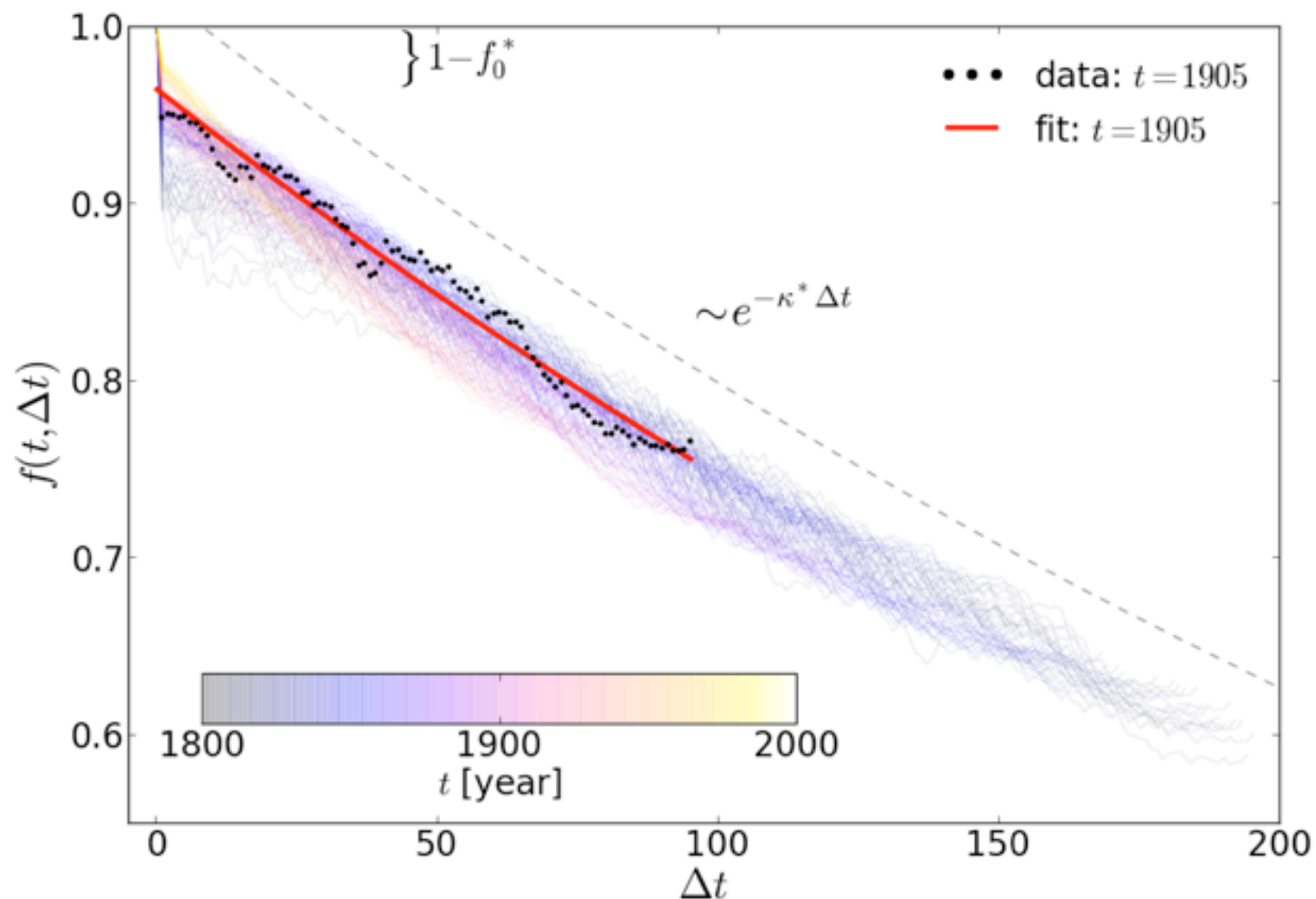
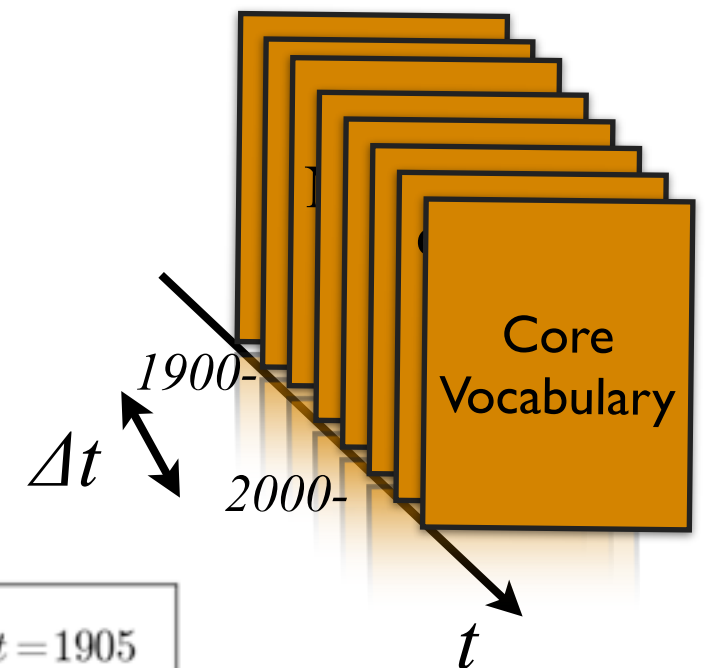
3. Text Analysis

