# Information Rate and Branching Processes of Scientific Fields

FranČesko, Henry, Jan, MJ
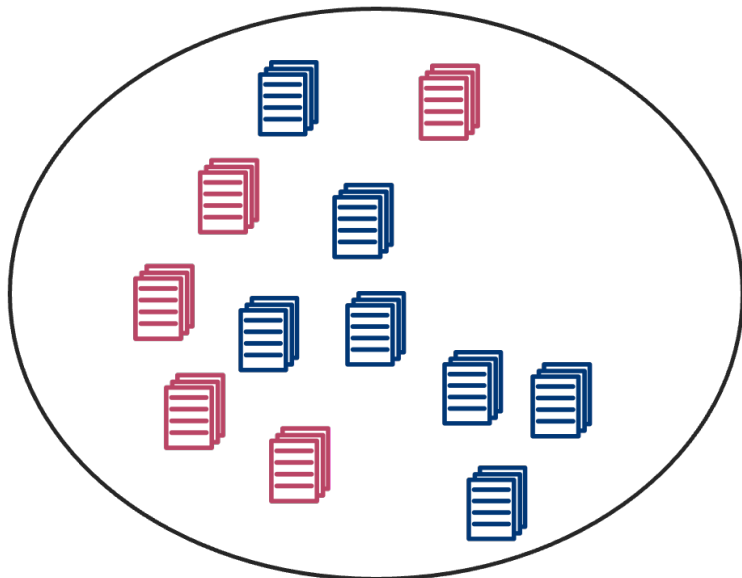
(Advisory Board: Jana)

# *Do Societies have a limited Information processing rate?*

For Sciences:

Is there a correlation between reaching a certain entropy-rate and the emergence of new subfields?
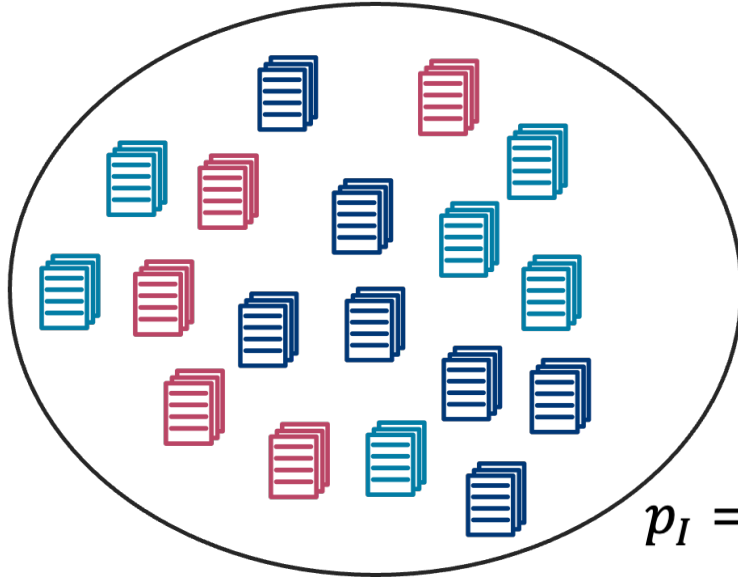
# Field A



Topic I

Topic II

Year i

$$p_I = \frac{5}{12} \quad p_{II} = \frac{7}{12}$$

$$H = \sum_i p_i \log p_i = \frac{5}{12} \log \frac{5}{12} + \frac{7}{12} \log \frac{7}{12}$$

