



# UNIVERSITÀ DEGLI STUDI DI PADOVA

## **Network Science**

A.Y. 23/24

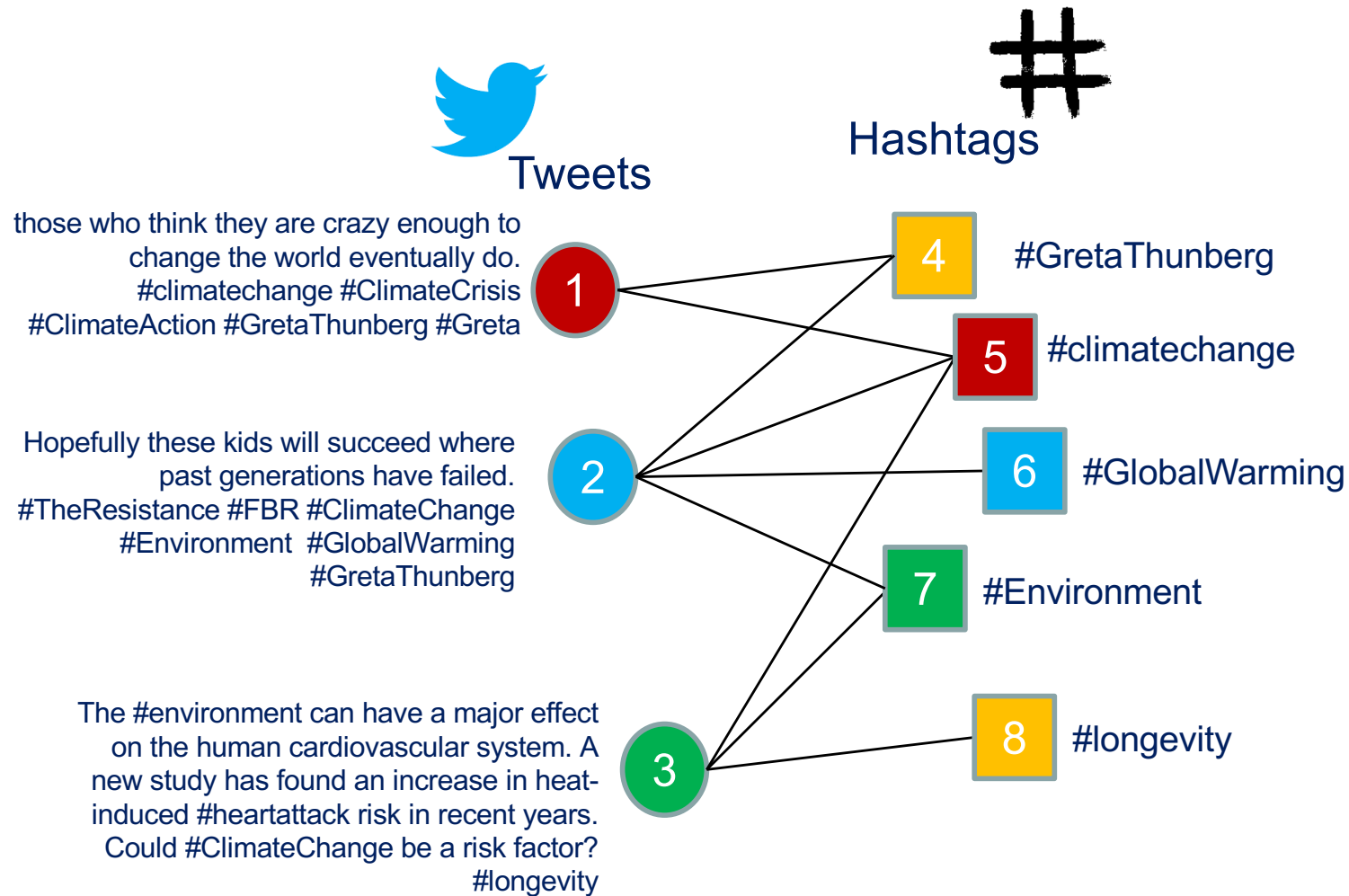
ICT for Internet & multimedia, Data science, Physics of data

# Semantic networks

network science tools for their study



# Conceptual picture of a semantic network on Twitter





- ❑ Data collection + polishing
- ❑ Building the semantic network (bipartite/projections)
- ❑ Topic (i.e., community) detection
  - ❑ **Modularity & InfoMap**
  - ❑ Non-negative matrix factorization (**NMF**)
  - ❑ Latent Dirichlet allocation (**LDA**)
  - ❑ Variational auto-encoders (**VAE**)
  - ❑ Embeddings and **BERTopic**

# Data collection

how to get data from the Internet using APIs



no longer available unless you pay 5k\$ per month ☹️

<https://developer.twitter.com/en/portal/dashboard>

## Twitter’s plan to cut off free data access evokes ‘fair amount of panic’ among scientists

Social media platform’s intent to increase revenue could end or limit many research projects

8 FEB 2023 · 4:35 PM ET · BY [KAI KUPFERSCHMIDT](#)

## Twitter’s plan to charge researchers for data access puts it in EU crosshairs

Elon Musk’s social media giant plans to charge academics to access its data – in potential violation of Europe’s content rules

BY MARK SCOTT  
MARCH 22, 2023



## Academic researchers blast Twitter’s data paywall as ‘outrageously expensive’



By [Brian Fung](#), CNN

Published 11:40 AM EDT, Wed April 5, 2023



## Reddit

## Subreddit

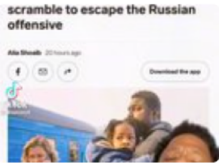
Search: ukrainer

Sort: Time

**r/blackoutukraine** · Posted by u/One-Designer-9406 2 years ago

**SPOILER**

Racism in Ukraine is getting out of hand  
#racisminukrain #ukrainewar #blackoutukraine




0 upvotes 0 comments

---

**r/MapPorn** · Posted by u/fpl123999 9 months ago

Ukraine war last 10 months in 15 seconds



2.3k upvotes 245 comments

Search: ukrainer

Communities


- r/ukrainewar** · 3.5k Members  
reddit for events of the Ukranian-Russian war ongoing since February 24, 2022. [Join](#)
- r/UkraineWarVideoReport** · 722k Members  
Community Driven Videos/Photos/Updates and Discussion [Join](#)
- r/RuZZiaUkrainewar** · 321 Members  
Unbiased reporting of the Russian Ukraine war. Please be respectful of Pro-Russian views as well as Pro-Ukrainian views. feel free to contribute [Join](#)
- r/UkrainianConflict** · 453k Members  
News, analysis, discussion and investigative journalism documenting the ongoing conflict in Ukraine. [Join](#)



GET NEW REDDIT | MY SUBREDDITS | HOME - POPULAR - ALL - RANDOM - USERS | ASKREDDIT - GAMING - PICS - TODAYILEARNED - FUNNY - WORLDNEWS - NEWS - MOVIES - M

reddit | PREFERENCES | options | **apps** | RSS feeds | friends | blocked | password/email | delete | Upbeat-Lychee-6630 (1) |

## developed applications



change icon

**ns2023**  
personal use script

**client\_id**  
Qbdk-FkA9jSQB9T7drY8UQ

download reddit context for the network science course, at the unioversity of Padova, of which I am the instructor

**secret** jtGPdqiaTj6hCWcvRPeS\_nMNEVkwXw **client\_secret**

**name**

**description**

**about url**

**redirect uri**

[delete app](#)

**username**

**developers** Upbeat-Lychee-6630 (that's you!) [remove](#)

add developer:

register  
asap... will  
be using this  
in the 1<sup>st</sup> lab





# Python Reddit API Wrapper – PRAW

<https://praw.readthedocs.io/en/stable/> (v7.7.1)

```
!pip install praw
```

```
import pandas as pd
import praw
reddit = praw.Reddit(client_id='Qbdk-FkA9jSQB9T7drY8UQ',
                    client_secret="jtGPdqiaTj6hCWcvRPeS_nMNEVkw",
                    user_agent='reddit scraper 1.0 by u/Upbeat-Lychee-6630',
                    check_for_async=False)
print(reddit.read_only)
```

True

```
df = pd.DataFrame([vars(post) for post in reddit.subreddit("all")
                  .search("#ukrainewar", sort='top', limit=10)])
```

```
df.to_excel('drive/MyDrive/Colab Notebooks/samples.xlsx',
           index=True)
```

your own description of your app,  
including version and username

from Reddit apps

"relevance", "hot", "top",  
"new", or "comments"

max 250  
per call

can also add `time_filter =`  
"all", "day", "hour", "month", "week", or "year"



there is a list of 116 entries per post,  
on which you can choose!!!

from this you extract the date

	title	created	score	upvote_ratio	ups	num_comments	selftext
0	Damn...we blinked and missed the T-34 stage of...	1.666899e+09	10394	0.99	10394	738	
1	Finnish🇫🇮 volunteer sends greetings home from ...	1.680237e+09	2095	1.00	2095	57	
2	Guess having 5 trucks fall into your office ca...	1.663341e+09	1974	1.00	1974	88	
3	[META] Important - Russia-Ukraine Crisis/War: ...	1.645712e+09	1284	0.89	1284	1	Hi, /u/Anonim97 here.\n\nWe - as a mods of 40k...
4	V*tniks coping hard Over the counter offensive...	1.662956e+09	1081	1.00	1081	86	
...	...	...	...	...	...	...	...
0	Russia Ukraine War.	1.695349e+09	1	1.00	1	0	





Content Posting API

Display API

Research API

About Research API

Getting Started

Frequently Asked Questions

Codebook

API Reference

Query Videos

Query User Info

Query Video Comments

## Query Videos

### Request

HTTP URL	https://open.tiktokapis.com/v2/research/video/query/
HTTP Method	POST
Scopes	research.data.basic

### Headers

Key	Type	Description
-----	------	-------------



### TikTok Research API application received

Thank you for applying to gain access to TikTok's research APIs. We have received your application and will proceed with the review process. You will be notified of the result via email within 3-4 weeks. If you have any questions, please contact our Support team.