Single book

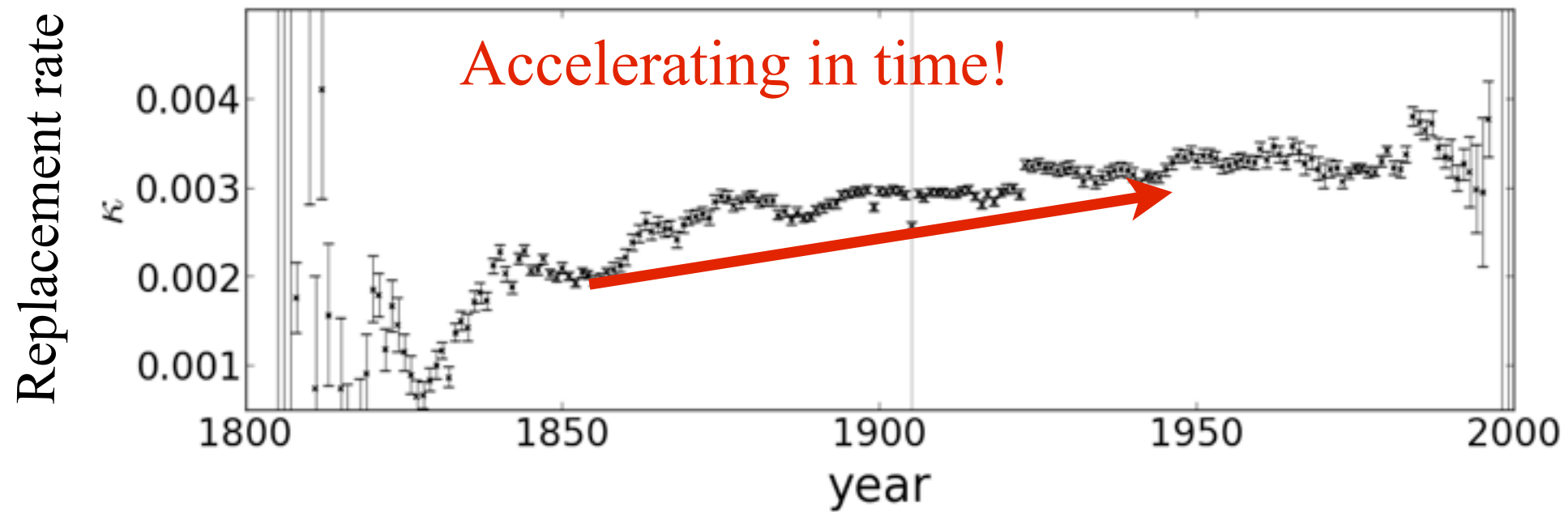


# Change in the core vocabulary

$f(t, \Delta t)$ : fraction of core words at time  $t$   
which remain core at time  $t + \Delta t$



# Change in the core vocabulary



# Change in the core vocabulary

**1900**

*majesty, doubtless,  
furnished, monsieur,  
Napoleon, hitherto*

Most frequent  
replaced words

**2000**

*cultural, context,  
technology, programs,  
environmental, computer*

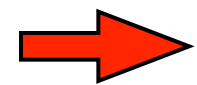
## Plan:

1. Vocabulary Growth

Centuries / millions of books



2. Innovations and Change



3. Text Analysis

Single book



E. G. Altmann, **Giampaolo Cristadoro, Mirko Degli Esposti**, "*On the origin of long-range correlations in texts*", PNAS 109, 11582 (2012)

# Observation

Symbolic sequence  $s$

War and Peace, by Leo Tolstoy  
BOOK ONE: 1805  
CHAPTER I

"Well, Prince, so Genoa and Lucca are now just family estates of the Buonapartes. But I warn you, if you don't tell me that this means war, if you still try to defend the infamies and horrors perpetrated by that Antichrist--I really believe he is Antichrist--I will have nothing more to do with you and you are no longer my friend, no longer my 'faithful slave,' as you call yourself! But how do you do? I see I have frightened you--sit down and tell me all the news."

...

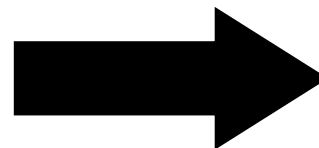
Numeric sequence

01001000011010000011010100010  
0110010100000  
001001000

01000001001001001011010000100101010010  
01000010100010010100100001001101010010001001  
01000011000001100100001000010001000010  
011000010010001100010000000010010100001  
100101100100001001000010010010100000010  
1001000100000011000010110100101001001000100001  
010000101000100100001010010010010000110100  
0110101001001001000000011000010010010000  
01100010001010100011010000110010000100010  
0100110010101100001010001000101000110010  
01000100001000010100001001000

...

$f(s)$



text  $s$ =War and Peace by Leo Tolstoy Book one 1805 Chapter 1 Well Prince

vowels=010010001101000001100100010011001000000000001001000001000001001...

t

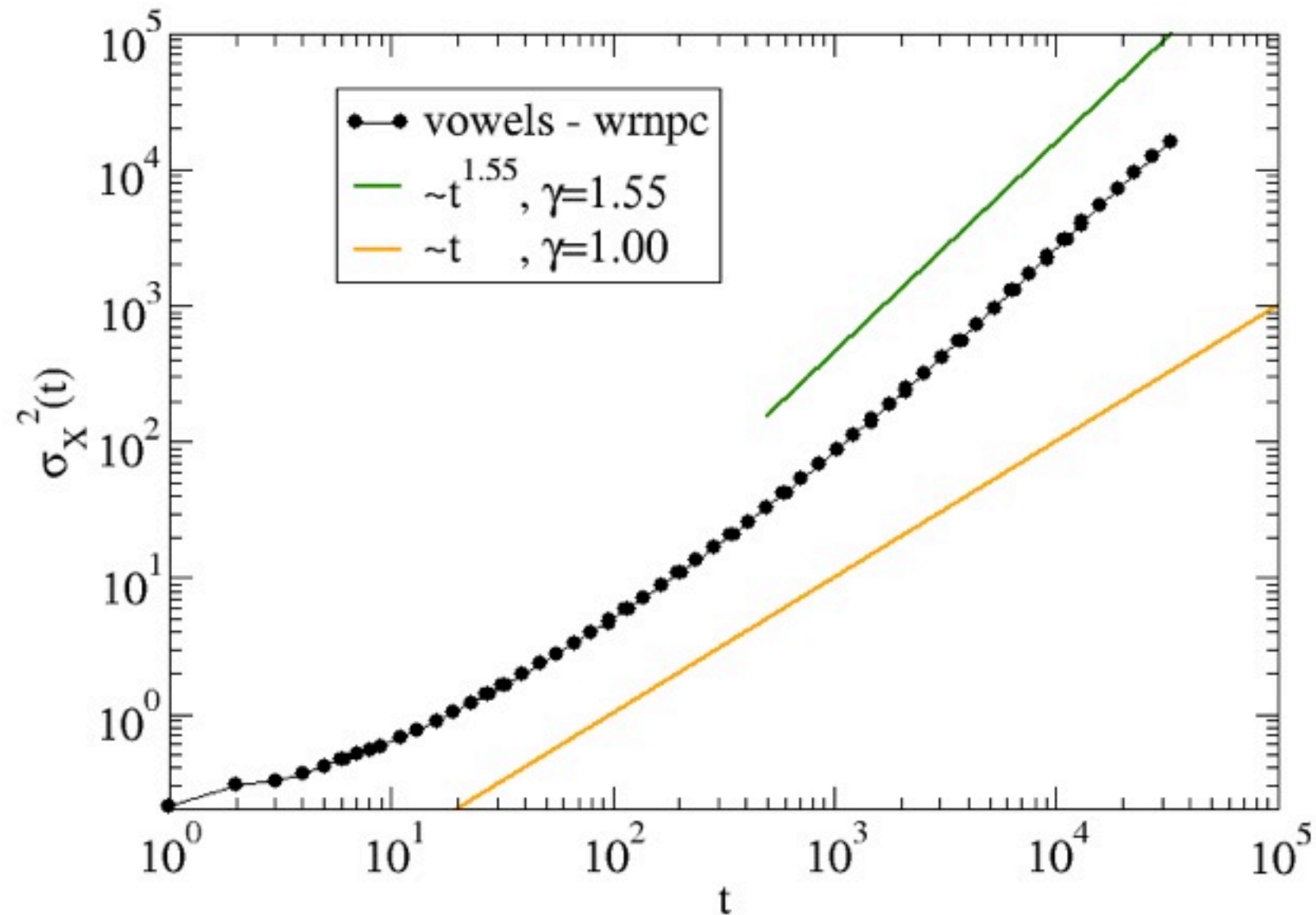
Random walk  $X(t)$ :