$$= 0.68$$

# Field A
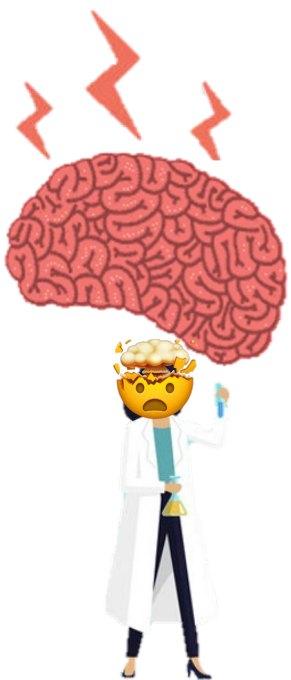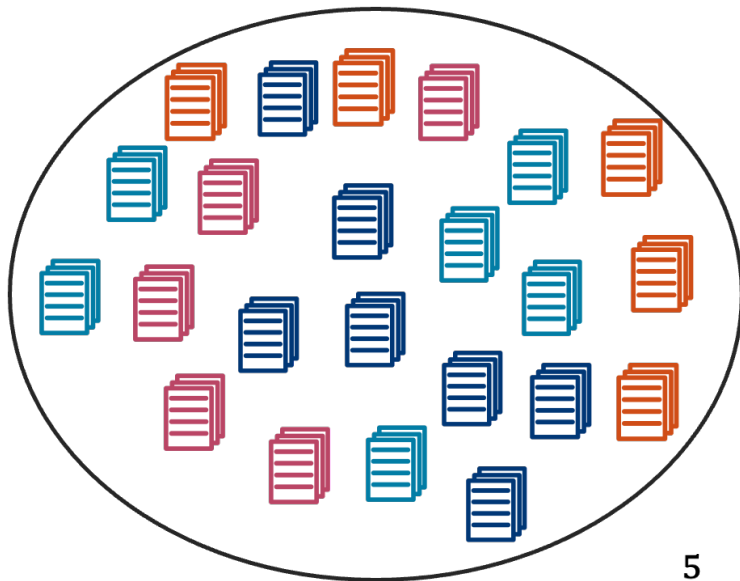


Topic I
Topic II
Topic III

$$p_I = \frac{5}{18} \quad p_{II} = \frac{7}{18} \quad p_{III} = \frac{6}{18}$$

Year i+1

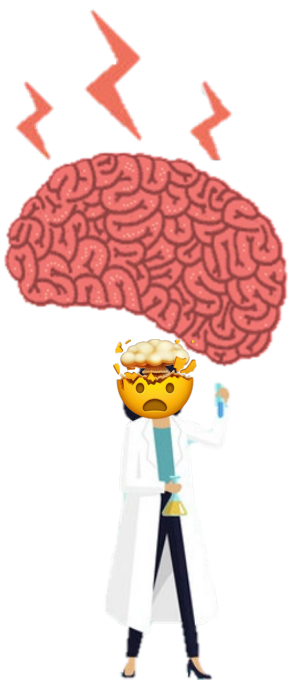$$H = \frac{5}{18}\log\frac{5}{18} + \frac{7}{18}\log\frac{7}{18} + \frac{6}{18}\log\frac{6}{18} =$$
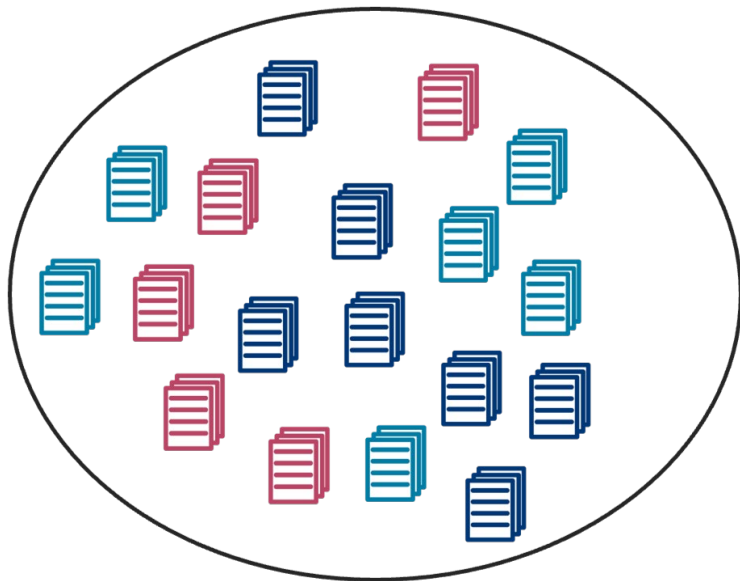
$$= 1$$

**Field A**

Topic I
Topic II
Topic III
Topic IV

Year i+2

$$p_I = \frac{5}{23} \quad p_{II} = \frac{7}{23} \quad p_{III} = \frac{6}{23} \quad p_{IV} = \frac{5}{23}$$

$$H = \frac{5}{23}\log\frac{5}{23} + \frac{7}{23}\log\frac{7}{23} + \frac{6}{23}\log\frac{6}{23} + \frac{5}{23}\log\frac{5}{23} =$$
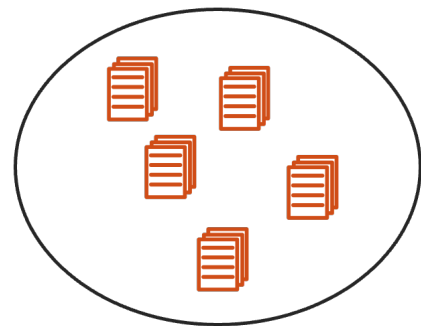
$$= 1.3$$