ENSEMBLEDATA Home Products Documentation Pricing

# Social media scraping through simple APIs

Fetch data from TikTok, Instagram, YouTube through simple APIs.  
Real-time, fast, reliable and easy to integrate

Get started Documentation

1 week trial  
period...  
register  
later on

EnsembleData / TikTokScraper Public

Code Issues Pull requests Actions

main 1 branch 0 tags

- fracogno Fix
- images upload tt logo
- src Fix
- Update README.md

<https://github.com/EnsembleData/TikTokScraper>

# Data preprocessing

how to polish raw data from the Internet



## 1. Superficial cleaning

Removing website links  
Removing accented characters  
Removing text inside square brackets  
Removing moderator messages  
Removing double spaces  
Removing non-text special words and characters  
Removing extra-used new lines  
Limiting all the repetitions to two characters and removing the extra characters  
Removing punctuation except main sentence punctuation  
Removing sentences that represent the rules of the community

Fixing **contractions**  
Removing **emoji**  
Removing **hashtags** and **mentions**  
Removing **numbers**  
**Lowercasing**  
Correct **spellings**

## 2. Subsentence

Tokenise subsentences

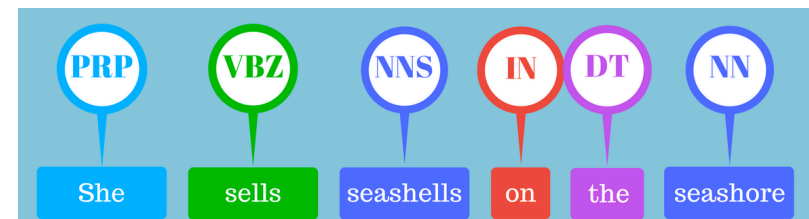
useful for long  
text samples  
(e.g., Reddit)

the bare minimum to  
polish the text,  
useful as an input to  
sentiment analysis

## 3. Deep cleaning

Stop word removal  
Word tokenization  
POS tagging  
Lemmatization

truly polished  
text, useful for  
building a  
semantic  
network





# SpaCy part-of-speech (POS) tags

<https://spacy.io/>

POS	description	example	POS	description	example
<b>ADJ</b>	<b>adjective</b>	big, old, green, incomprehensible, first	PART	particle	's, not,
ADP	adposition	in, to, during	<b>PRON</b>	<b>pronoun</b>	I, you, he, she, myself, themselves, somebody
<b>ADV</b>	<b>adverb</b>	very, tomorrow, down, where, there	<b>PROPN</b>	<b>proper noun</b>	Mary, John, London, NATO, HBO
AUX	auxiliary	is, has (done), will (do), should (do)	PUNCT	punctuation	., (, ), ?
CONJ	conjunction	and, or, but	SCONJ	subordinating conjunction	if, while, that
CCONJ	coordinating conjunction	and, or, but	SYM	symbol	\$, %, §, ©, +, -, ×, ÷, =, :, 😊
DET	determiner	a, an, the	<b>VERB</b>	<b>verb</b>	run, runs, running, eat, ate, eating
INTJ	interjection	psst, ouch, bravo, hello	X	other	sfpkdspxmsa
<b>NOUN</b>	<b>noun</b>	girl, cat, tree, air, beauty	SPACE	space	
NUM	numeral	1, 2017, one, seventy-seven, IV, MMXIV			

spaCy



title	true date	score	upvote_ratio	selftext	superficial cleaning	deep cleaning	deep cleaning
	created				title_sup_clean	title_deep_clean	title_deep_clean_pos
Damn...we blinked and missed the T-34 stage of...	2022-10-27	10390	0.99	NaN	damn we blinked and missed the t stage of the ...	damn blink miss t stage war	[damn ADV, blink VERB, miss VERB, t PROPN, sta...
Finnish🇫🇮 volunteer sends greetings home from ...	2023-03-31	2095	1.00	NaN	finnish volunteer sends greetings home from so...	finnish volunteer send greeting home	[finnish ADJ, volunteer NOUN, send VERB, greet...
Guess having 5 trucks fall into your office ca...	2022-09-16	1980	1.00	NaN	guess having trucks fall into your office can ...	guess have truck fall office significant emoti...	[guess VERB, have VERB, truck NOUN, fall VERB,...
[META] Important - Russia-Ukraine Crisis/War: ...	2022-02-24	1280	0.89	Hi, /u/Anonim97 here.\n\nWe - as a mods of 40k...	important russia ukraine crisis war info and ...	important russia ukraine crisis war info way help	[important ADJ, russia PROPN, ukraine PROPN, c...
V*tniks coping hard Over the counter offensive...	2022-09-12	1076	1.00	NaN	v tniks coping hard over the counter offensive...	tnik cope hard counter offensive traitor pfp lmao	[tnik NOUN, cope VERB, hard ADJ, counter NOUN,...
...	...	...	...	...	...	...	...

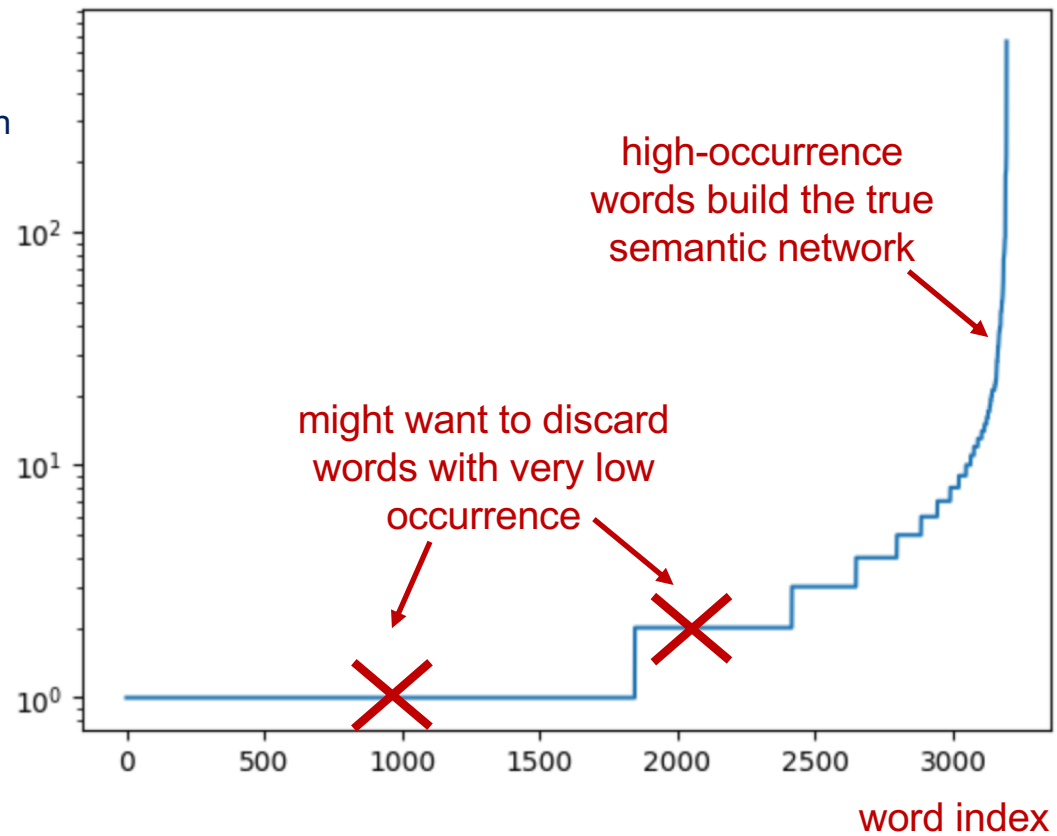


only ADJ, ADV, NOUN, PRON, PROPN, VERB kept





**occurrence**  
(i.e., number of  
times it appears in  
the documents)



# Building the semantic network

bipartite and projected counterparts



# Probability matrices linking words to documents

number of occurrences  
of words in documents

$$N_{wd} = \begin{array}{cccc} 0 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 1 \end{array} \begin{array}{l} \text{\#globalwarming} \\ \text{\#climatechange} \\ \text{\#climateaction} \\ \text{\#gretathunberg} \\ \text{\#environment} \end{array}$$

probability of words  
given a documents

$$P_{w|d} = \begin{array}{cccc} 0 & \frac{1}{2} & \frac{1}{5} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{2} & \frac{1}{5} & \frac{1}{4} \\ \frac{1}{3} & 0 & \frac{1}{5} & 0 \\ 0 & 0 & \frac{1}{5} & \frac{1}{4} \\ \frac{1}{3} & 0 & \frac{1}{5} & \frac{1}{4} \end{array}$$

we identify a **document** probability  $p_d = \begin{cases} \frac{1}{D} & \text{equally likely} \\ \frac{n_d}{\sum_d n_d} & \text{custom} \end{cases}$