0 => left

1 => right

$$\sigma_X^2(t) := \langle X(t)^2 \rangle - \langle X(t) \rangle^2$$

# Transport and long correlations



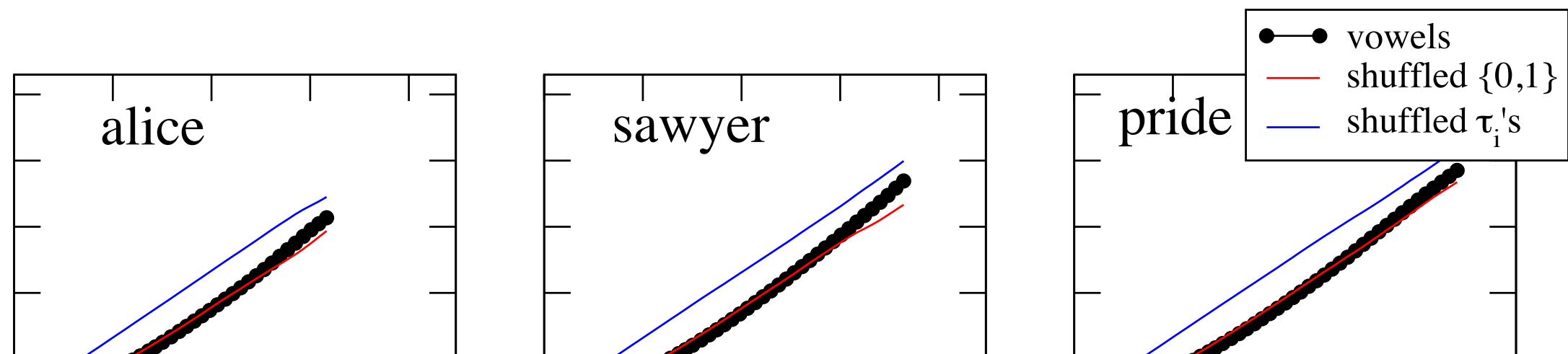
Super-diffusion:  $\sigma_X^2(t) \simeq t^\gamma$ ,  $1 < \gamma < 2$

Long-range correlation:  $C_{\mathbf{x}}(t) \simeq t^{-\beta}$ ,  $0 < \beta < 1$

$$\gamma = 2 - \beta$$

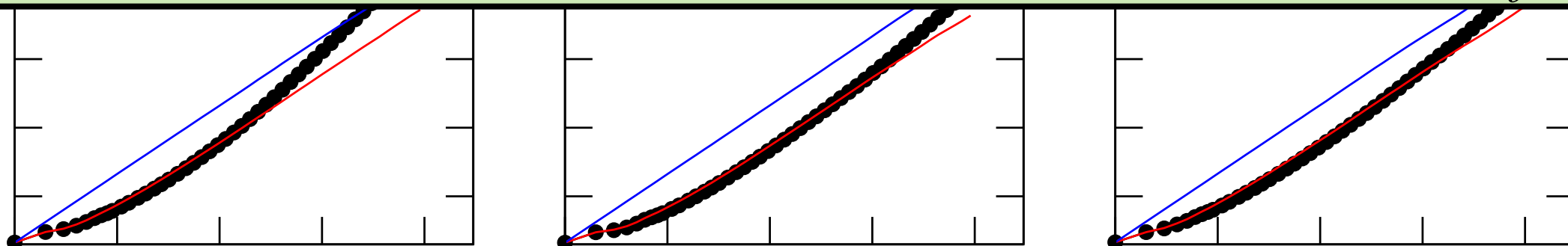


# *Different books*



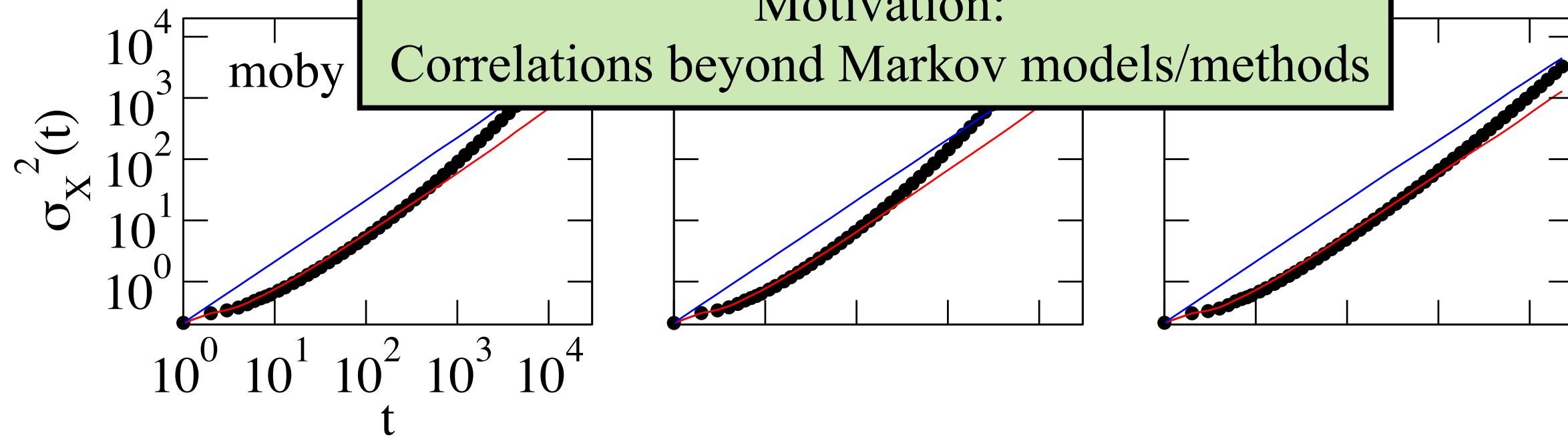
Basic questions are still open:

- What is the origin of the long-range correlation?
- How is it connected to the semantics of the text (story)?
- What is the role of the observable  $f$ ?