# Field A- Subfield 1

Topic I

Topic II

Topic III

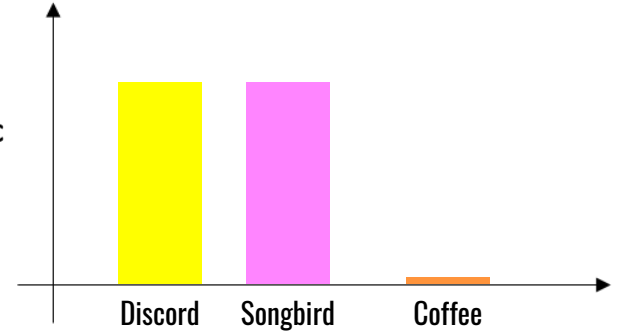Topic IV

Year i+2

# Field A- Subfield 2

# Method

Constantly **checking Discord** has become second nature, like a **songbird** warily scanning the skies for signs of danger before taking flight.
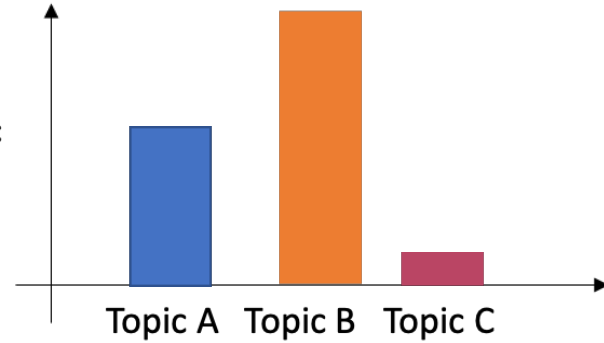
(ChatGPT)

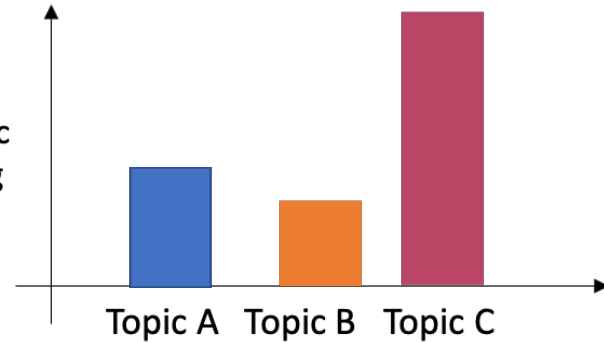→ Probabilistic Embedding

# Method



Paper i

Word Embedding
+ Topic Modelling

Probabilistic
Embedding
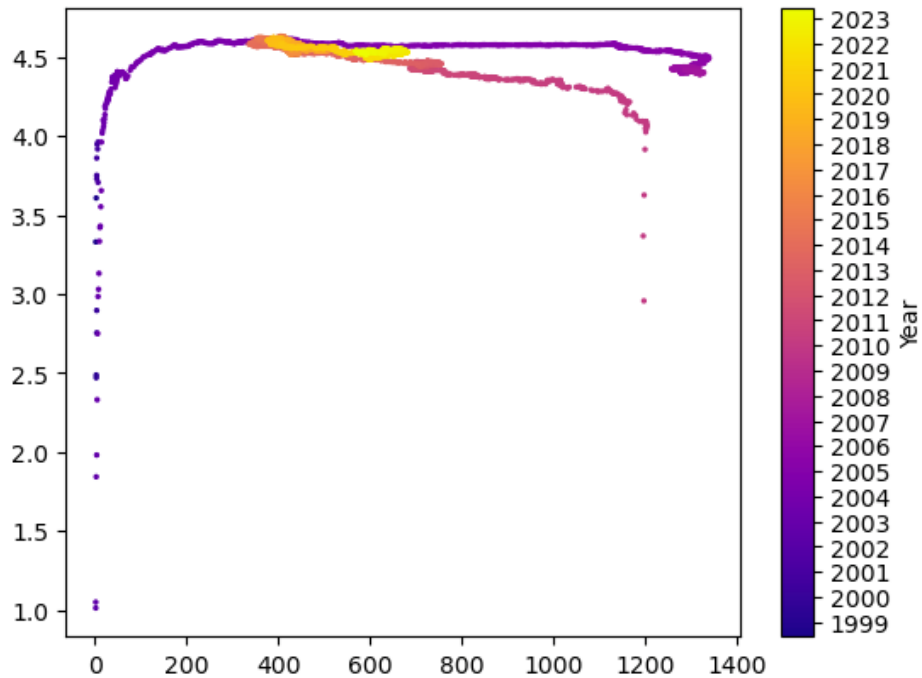
Topic A   Topic B   Topic C

Paper j

Probabilistic
Embedding

Topic A   Topic B   Topic C

# Results

# Entropy rates

# Entropy rates

# What does that tell us? 🙈🙈🙈



physics.acc-ph

q-bio.PE

# Branching process

## Speciation



[Morris, et al., 2023]

# Branching process

## Speciation



[Morris, et al., 2023]