we capture the **statistical** properties by normalizing by columns



# Probability matrices projecting to words or documents

bipartite  
network

joint probability of words  
and documents

$$P_{wd} = P_{w|d} \text{diag}(p_d)$$

0	1/8	1/20	1/16
1/12	1/8	1/20	1/16
1/12	0	1/20	0
0	0	1/20	1/16
1/12	0	1/20	1/16

marginal probabilities

$$p_w = P_{wd} \mathbf{1} \quad p_d = P_{wd}^T \mathbf{1}$$

$$p_{w_1, w_2} = \sum_d p_{w_1|d} \cancel{p_{w_2|d}} p_{w_2, d}$$

$$P_{ww} = P_{wd} \text{diag}(p_d)^{-1} P_{wd}^T$$

$$p_w = P_{ww} \mathbf{1}$$

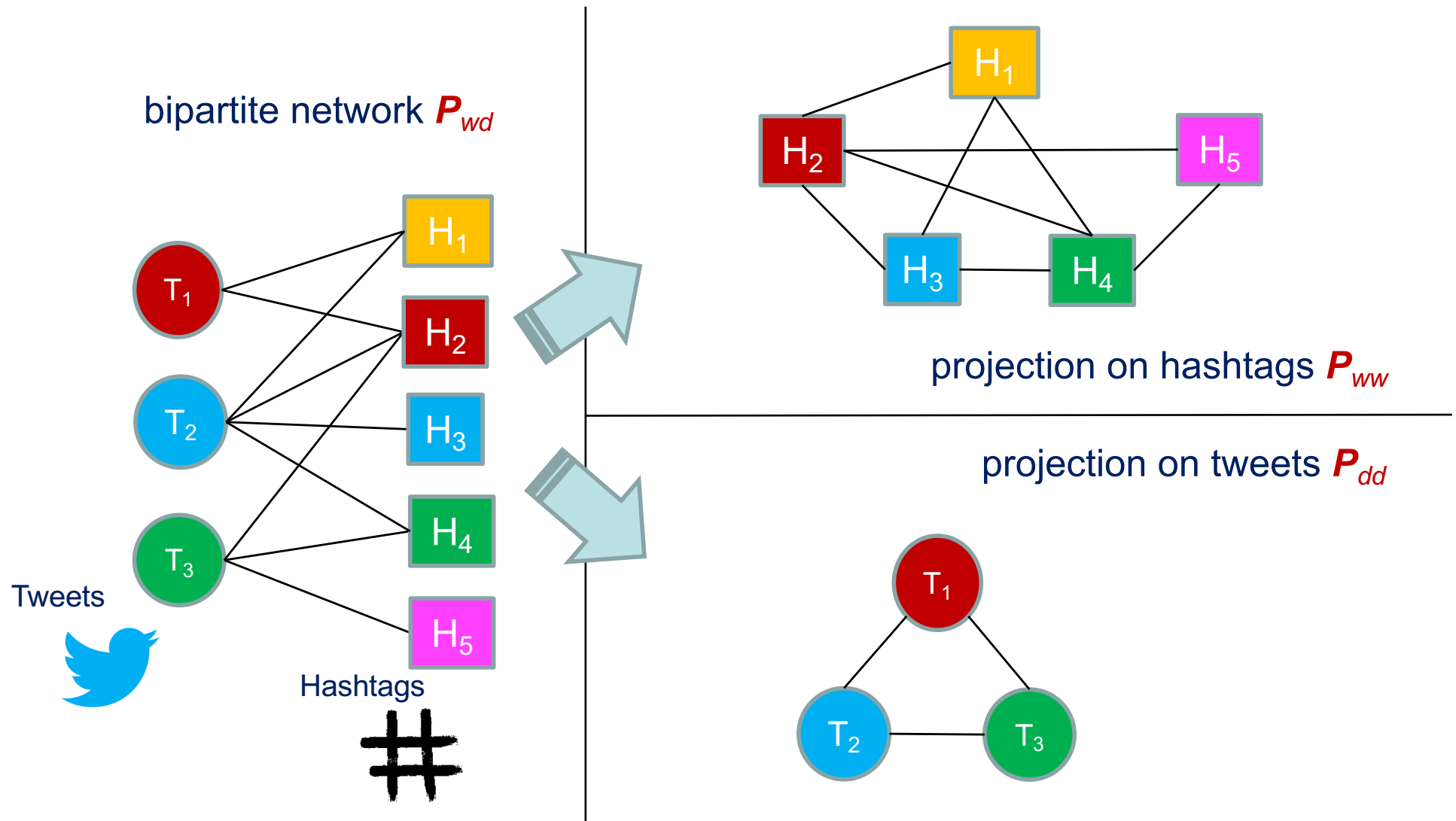
projection on words

projection on documents

$$p_{d_1, d_2} = \sum_w p_{d_1|w} \cancel{p_{d_2|w}} p_{d_2, w}$$

$$P_{dd} = P_{wd}^T \text{diag}(p_w)^{-1} P_{wd}$$

$$p_d = P_{dd} \mathbf{1}$$





# The role of TF-IDF

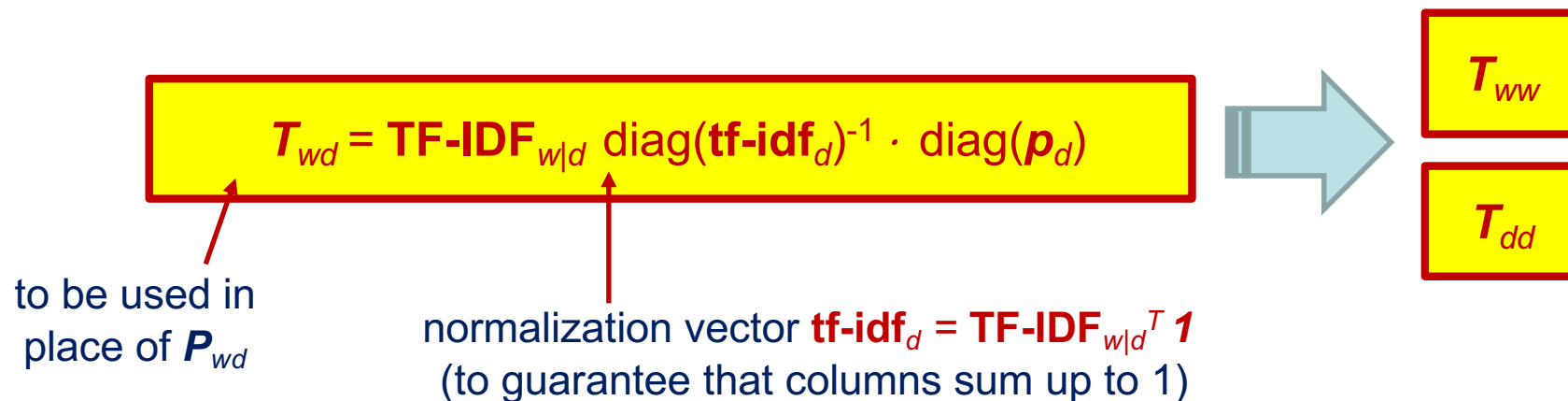
term frequency – inverse document frequency

term frequency =  
frequency (probability) of  
the word in the document

inverse document frequency =  
(log) fraction of documents  
that contain the word

$$\text{TF-IDF}_{w|d} = p_{w|d} \cdot -\log \left( \frac{\sum_d (n_{wd} > 0)}{D} \right)$$

- ❑ An heuristic
- ❑ **Punishes** words that appear in many documents
- ❑ **Enhances** words that are document specific



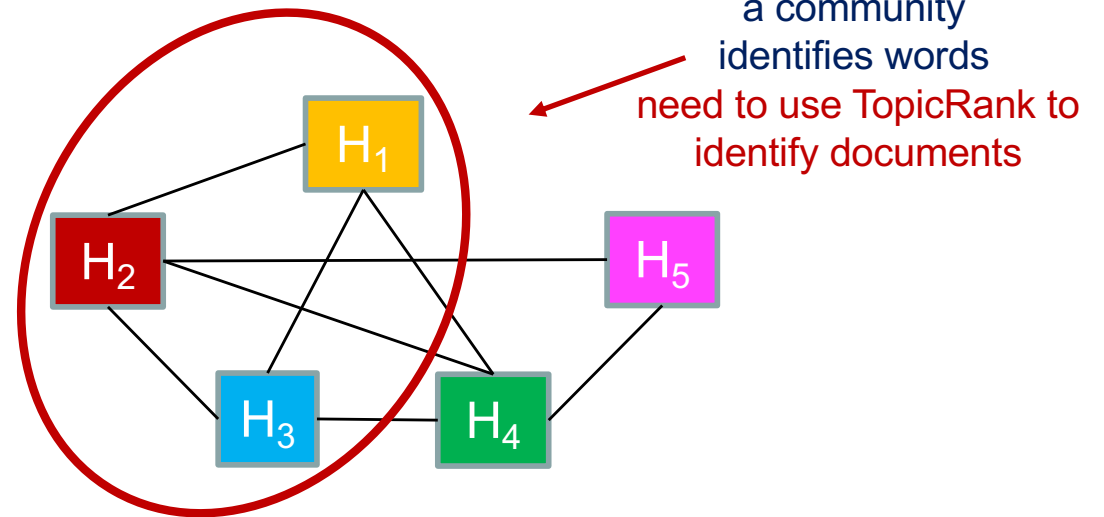
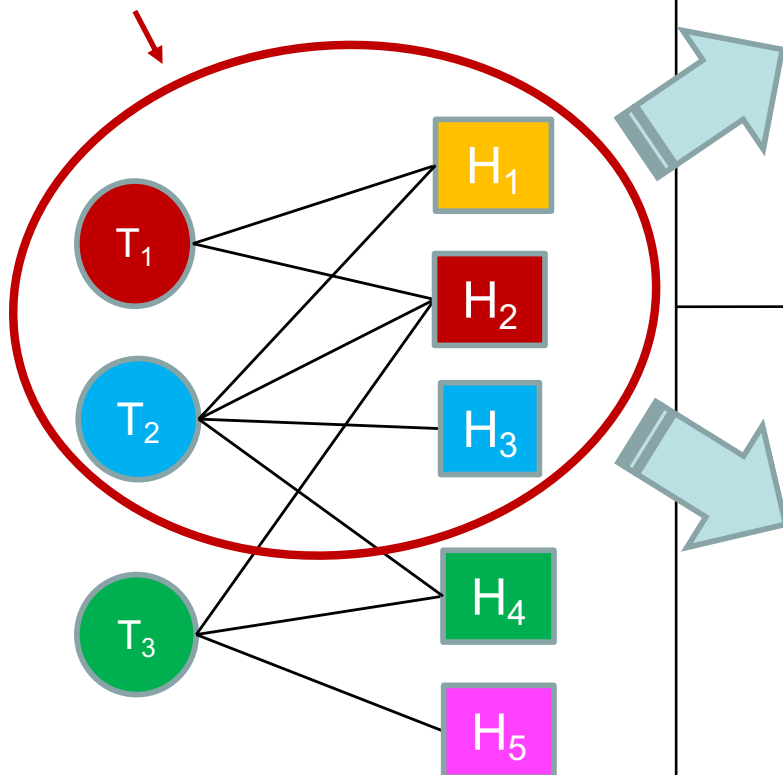
# Topic detection

i.e., community detection in semantic networks



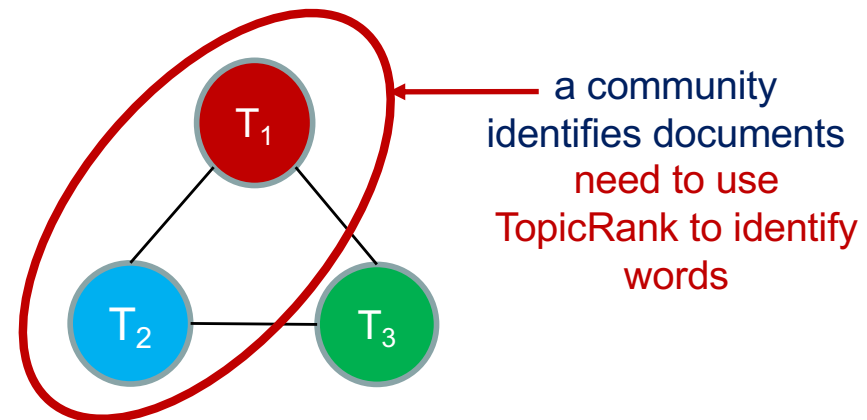
bipartite network  $P_{wd}$  or  $T_{wd}$

a community identifies both documents and words



projection on words  $P_{ww}$  or  $T_{ww}$

projection on documents  $P_{dd}$  or  $T_{dd}$

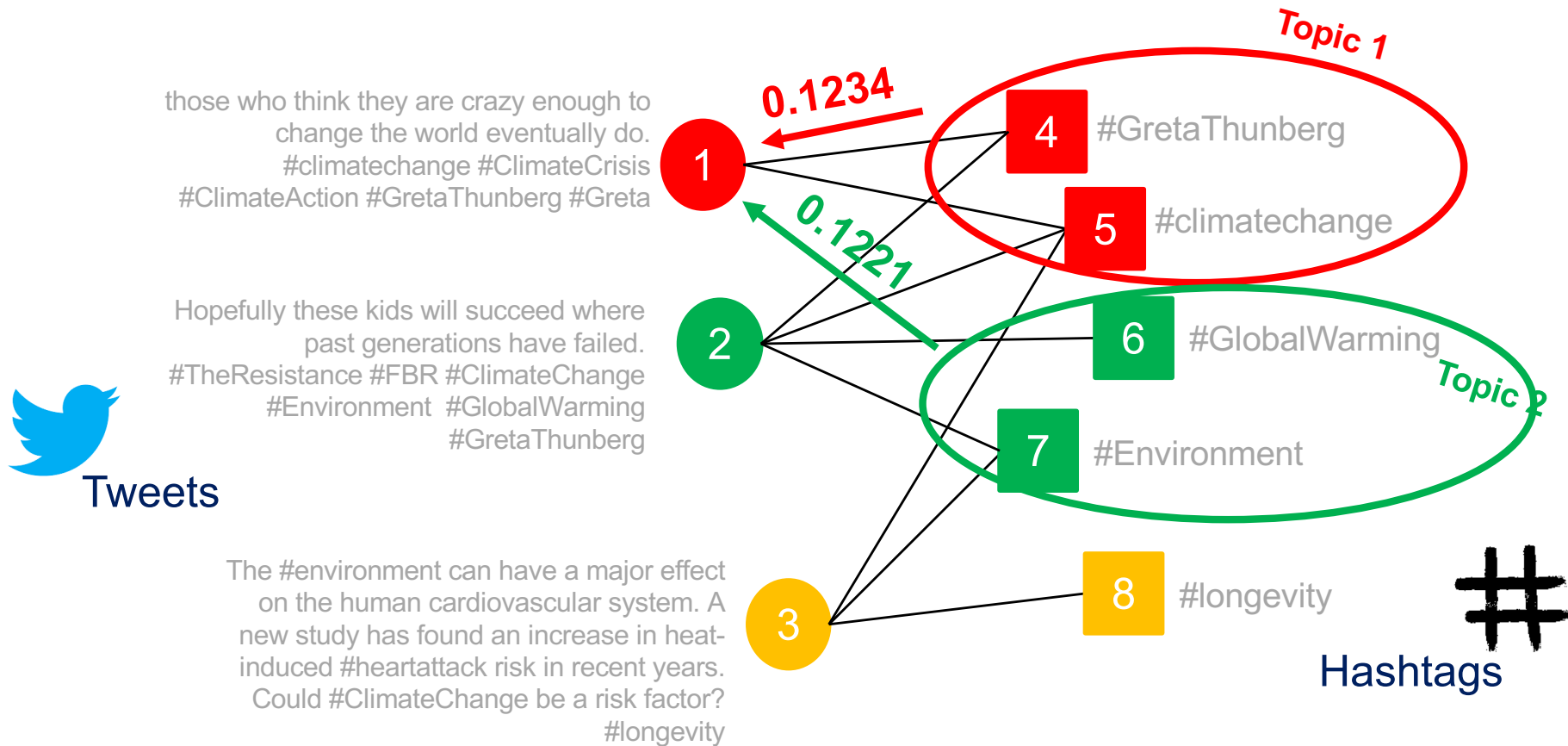








**Tweet 1** is assigned  
to **Topic 1** !!!





# Normalized mutual information

a wrap-up in topic detection

statistical dependencies about  
words and topics

$$P_{wt} = P_{wd} C^T$$



probability of a  
topic

$$p_t = P_{wt}^T \mathbf{1}$$

fraction of knowledge related to  
the topic that is explained by  
words (equal to 1 if topics use  
different words)

$$\text{NMI} = \frac{I(W;T)}{H(T)}$$





**C** topic assignment to be assessed for quality

$$P_{tt} = C P_{dd} C^T$$

can be interpreted as a probability matrix linking topics, its entries are the sum of the links of **A** from topic *i* to topic *j*