Motivation:

Correlations beyond Markov models/methods



# Burstiness and correlation

Inter-event times:  $\mathbf{x} = \textcolor{red}{1}00\textcolor{red}{1}000\textcolor{red}{1}10\textcolor{red}{1}\dots$   
 $\tau = 4$

$$\tau_k = 3, 4, 1, 2, \dots$$

Transport

$$\sigma_X^2(t) := \langle X(t)^2 \rangle - \langle X(t) \rangle^2 \simeq t^\gamma,$$

$$\gamma = 1$$

$$1 < \gamma \leq 2$$

short range

long range

Power spectrum

$$S(\omega) := \int_{-\infty}^{+\infty} C_{\mathbf{x}}(t) e^{2\pi i t \omega} dt$$

$$S(0) = \frac{\sigma_\tau^2}{\langle \tau \rangle^3} \left( 1 + 2 \sum_k C_\tau(k) \right)$$

$$S(0) \rightarrow \infty$$

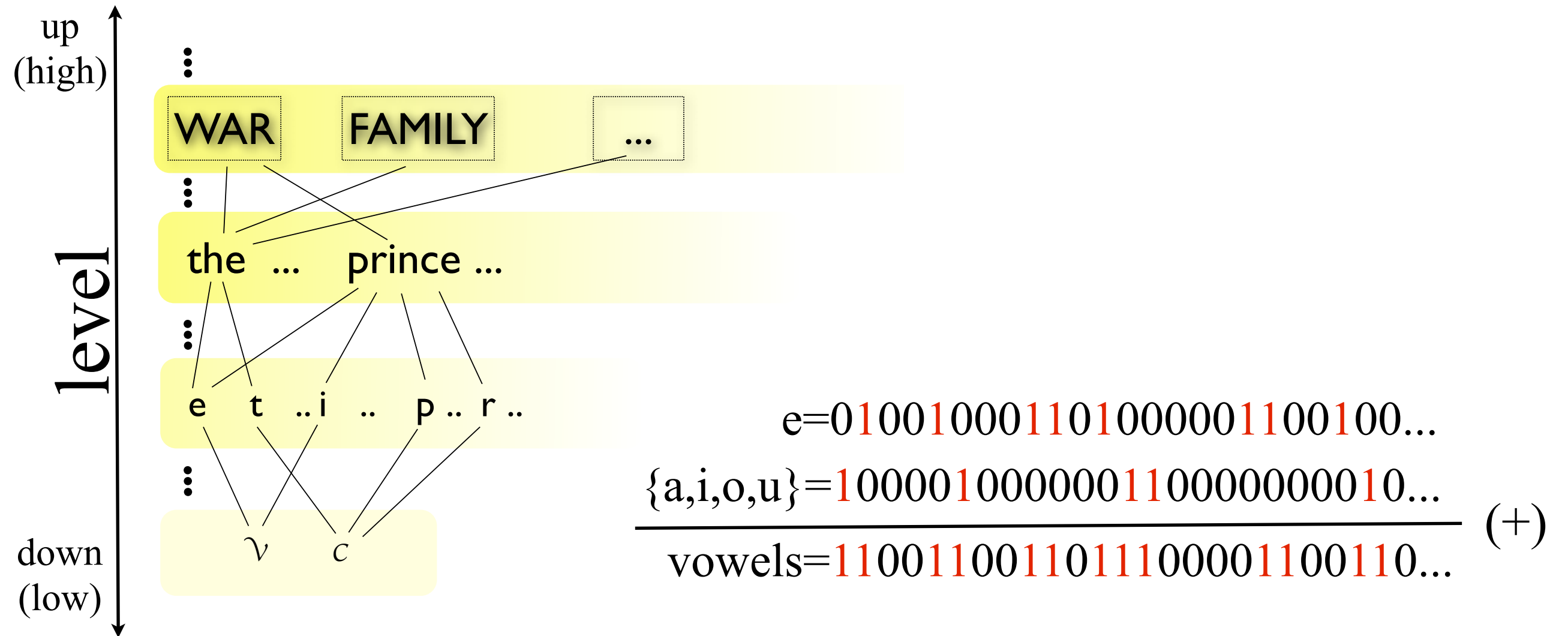
Burstiness

$$\frac{\sigma_\tau}{\langle \tau \rangle} \rightarrow \infty$$

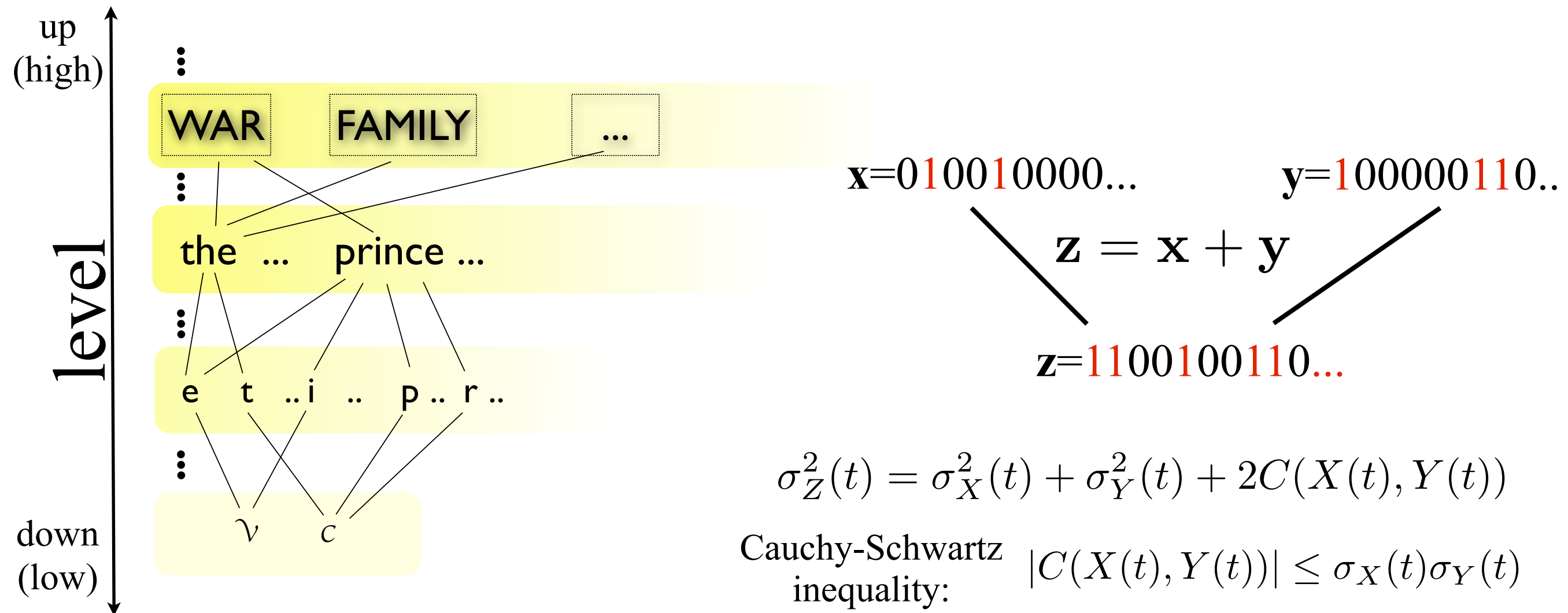
Correlation in  $\tau_k$

$$\sum_{k=1}^{\infty} C_\tau(k) \rightarrow \infty$$

# Model

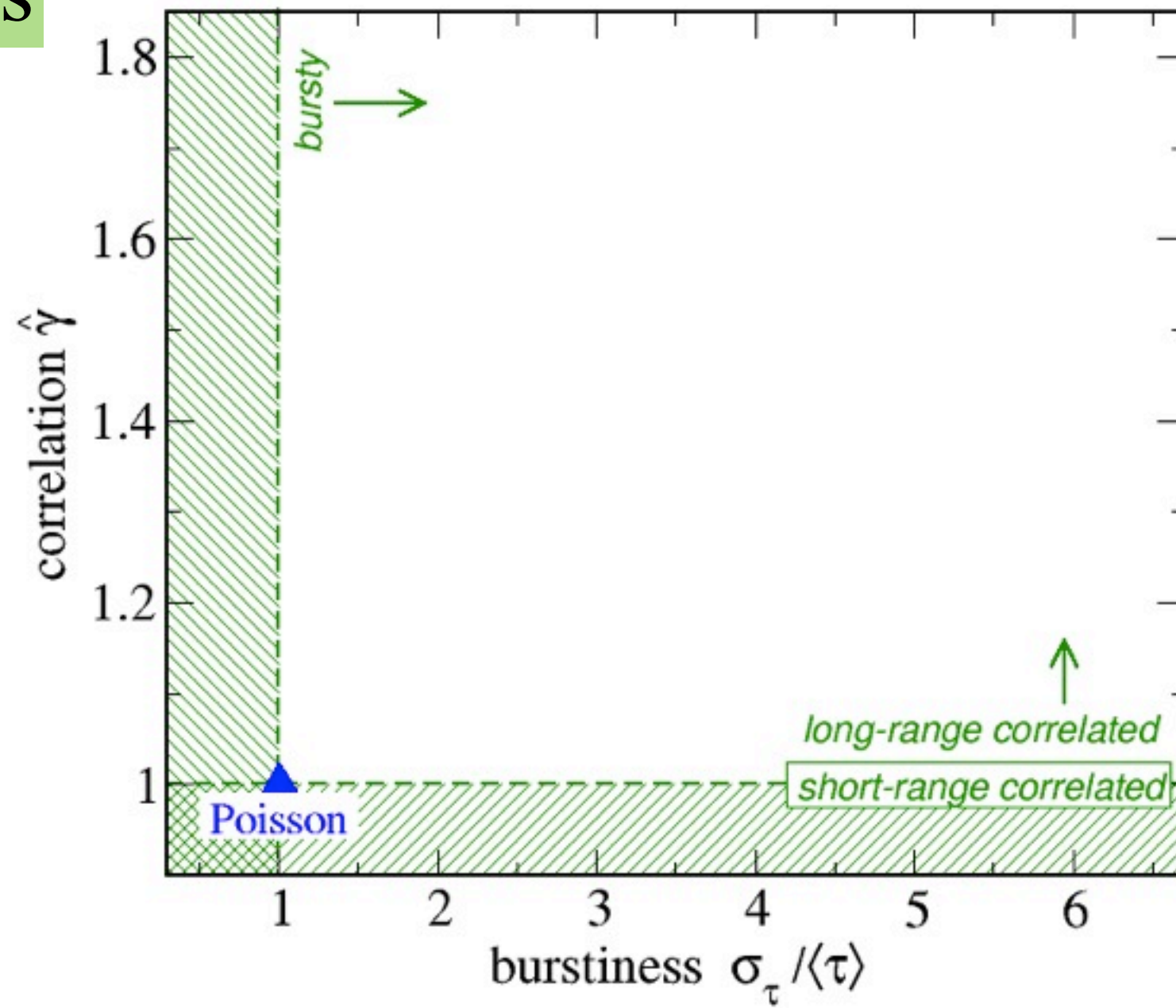