## Emergence of scientific fields


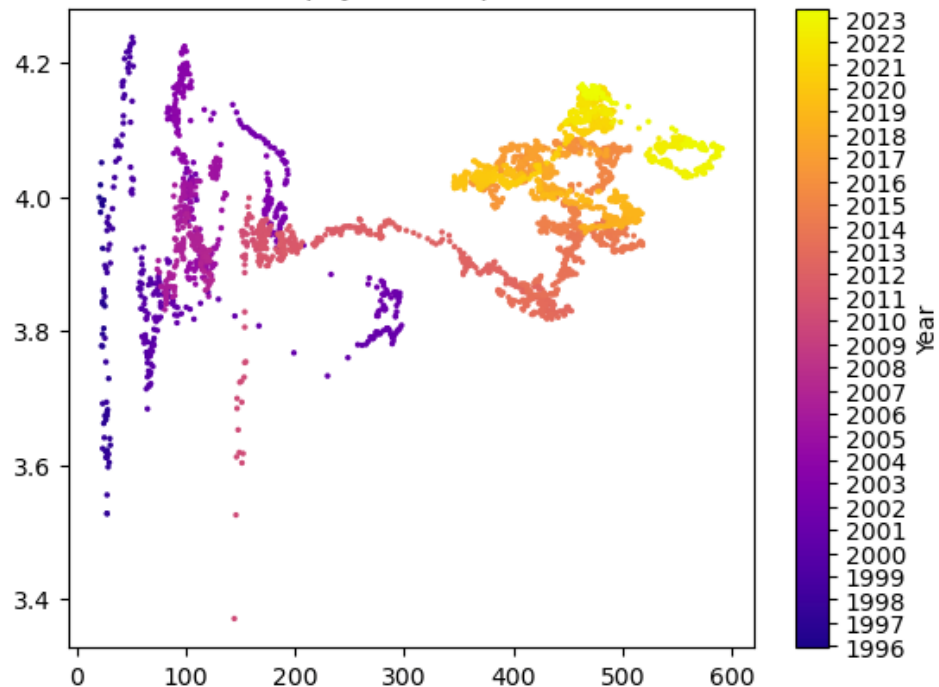
- Computer Science
- Economics
- Electrical Engineering
- Mathematics
- Quantitative Biology
- Quantitative Finance
- Statistics
- Physics
- Other

1986 - 2023

# Branching process



1986 - 2008

- ● Computer Science
- ● Economics
- ● Electrical Engineering
- ● Mathematics
- ● Quantitative Biology
- ● Quantitative Finance
- ● Statistics
- ● Physics
- ● Other

Time (years)

# Branching process



1986 - 2008    2012

Time (years)

Computer Science
Economics
Electrical Engineering
Mathematics
Quantitative Biology
Quantitative Finance
Statistics
Physics
Other

# Branching process



1986 - 2008    2012    2015

Time (years)