$P_{11}$	$P_{12}$	$P_{13}$
$P_{21}$	$P_{22}$	$P_{23}$
$P_{31}$	$P_{32}$	$P_{33}$

$$p_t = P_{tt} \mathbf{1}$$

can be interpreted as the probability vector of topics

modularity

$$Q = \sum_t (P_{tt} - p_t^2) < 1$$

to be maximized

normalized cut

normalized version

$$Ncut = 1 - \frac{\sum_t P_{tt}/p_t}{\sum_t 1} > 0$$

to be minimized



PageRank vector (ranking of documents)

$$\mathbf{r} = (1-c) \mathbf{P}_{d|d} \mathbf{r} + c \mathbf{1}/N$$

Here  $c_i$   
is the  $i$ th  
row of  $\mathbf{C}$

$$\mathbf{P}_{d|d} = \mathbf{P}_{dd} \text{diag}^{-1}(\mathbf{p}_d)$$

$$q_i = \left( 1 - (1-c) \frac{c_i \mathbf{1}}{N} \right) \mathbf{z}_i \mathbf{1} - c c_i \mathbf{P}_{d|d} \mathbf{z}_i^T$$

$$\mathbf{z}_i = c_i \text{diag}(\mathbf{r})$$

$$\text{InfoMap} = f(\mathbf{q}) + \sum_i f([q_i, \mathbf{z}_i])$$

normalized version

$$\frac{\text{InfoMap}}{f(\mathbf{r})} - 1$$

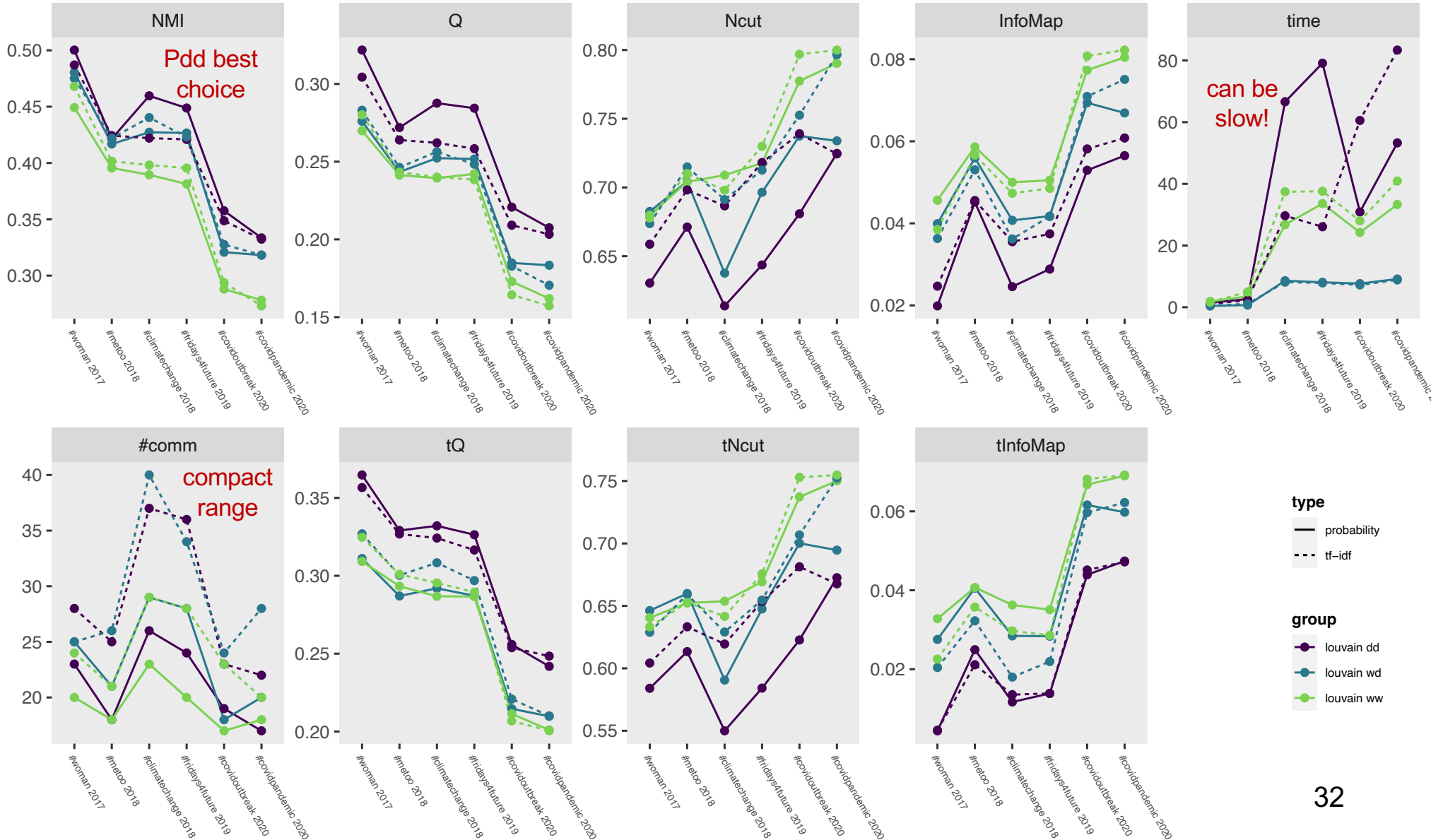
to be minimized

entropy function

$$f(\mathbf{x}) = - \sum_i x_i \log \left( \frac{x_i}{\sum_j x_j} \right)$$

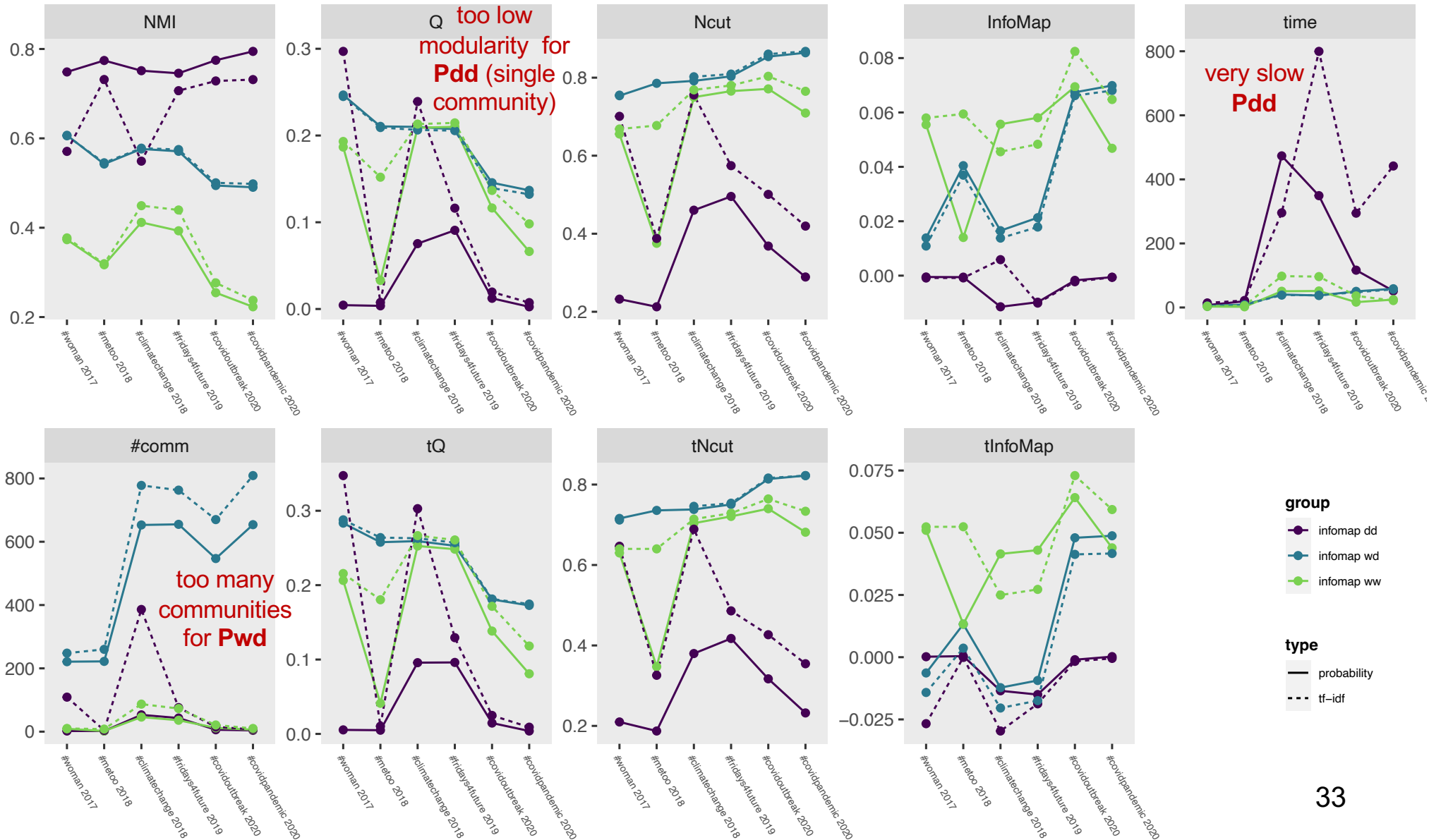


# A comparison of the different approaches - Louvain



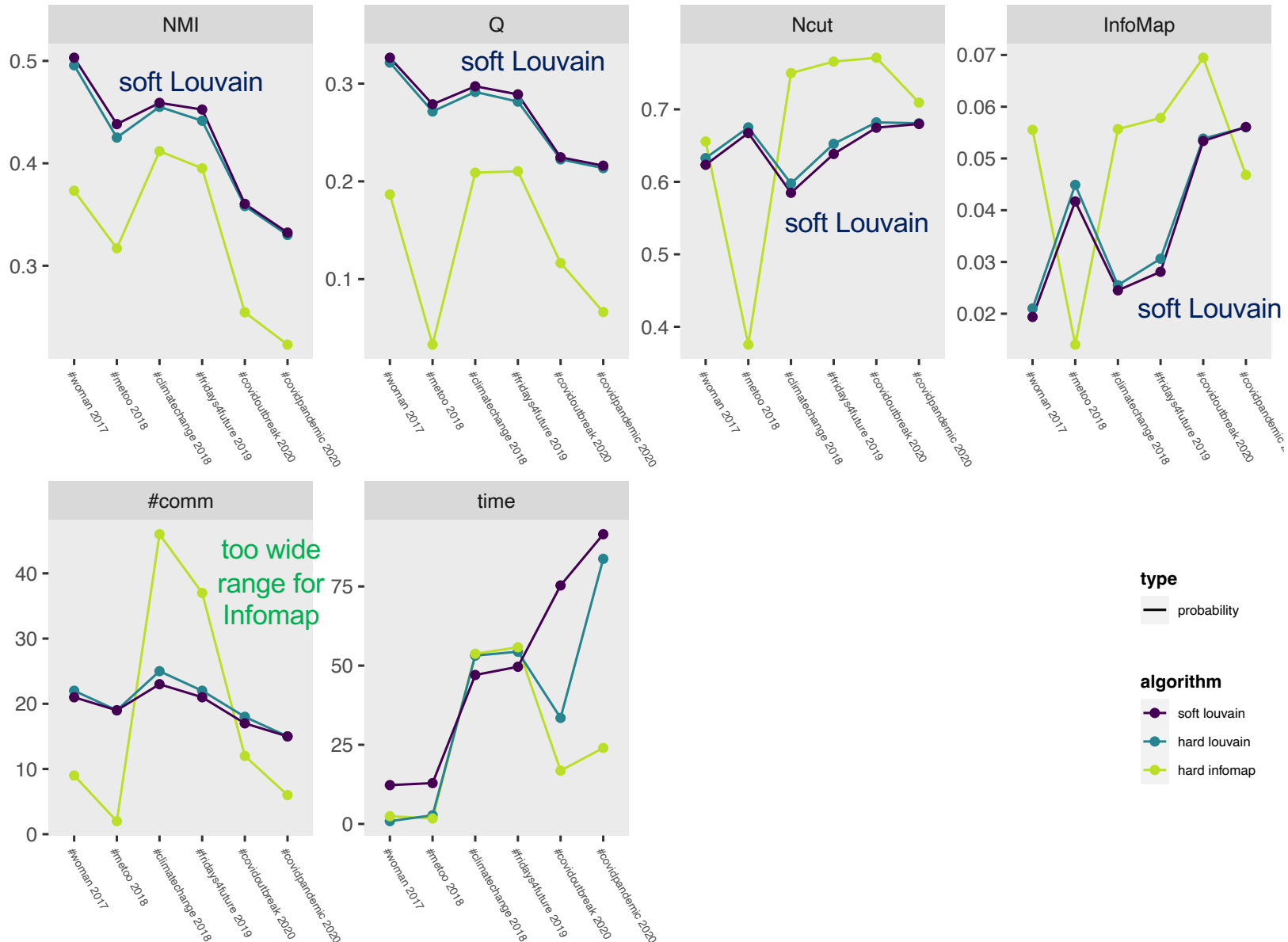


# A comparison of the different approaches - Infomap

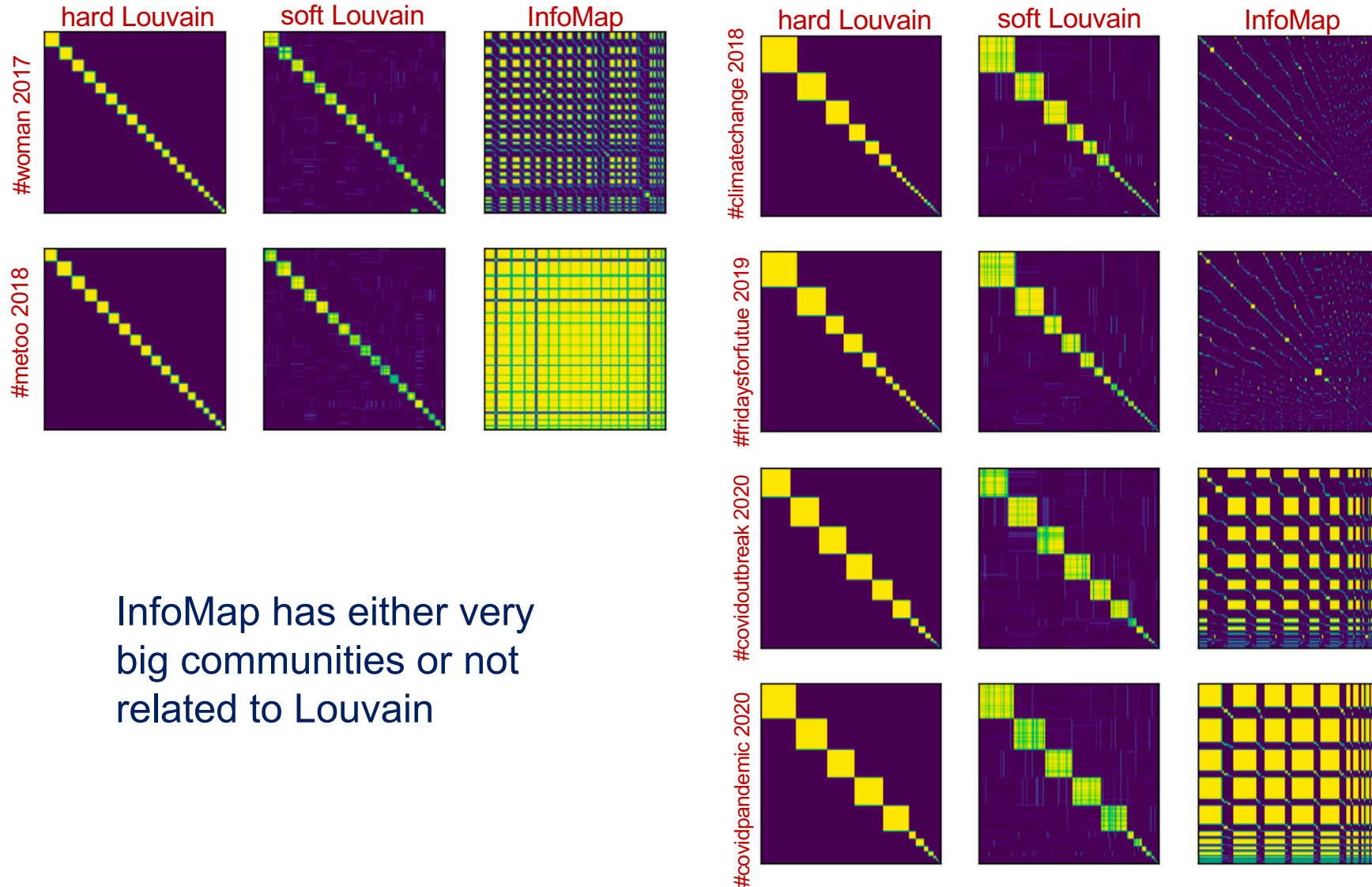