# Model

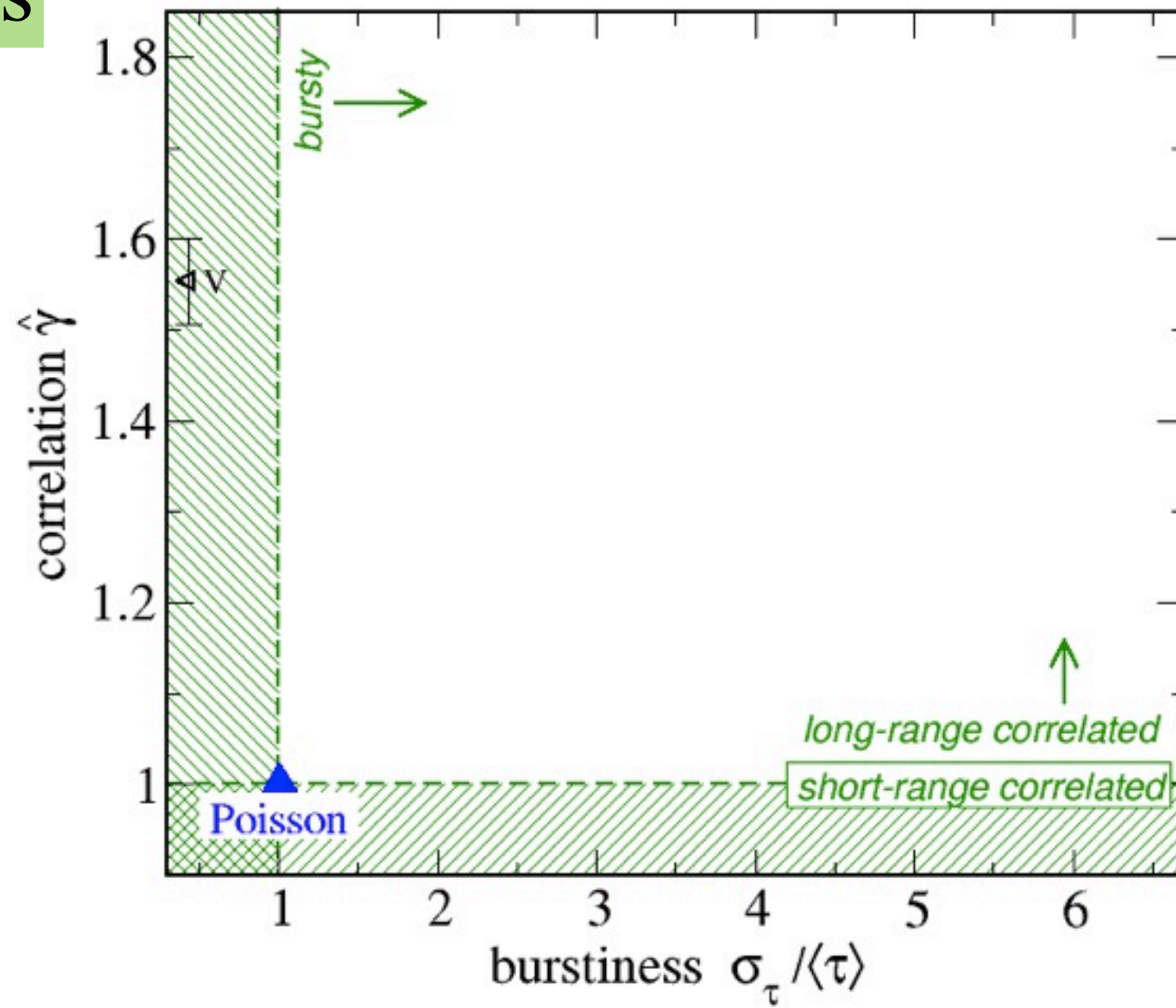


	Long-range correlation
Moving up:	necessarily preserved

# Observations

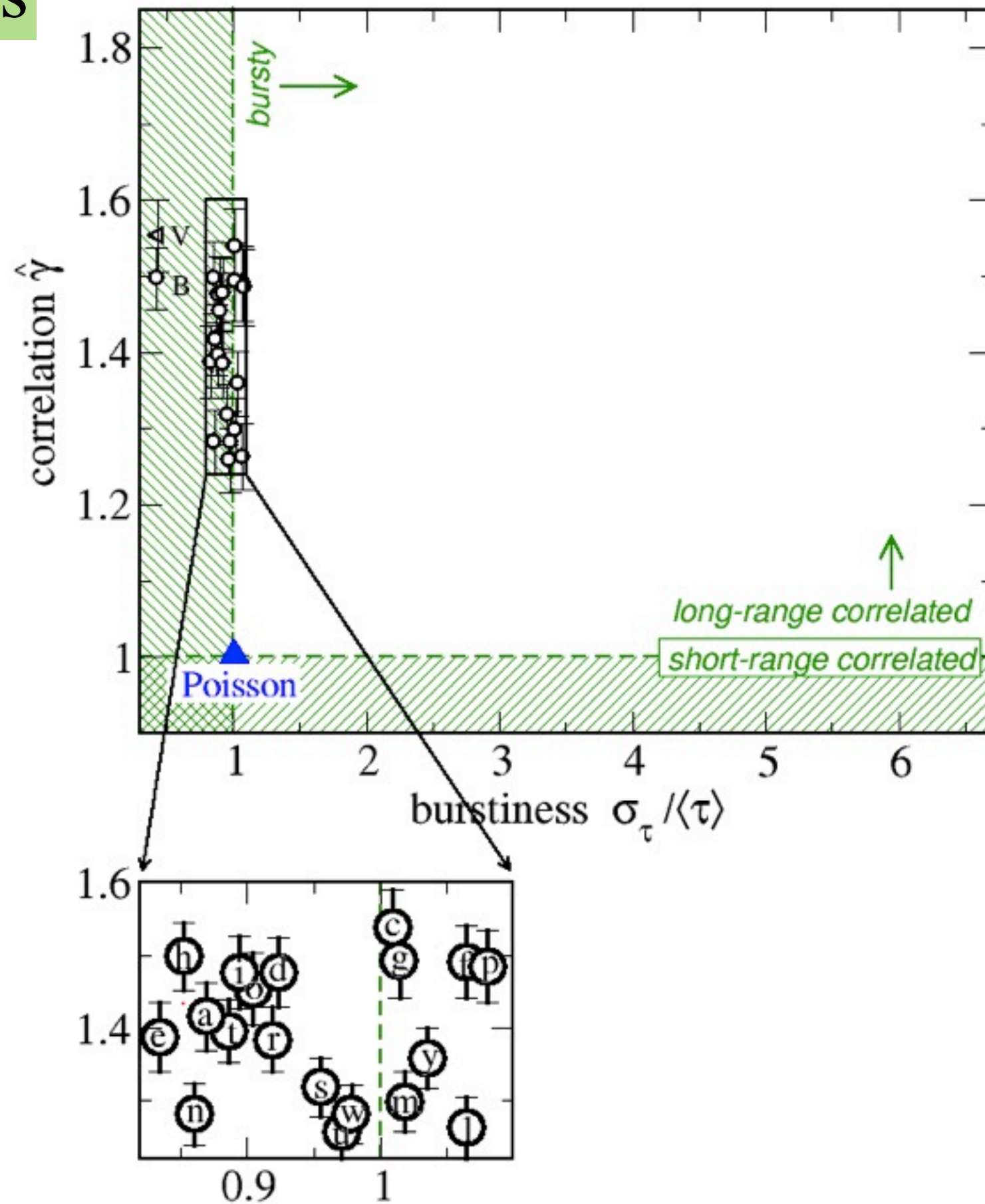


# Observations



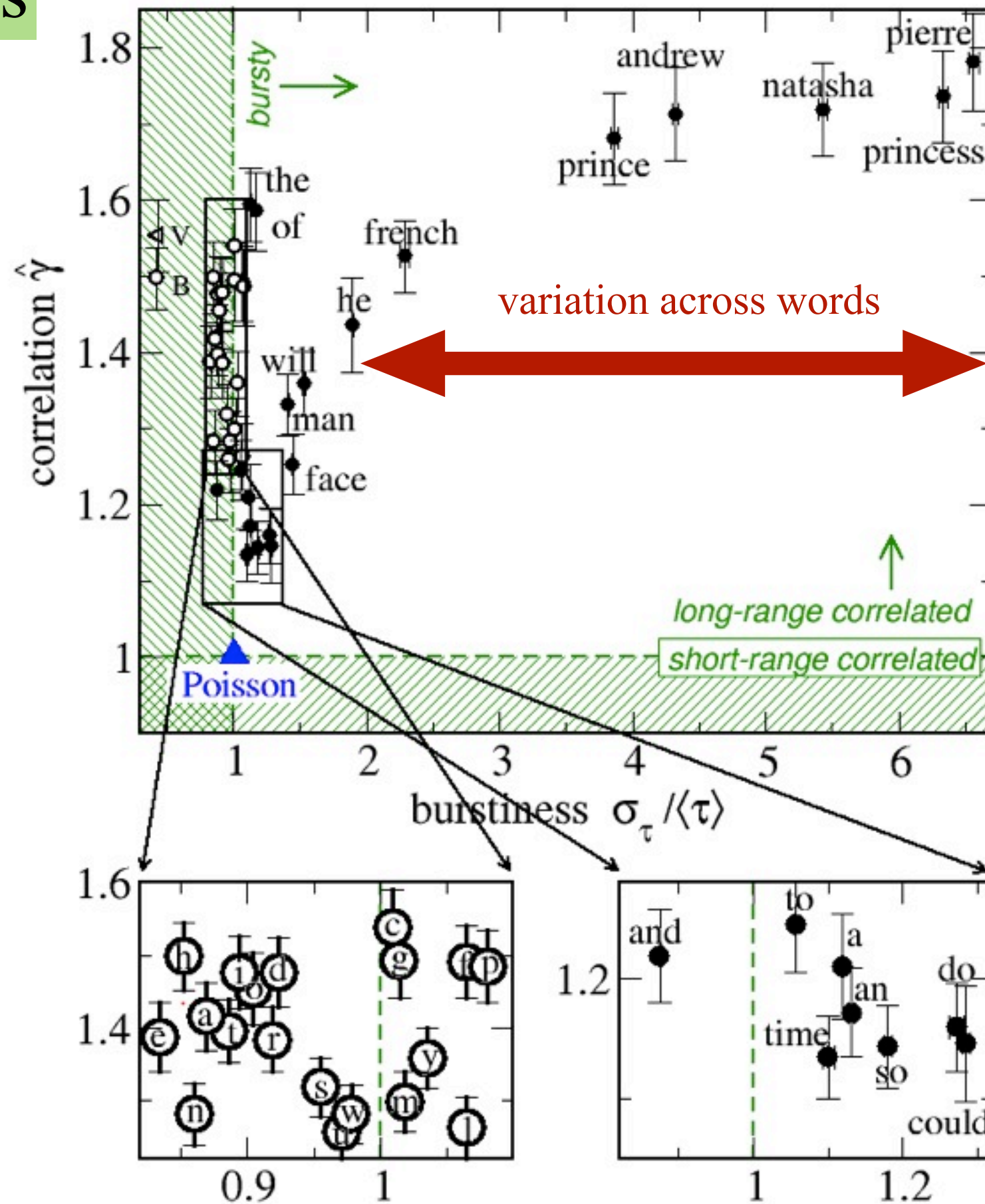


# Observations





# Observations



# Burstiness and word types

Words with similar frequency:  $\langle \tau \rangle \approx 820$

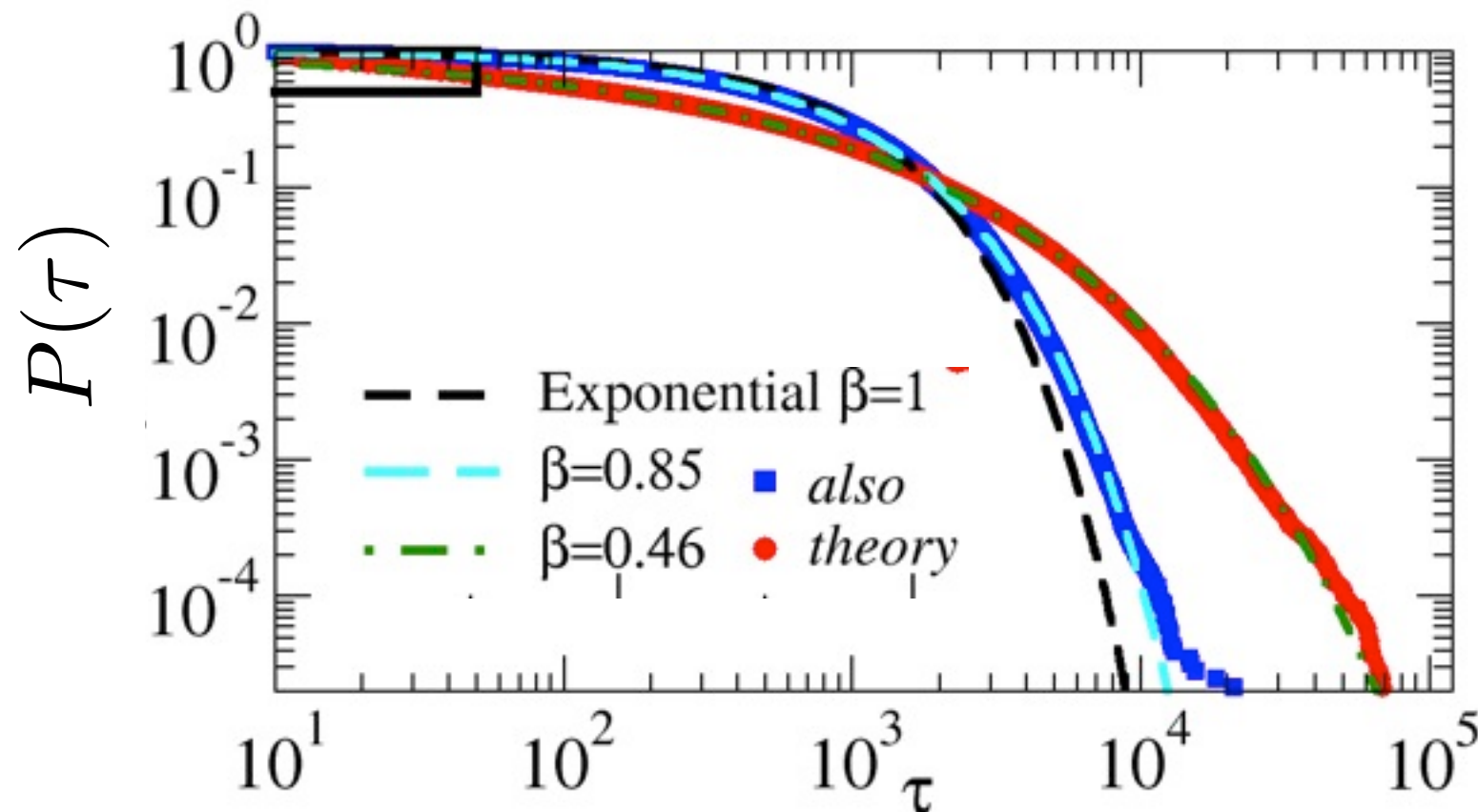
word: *theory*



word: *also*



Poisson/random



Weibull distribution

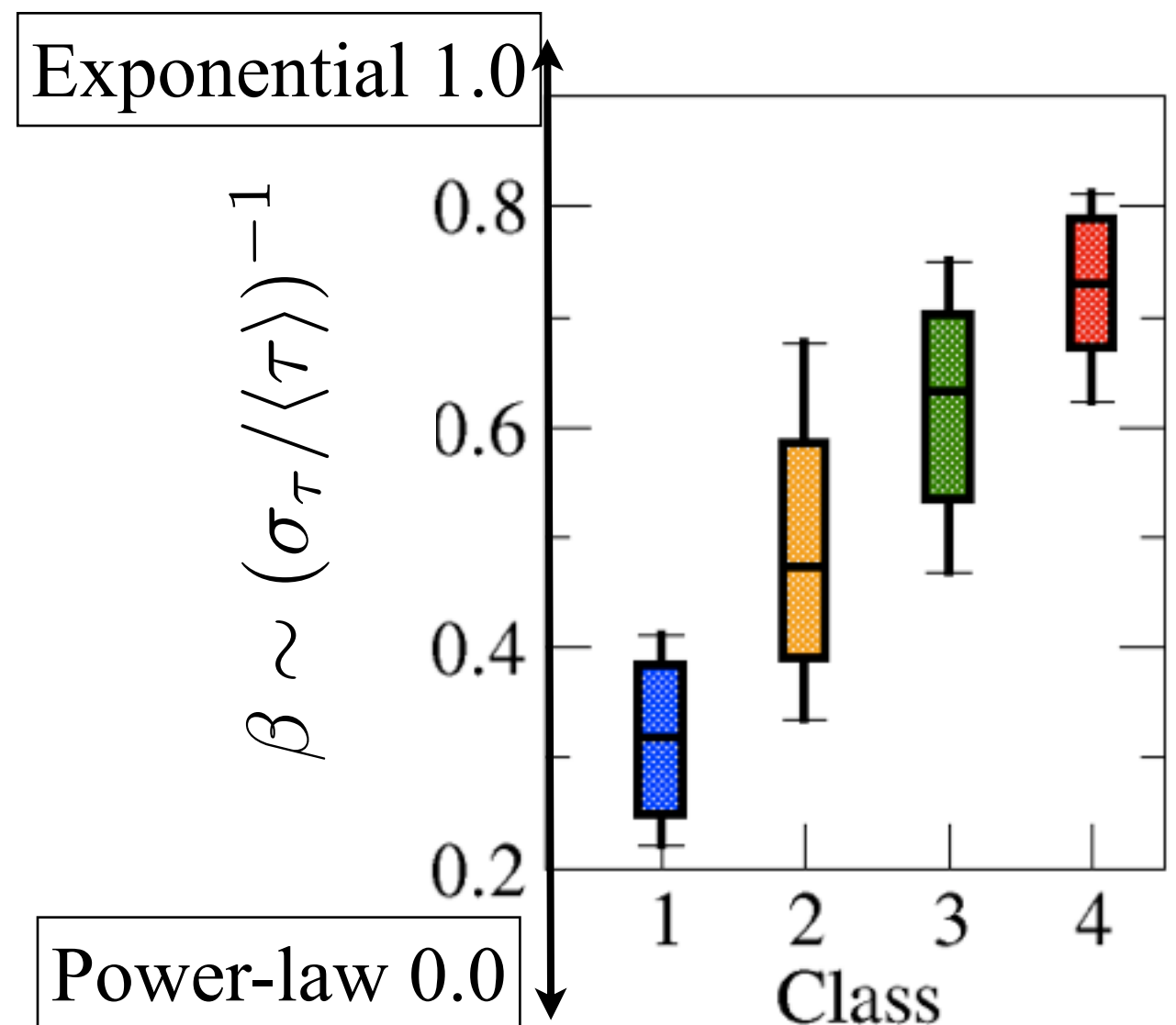
$$P_{\beta}(\tau) = e^{-a\tau^{\beta}}$$

# Burstiness and word types

Semantic Classes

Class	Name	Examples of words
1	Entities	Africa, Bible, Darwin prince e,t
2	Predicates and Relations	blue, die, in, religion theory <e,t>
3	Modifiers and Operators	believe, everyone, forty also <<e,t>,t>
4	Higher Level Operators	hence, let, supposedly the <<<>,>>...