- Computer Science
- Economics
- Electrical Engineering
- Mathematics
- Quantitative Biology
- Quantitative Finance
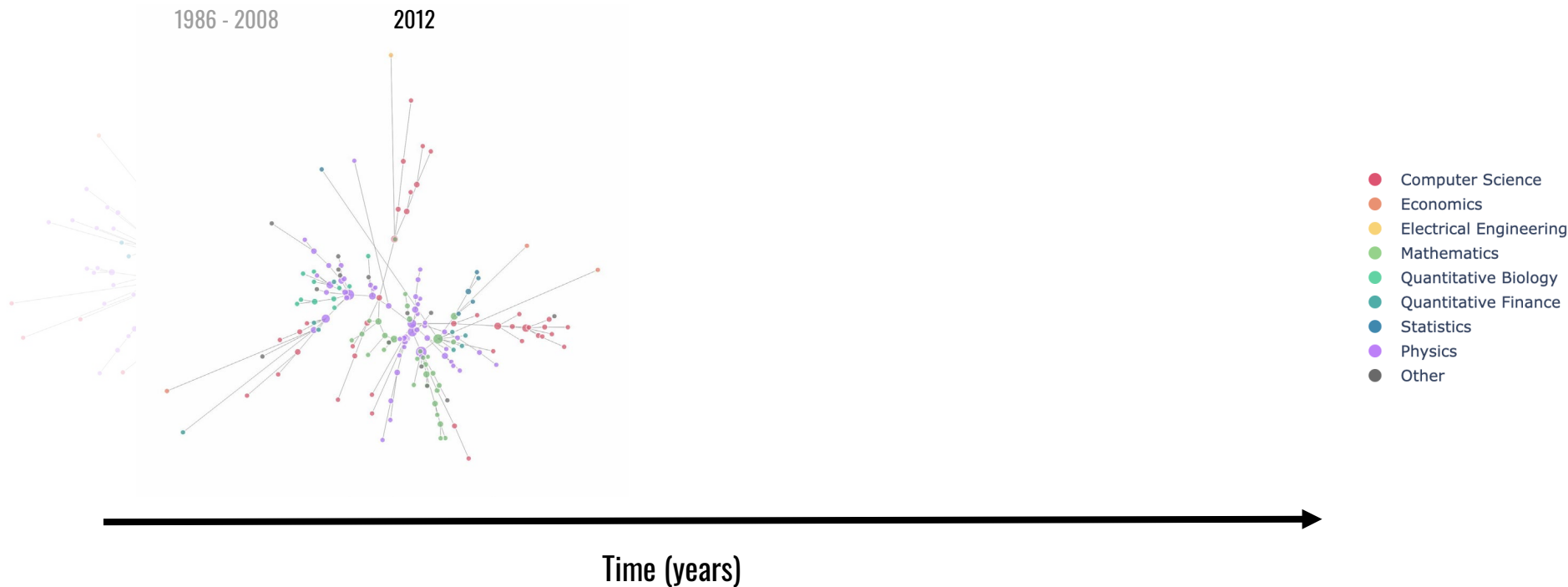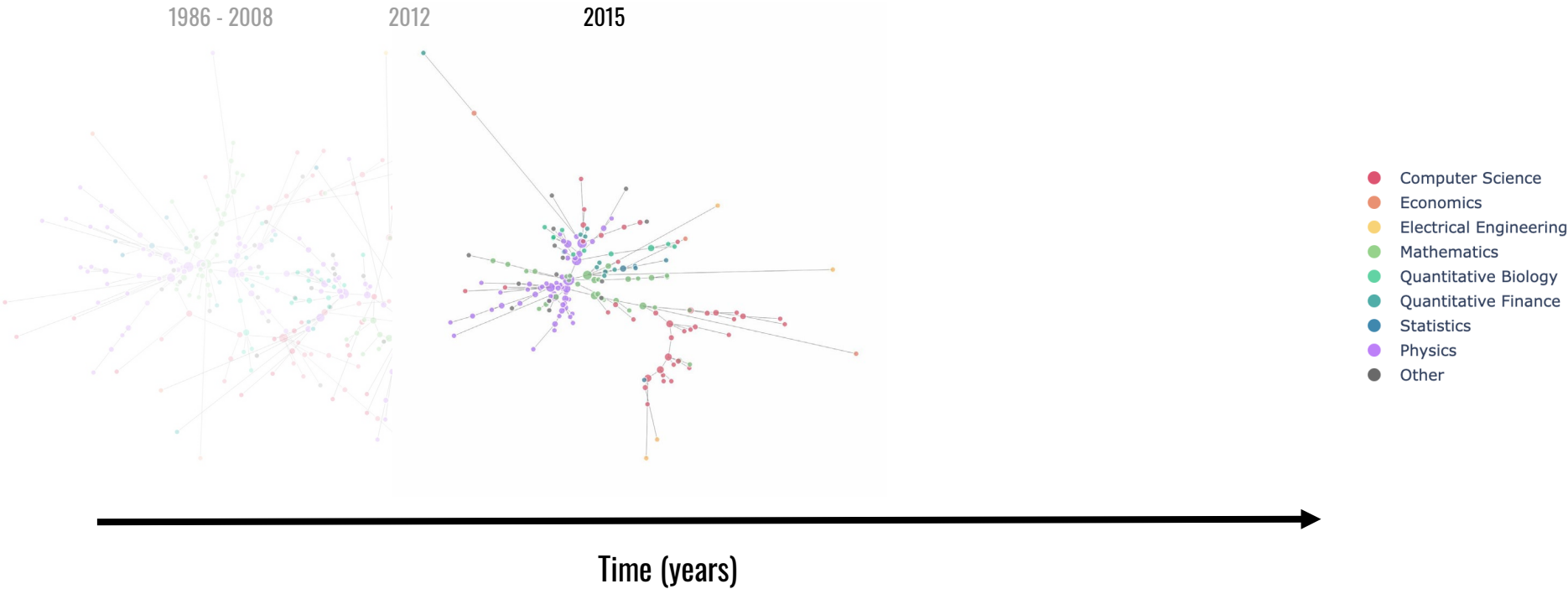- Statistics
- Physics
- Other