# A comparison hard/soft Louvain Pdd versus InfoMap Pww







InfoMap has either very big communities or not related to Louvain



- ❑ Louvain Pdd – provides the best results  
produces balanced clusters
- ❑ Louvain soft – slightly strengthens the result
- ❑ Bipartite networks – run much faster  
but performance deteriorates
- ❑ InfoMap – not robust 😞  
would be nice to see **BigCLAM** and **SBMs**  
... your task! 😊

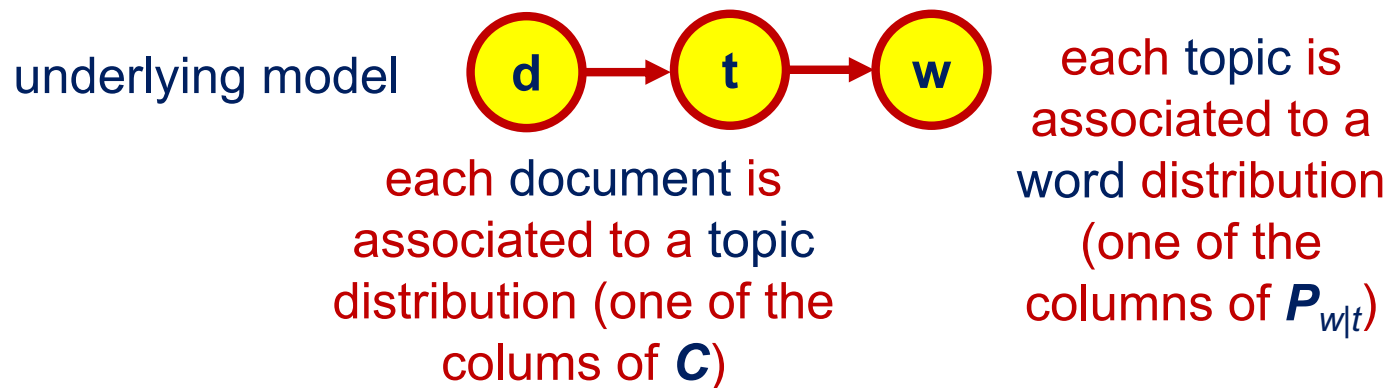
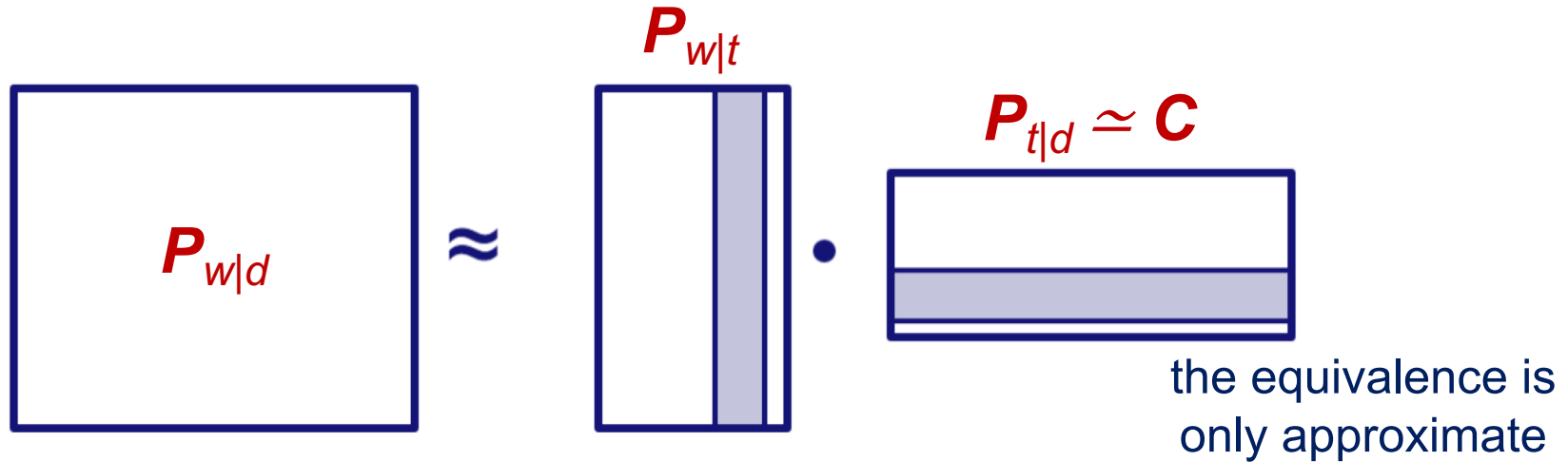
# Non-negative Matrix Factorization

and its application to topic detection



# NMF = nonnegative matrix factorization

rationale





$A = P_{w|d}$  is column stochastic

$$\operatorname{argmin}_{W \geq 0, H \geq 0} \sum_{ij} |A_{ij} - [WH]_{ij}|^2$$

minimizing the Frobenius norm does not ensure a column stochastic product  $WH$

$$\operatorname{argmin}_{W \geq 0, H \geq 0} \sum_{ij} A_{ij} \log \left( \frac{A_{ij}}{[WH]_{ij}} \right) - A_{ij} + [WH]_{ij}$$

minimizing the generalized Kullback-Leibler divergence ensures a column stochastic product  $WH$

$$f(y) = x \log \left( \frac{x}{y} \right) - x + y$$

$$f'(y) = -\frac{x}{y} + 1 = 0 \rightarrow y = x$$

Ho & Van Dooren. "Non-negative matrix factorization with fixed row and column sums." (2008)



```
from sklearn.decomposition import NMF
Pwgd = Pwd/Pwd.sum(axis=0).flatten()
```