Syntactic (verbs, adjectives, nouns, etc..) relations act on shorter time scales:  
 $\langle \tau \rangle \gg$  sentence length





# Keywords of unknown texts: the Voynich Manuscript

A botanical illustration of a plant, likely a species of Euphorbia, shown from a top-down perspective. The plant has a thick, brown, fibrous root system that spreads out at the base. Three main stems emerge from the center of the root system. The central stem is upright and bears a dense, elongated cluster of small, green, oval-shaped leaves. The two side stems are more branched and bear clusters of small, red, star-shaped flowers. The plant is drawn on aged, slightly textured paper.



# Keywords of unknown texts: the Voynich Manuscript

## Keywords

### New Testament

Portuguese	English	German	Voynich
nasceu	begat	zeugete	cthy
Pilatos	Pilates	zentner	qokeedy
céus	talents	himmelreich	shedy
bem-aventurados	loaves	pilatus	qokain
Isabel	Herod	schwert	chor
anjo	tares	Maria	lkaiin
menino	vineyard	Elisabeth	qol
vinha	shall	Etliches	lchedy
sumo	boat	unkraut	sho
sepulcro	demons	euch	qokaiin
joio	five	schiff	olkeedy
Maria	pay	ihn	qokal
portanto	sabbath	weden	qotain
Herodes	hear	heuchler	dchor
talentos	whosoever	tempel	otedy

D. R. Amancio, E. G. Altmann, D. Rybski, O. N. Oliveira Jr., L. da Fontoura Costa, "*Probing the statistical properties of unknown texts: application to the Voynich Manuscript*", *PLoS ONE* 8, e67310 (2013)

# Summary of conclusions

Thank you,  
for your attention!

## 1. Vocabulary Growth

$$N_{dp}(N_c) = \begin{cases} M, & M \ll M_b, \\ M^{1/\gamma}, & M \gg M_b \end{cases}$$

## 2. Innovations and Change

- accelerated core-vocabulary replacement
- quantification of exogenous/endog. factors

## 3. Text Analysis

Burstiness



Long-range correlation

- \* **Martin Gerlach**, José M. Miotto, Fakhreh Ghanbarnejad (MPIPKS, Dresden)
- \* **Giampaolo Cristadoro**, Mirko Degli Esposti (Univ. Bologna)
- \* Adilson E. Motter, Janet Pierrehumbert (Northwestern Univ., USA)
- \* Diego R. Amancio, Osvaldo N. Oliveira Jr., Luciano da Fontoura (USP, Brazil)
- \* Diego Rybski (PIK, Potsdam)

## References:

- M. Gerlach, E. G. Altmann, *Stochastic model for the vocabulary growth in natural languages*, Phys. Rev. X 3, 021006 (2013)
- D. R. Amancio, E. G. Altmann, D. Rybski, O. N. Oliveira Jr., L. da Fontoura Costa, "Probing the statistical properties of unknown texts: application to the Voynich Manuscript", PLoS ONE 8, e67310 (2013)
- E. G. Altmann, G. Cristadoro, M. Degli Esposti, "On the origin of long-range correlations in texts", PNAS 109, 11582 (2012)
- E. G. Altmann, J. B. Pierrehumbert, A. E. Motter, "Beyond word frequency: Bursts, lulls, and scaling in the temporal distributions of words", PLoS ONE 4 (11) e7678 (2009)