# Branching process

# Branching process

# Branching process



1986 - 2008    2012    2015    2018    2021    2023

- ● Computer Science
- ● Economics
- ● Electrical Engineering
- ● Mathematics
- ● Quantitative Biology
- ● Quantitative Finance
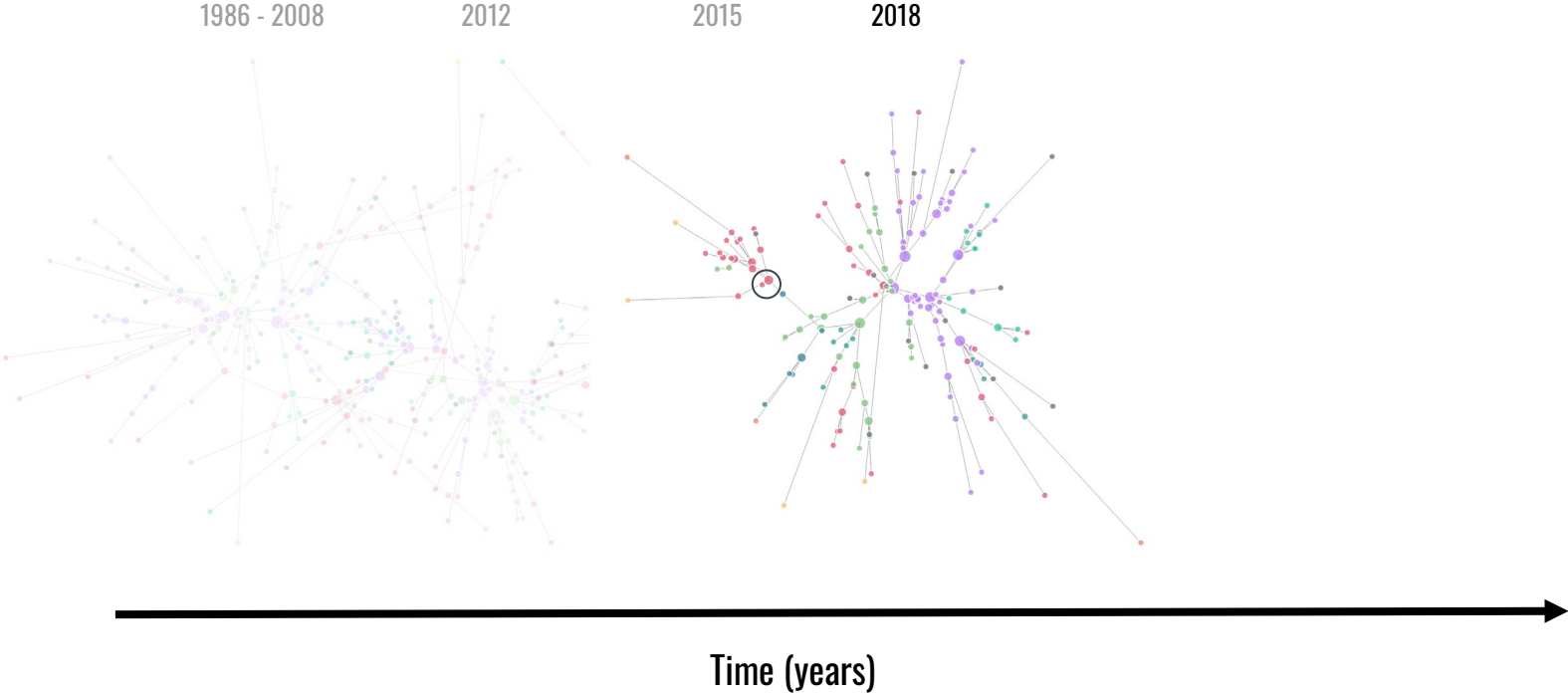- ● Statistics
- ● Physics
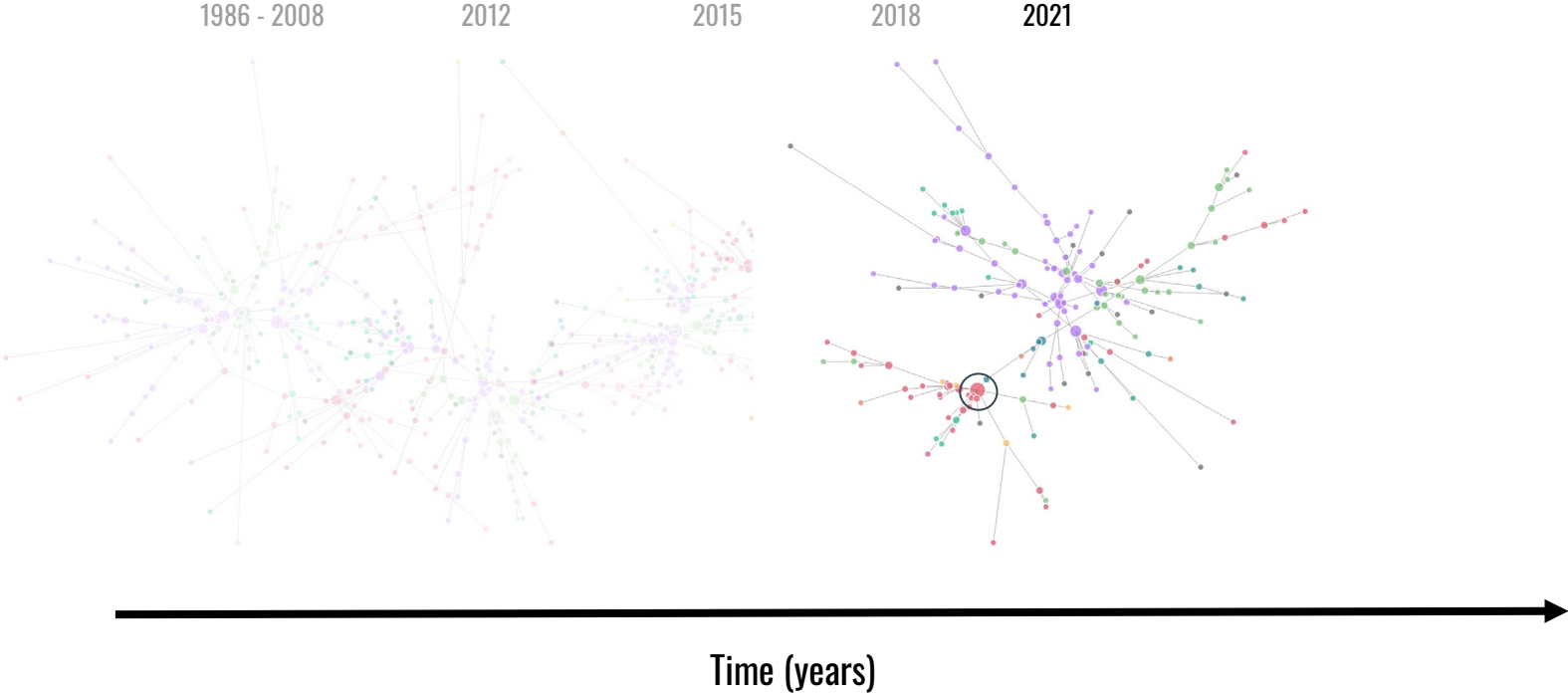- ● Other
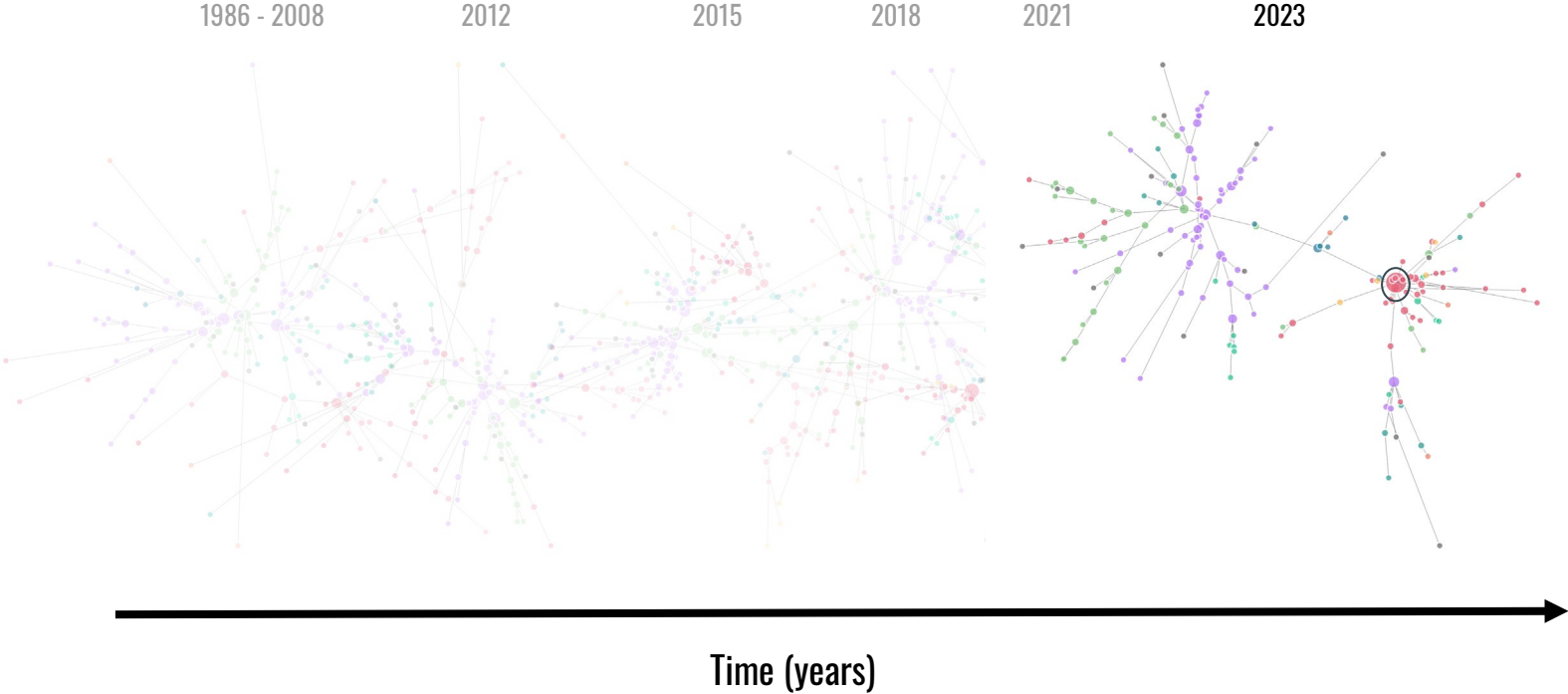
Time (years)