run on different number of topics, then choose  
the best fit, e.g., according to modularity

```
# fit nmf model  $X = W \cdot H$ 
model = NMF(n_components=i, init='nndsvd',
            solver='mu', beta_loss='kullback-leibler')
W = model.fit_transform(Pwgd)
H = sps.csr_matrix(model.components_)
# column normalized versions
H = sps.diags(W.sum(axis=0).flatten())*H # Ptgd
W = W/W.sum(axis=0).flatten() # Pwgt
# community assignment C
C = sps.csr_matrix(np.transpose(H/H.sum(axis=0).flatten()))
```

wisely initialize  
for best  
performance

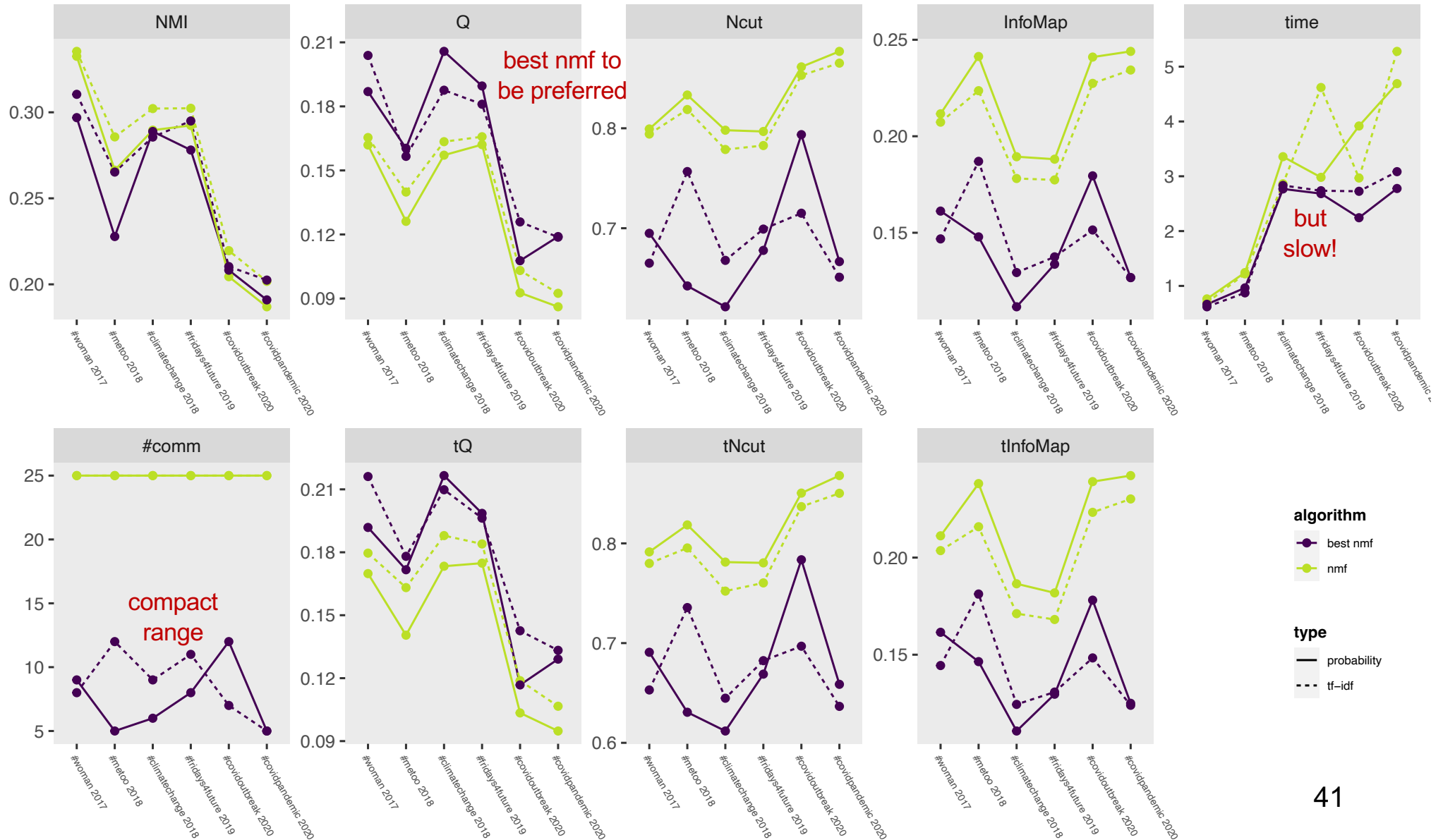
choose generalized  
Kullback-Leibler  
divergence, and the  
related solver

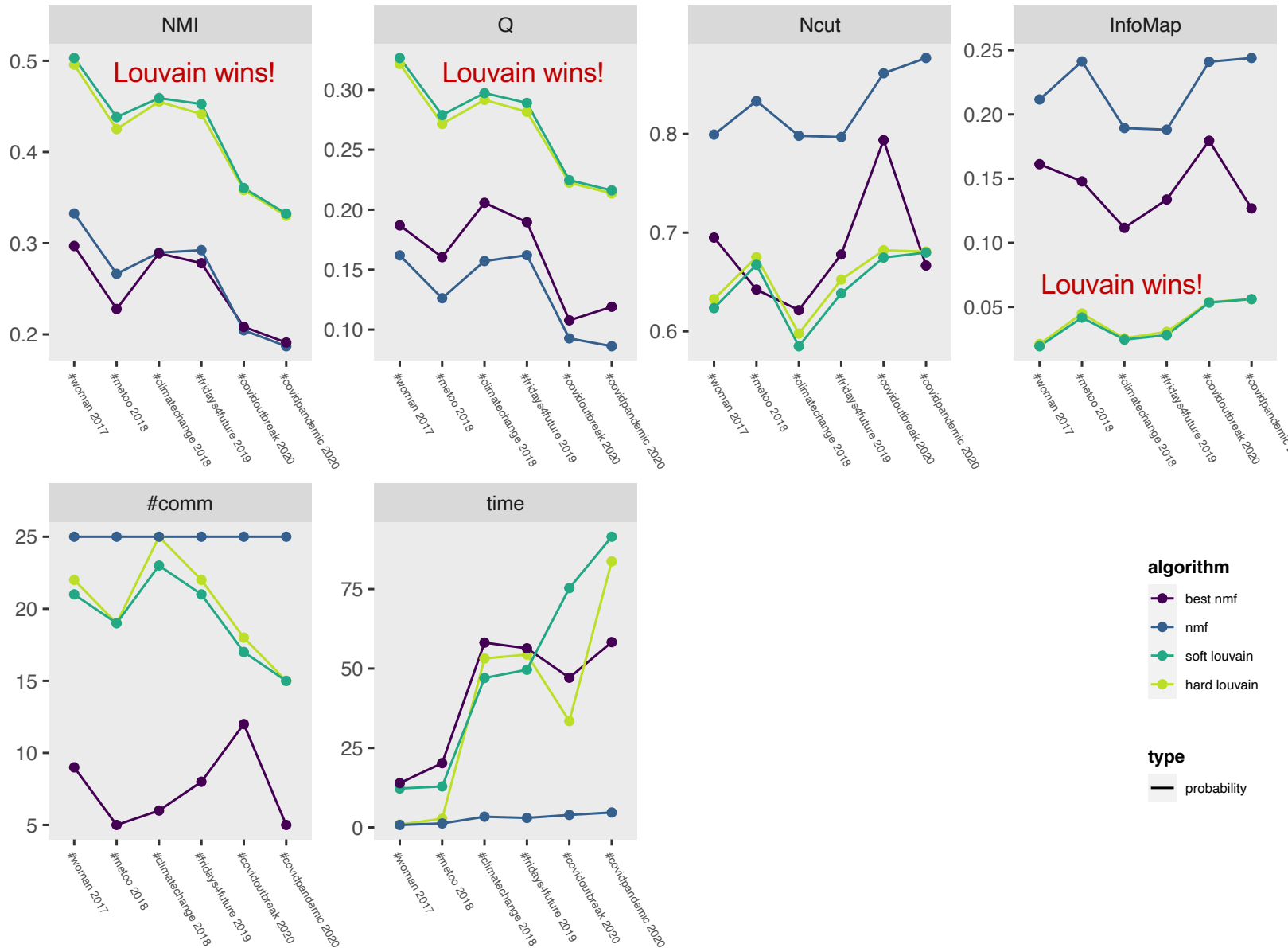
need to make  $W$   
column stochastic,  
to have  $H$  column  
stochastic too

force column stochasticity in  $H$  (not needed though)



# A comparison of the different approaches





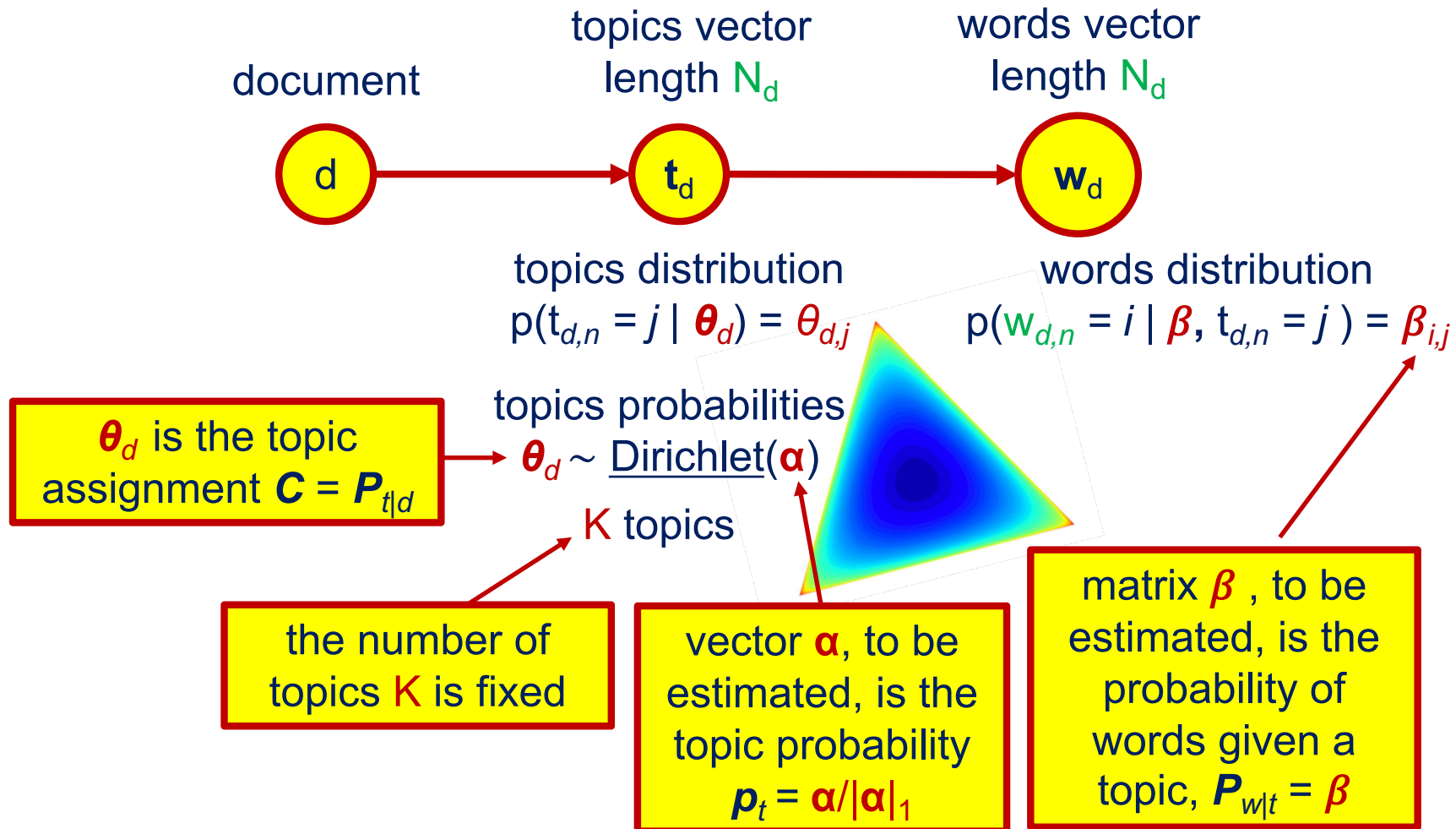




- ❑ Naturally provides a soft topic assignment
- ❑ NMF – not strikingly good
  - probably due to the fact that we want to express a sparse matrix through an eigenvector-like product with few eigenvectors (the fit is far from ideal)
- ❑ Comparison – with Louvain
  - much weaker
- ❑ Complexity – generally slow
  - need to test it for different numbers of topics ☹
  - fast for fixed topic number