# Branching process

# Branching process

# Future work

- Correlate the entropy rate against other indicators (e.g.: number of distinct authors,...)
- Expand methodology to other data sets
  - Reddit (subreddit), Patent, and APS dataset

Hope you liked our presentation!

3 Ì !ɖÌ −!! .†Í π!ɧÊđ!Ì ð. .π˝!! µ†ð†!!π∅Ì Â Â πÊµ†ð∅̃Ê⥮!!

FIN

# Embedding

# Integrating topic modeling and word embedding to characterize violent deaths

Alina Arseniev-Koehler[a,b,1] , Susan D. Cochran[b,c,d] , Vickie M. Mays[b,e,f], Kai-Wei Chang[b,g], and Jacob G. Foster[a,b,1]

# Latent embedding space

## WORD2VEC EMBEDDING



## TOPIC EMBEDDING

$$\Pr[w \text{ emitted at } t|\mathbf{c}_t] = \frac{\exp\left(\langle \mathbf{c}_t, \mathbf{w}\rangle\right)}{Z_{\mathbf{c}_t}}.$$

$$\Pr[w \text{ emitted at } t|\mathbf{c}_t] = \alpha p(w) + (1-\alpha)\frac{\exp\left(\langle \widetilde{\mathbf{c}}_t, \mathbf{w}\rangle\right)}{Z_{\widetilde{\mathbf{c}_t}}},$$
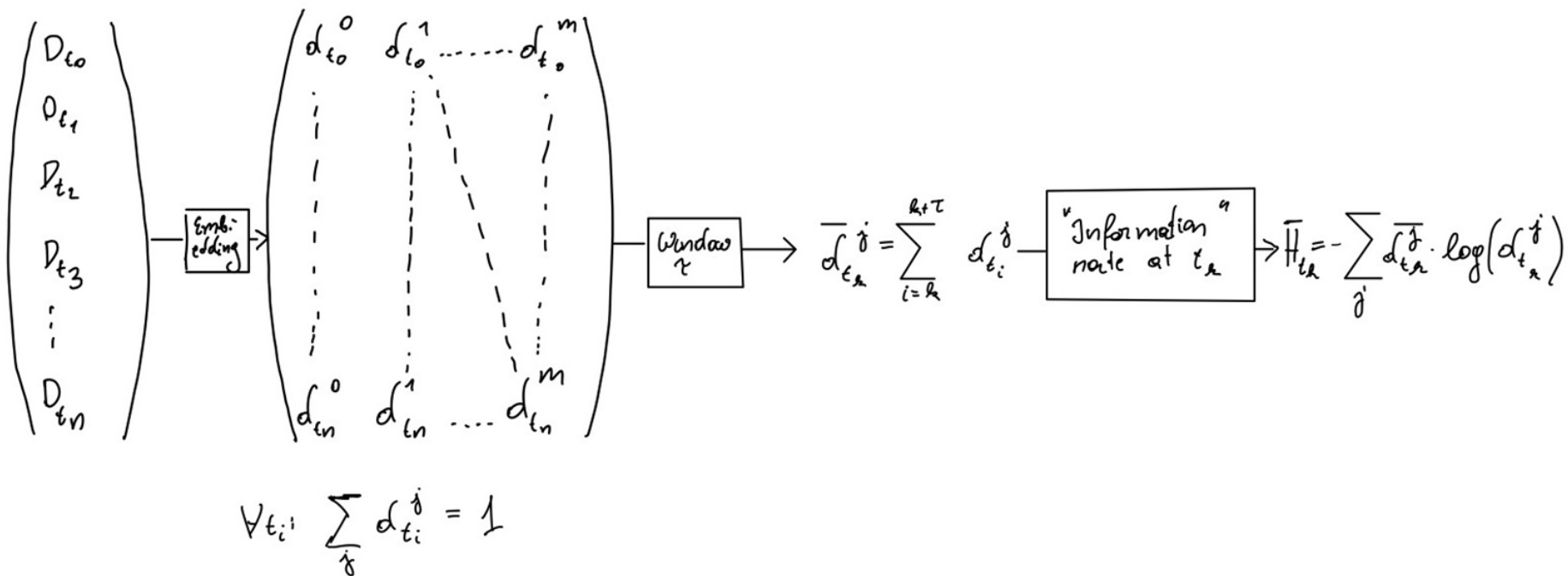
$$(\widetilde{\mathbf{c}}_t)_{\mathrm{MAP}} = \sum_{w\in\mathscr{C}} \frac{a}{p(w)+a}\mathbf{w}, \text{ where } a = \frac{1-\alpha}{\alpha Z}.$$

# Motivation

# Different languages, similar encoding efficiency: Comparable information rates across the human communicative niche

Christophe Coupé[1,2]*, Yoon Mi Oh[3,4]*, Dan Dediu[1,5], François Pellegrino[1†]

Language is universal, but it has few indisputably universal characteristics, with cross-linguistic variation being the norm. For example, languages differ greatly in the number of syllables they allow, resulting in large variation in the Shannon information per syllable. Nevertheless, all natural languages allow their speakers to efficiently encode and transmit information. We show here, using quantitative methods on a large cross-linguistic corpus of 17 languages, that the coupling between language-level (information per syllable) and speaker-level (speech rate) properties results in languages encoding similar information rates (~39 bits/s) despite wide differences in each property individually: Languages are more similar in information rates than in Shannon information or speech rate. These findings highlight the intimate feedback loops between languages' structural properties and their speakers' neurocognition and biology under communicative pressures. Thus, language is the product of a multiscale communicative niche construction process at the intersection of biology, environment, and culture.