# Latent Dirichlet allocation

LDA = a stochastic model for topic detection





topics assignment probability (Dirichlet)

$$p(\boldsymbol{\theta}_d | \boldsymbol{\alpha}) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K [\theta_{d,k}]^{\alpha_k - 1}$$

words probability

$$p(\mathbf{w}_d | \boldsymbol{\beta}, \boldsymbol{\theta}_d) = \prod_{n=1}^{N_d} [\boldsymbol{\beta} \boldsymbol{\theta}_d]_{w_{d,n}}$$

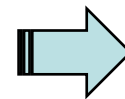
this dependence  
between  $\boldsymbol{\beta}$  and  $\boldsymbol{\theta}$   
is the trickiest part

overall probability

$$p(\text{corpus} | \boldsymbol{\alpha}, \boldsymbol{\beta}) = \prod_d \int p(\mathbf{w}_d | \boldsymbol{\beta}, \boldsymbol{\theta}_d) p(\boldsymbol{\theta}_d | \boldsymbol{\alpha}) d\boldsymbol{\theta}_d$$

target optimization

$$\operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}} p(\text{corpus} | \boldsymbol{\alpha}, \boldsymbol{\beta})$$



$$\mathbf{C} = \mathbf{P}_{t|d} = \boldsymbol{\theta}$$
$$\mathbf{P}_{wt} = \boldsymbol{\beta} \operatorname{diag}(\boldsymbol{\alpha} / |\boldsymbol{\alpha}|_1)$$

this is what we get



```
from sklearn.decomposition import LatentDirichletAllocation
```

```
# fit lda model
```

```
lda = LatentDirichletAllocation(n_components=i,  
                                learning_method="batch")
```

```
lda.fit(Mwd.T)
```

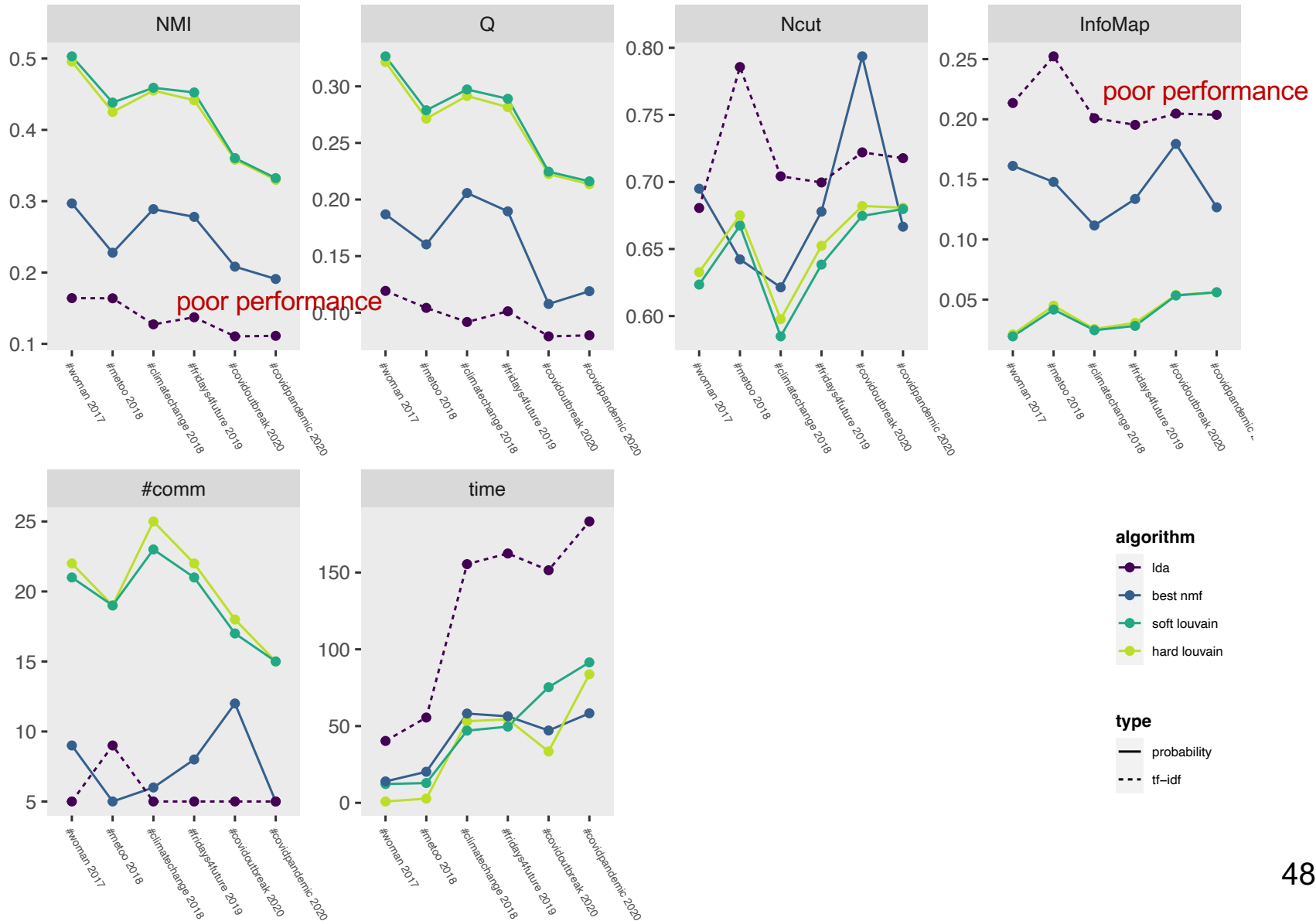
```
# community assignment C = Ptgd'
```

```
C = sps.csr_matrix(lda.transform(Mwd.T))
```

← initialise and fit model

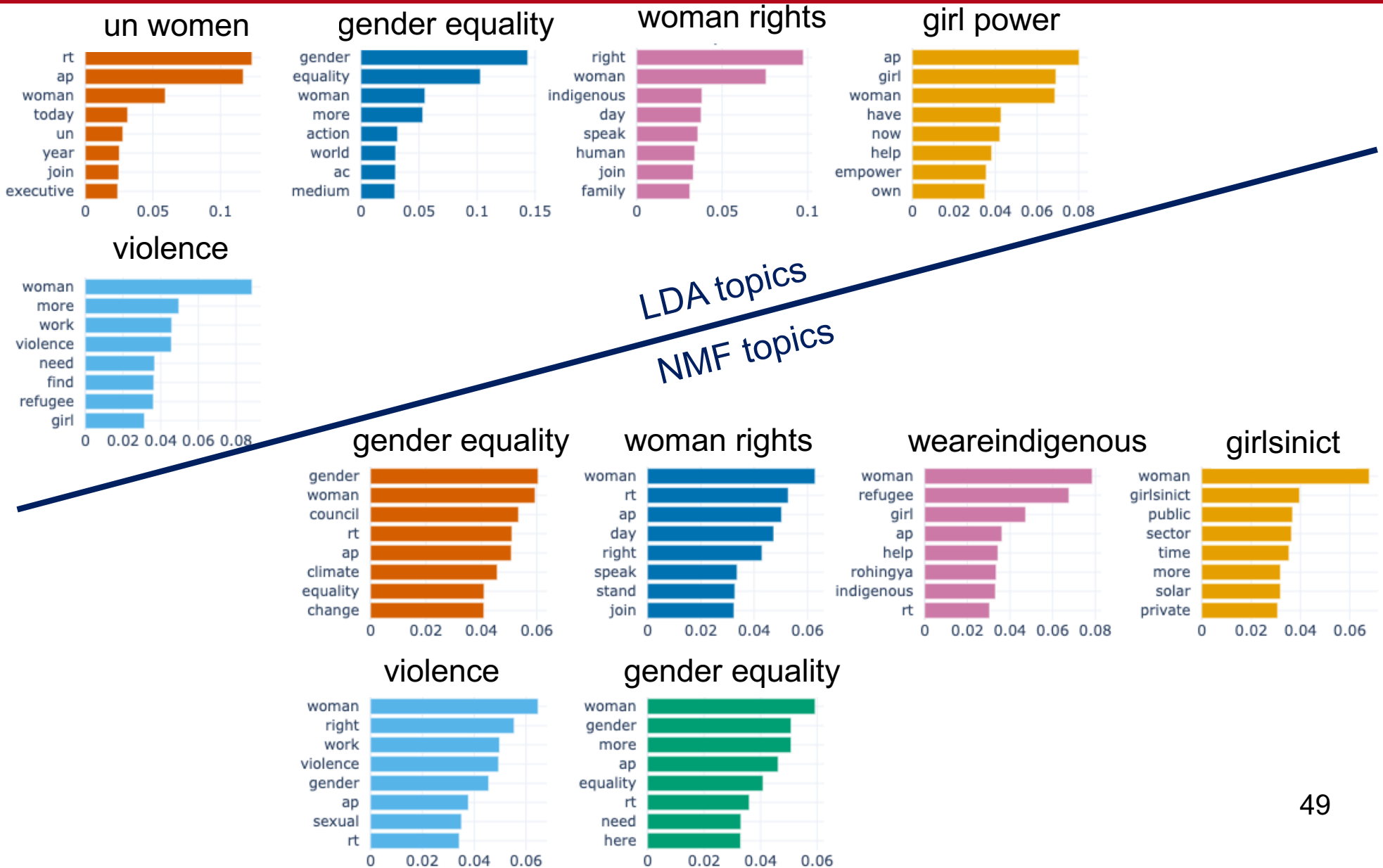
← extract topic assignment

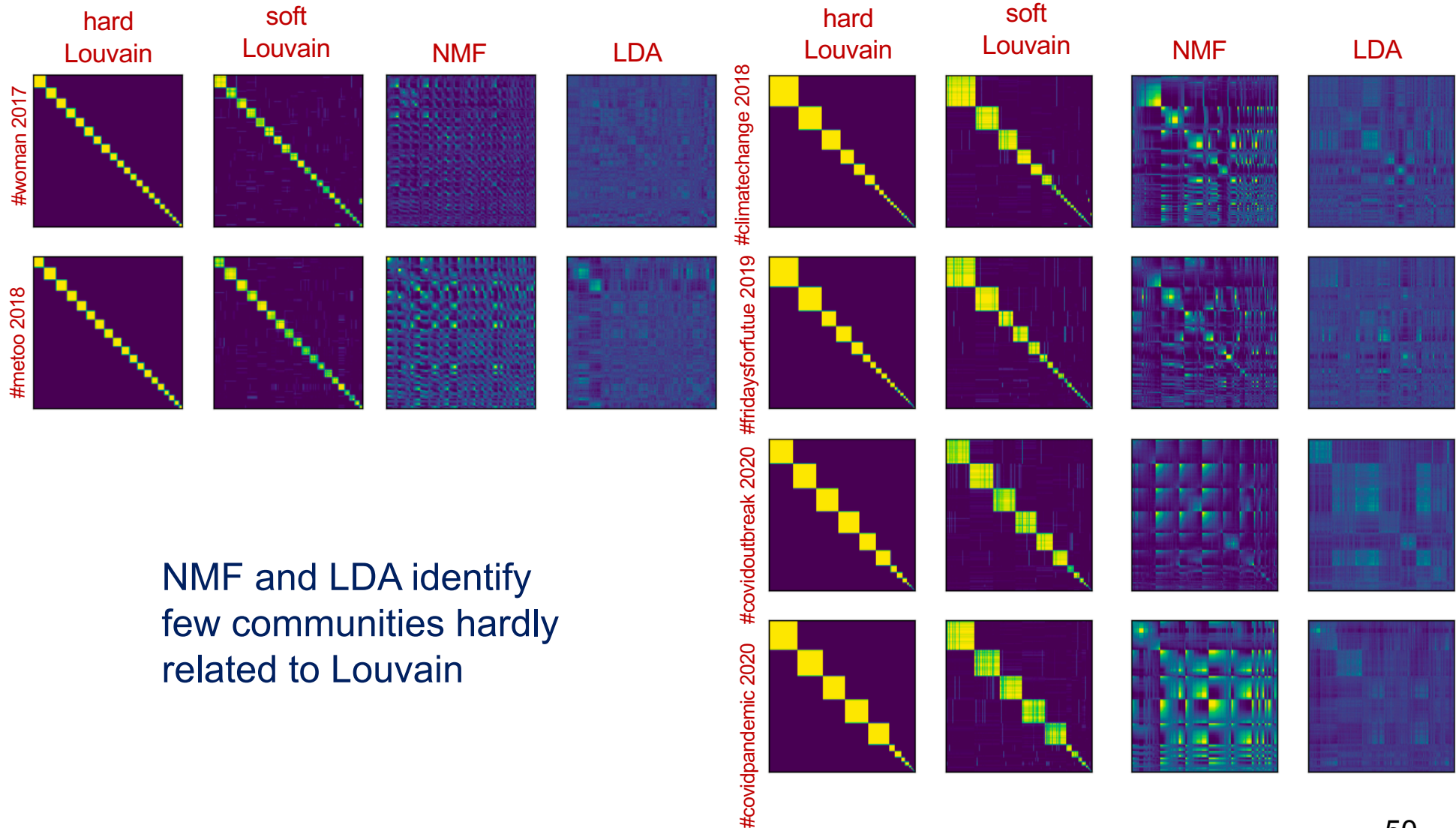
LDA





# A comparison NMF versus LDA topics





NMF and LDA identify  
few communities hardly  
related to Louvain





- ❑ Naturally provides a soft topic assignment
- ❑ LDA – not strikingly good
  - same eigenvector-like product as NMF
  - worse than NMF ... known issue ☹
  - probably due to the **Dirichlet** assumption (questionable)
  - and the **variational inference** (suboptimum approach)
- ❑ Comparison – with Louvain
  - much weaker
- ❑ Complexity – generally slow
  - need to test it for different numbers of topics ☹
  - fast for fixed topic number

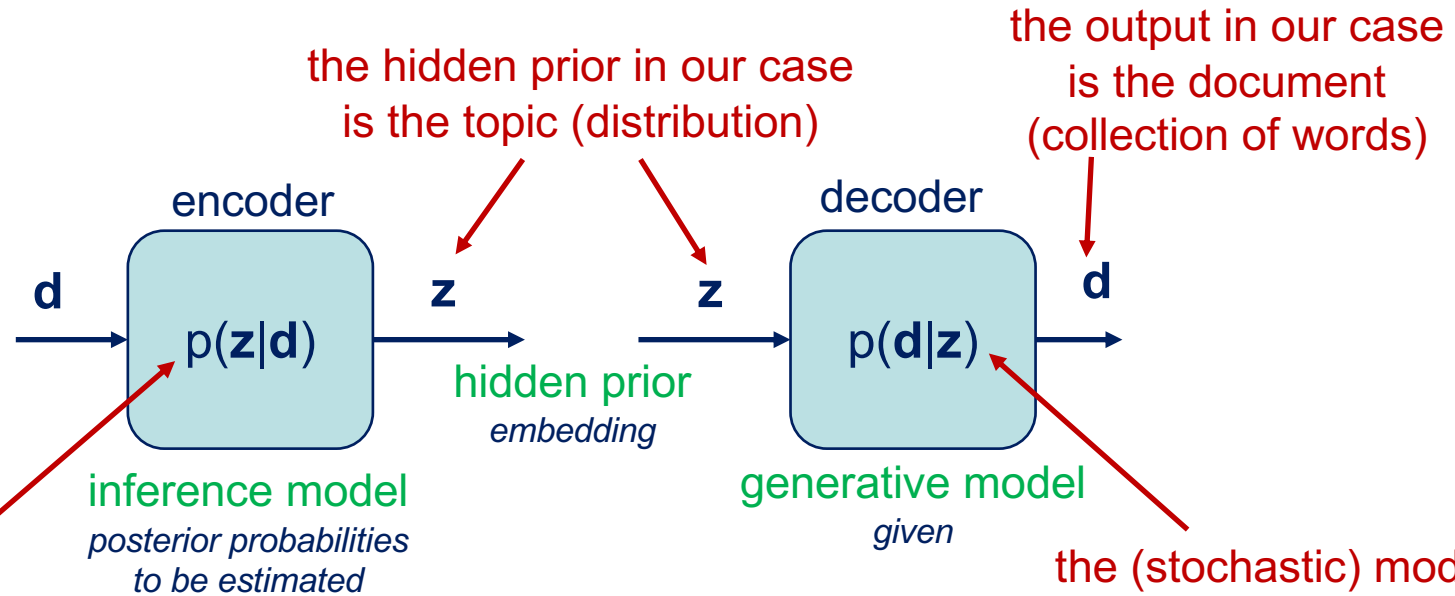
# Variational Auto Encoders

an application to topic analysis



# Variational Auto-Encoders

Kingma, Welling, "Auto-encoding variational Bayes," (2013)



but we are interested in the inverse link that, given a document tells what topic it is associated with

the (stochastic) model explains how a document is generated from a topic (distribution)

$$p(z|d) = \frac{p(d|z) p(z)}{p(d)} \cong q(z|d)$$

impossible to know in the closed form

needs an a-priori model for the embedding

is approximated by a simple alternative model

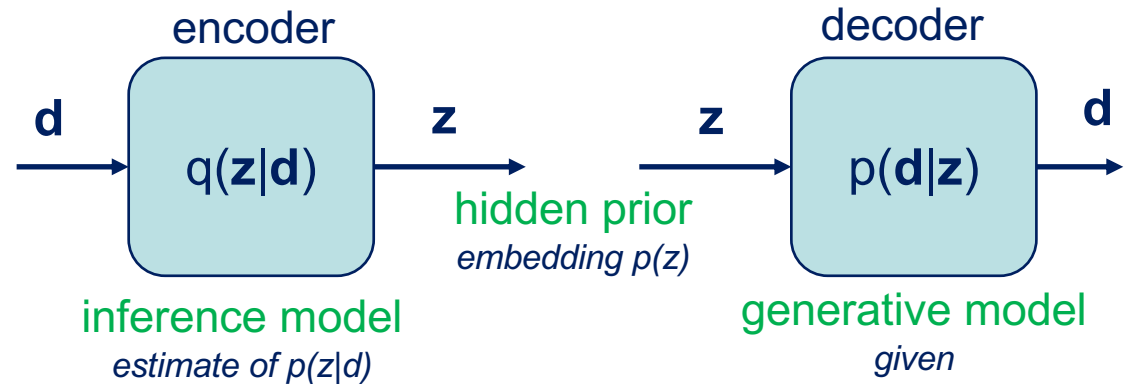


# VAE optimization rationale

ELBO = evidence lower bound

## ELBO

$$\mathcal{L}_{\theta, \phi}(\mathbf{d}) \leq \log p_{\theta}(\mathbf{d})$$



$$\begin{aligned} \mathcal{L}_{\theta, \phi}(\mathbf{d}) &= \log p_{\theta}(\mathbf{d}) - D_{\text{KL}}\left(q_{\phi}(z|\mathbf{d}) \parallel p_{\theta}(z|\mathbf{d})\right) \\ &= \int dz q_{\phi}(z|\mathbf{d}) \log \left( \frac{p_{\theta}(z, \mathbf{d})}{q_{\phi}(z|\mathbf{d})} \right) \\ &= \underbrace{\int dz q_{\phi}(z|\mathbf{d}) \log \left( p_{\theta}(\mathbf{d}|z) \right)}_{\mathcal{L}_1} + \underbrace{\int dz q_{\phi}(z|\mathbf{d}) \log \left( \frac{p_{\theta}(z)}{q_{\phi}(z|\mathbf{d})} \right)}_{\mathcal{L}_2} \end{aligned}$$

to be maximized wrt parameters  $\theta$  and  $\phi$  provides fitting on  $p(\mathbf{z})$ ,  $p(\mathbf{d}|\mathbf{z})$ , and  $q(\mathbf{z}|\mathbf{d})$

a-priori model (given)

inference model (approximate)      generative model (given)



$$\underbrace{\int dz q_{\phi}(z|\mathbf{d}) \log \left( \frac{p_{\theta}(z)}{q_{\phi}(z|\mathbf{d})} \right)}_{\mathcal{L}_2}$$

both should have a simple parametrization on  $\theta$  and  $\phi$

e.g., the Gaussian case

$$p_{\theta}(\mathbf{z}) = \frac{1}{\sqrt{\det(2\pi \text{diag}(\boldsymbol{\sigma}_{\theta}^2))}} \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu}_{\theta})^T \text{diag}^{-1}(\boldsymbol{\sigma}_{\theta}^2)(\mathbf{z} - \boldsymbol{\mu}_{\theta})\right)$$

$$q_{\phi}(\mathbf{z}|\mathbf{d}) = \frac{1}{\sqrt{\det(2\pi \text{diag}(\boldsymbol{\sigma}_{\phi}^2(\mathbf{d})))}} \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu}_{\phi}(\mathbf{d}))^T \text{diag}^{-1}(\boldsymbol{\sigma}_{\phi}^2(\mathbf{d}))(\mathbf{z} - \boldsymbol{\mu}_{\phi}(\mathbf{d}))\right)$$

$$\mathcal{L}_2(\theta, \phi) = \frac{1}{2} \sum_i \left[ 1 + \log\left(\frac{\sigma_{\phi,i}^2(\mathbf{d})}{\sigma_{\theta,i}^2}\right) - \frac{\sigma_{\phi,i}^2(\mathbf{d})}{\sigma_{\theta,i}^2} - \frac{(\mu_{\phi,i}(\mathbf{d}) - \mu_{\theta,i})^2}{\sigma_{\theta,i}^2} \right]$$



# L1 ELBO function

approximated through Monte Carlo estimation

$$\underbrace{\int dz q_\phi(\mathbf{z}|\mathbf{d}) \log(p_\theta(\mathbf{d}|\mathbf{z}))}_{\mathcal{L}_1}$$

mostly too complex to be written in the closed form

solution: Monte Carlo approximation

$$\mathcal{L}_1(\theta, \phi) = \frac{1}{L} \sum_{\ell=1}^L \log(p_\theta(\mathbf{d}|\mathbf{z}_\ell))$$

samples generated according to the correct distribution

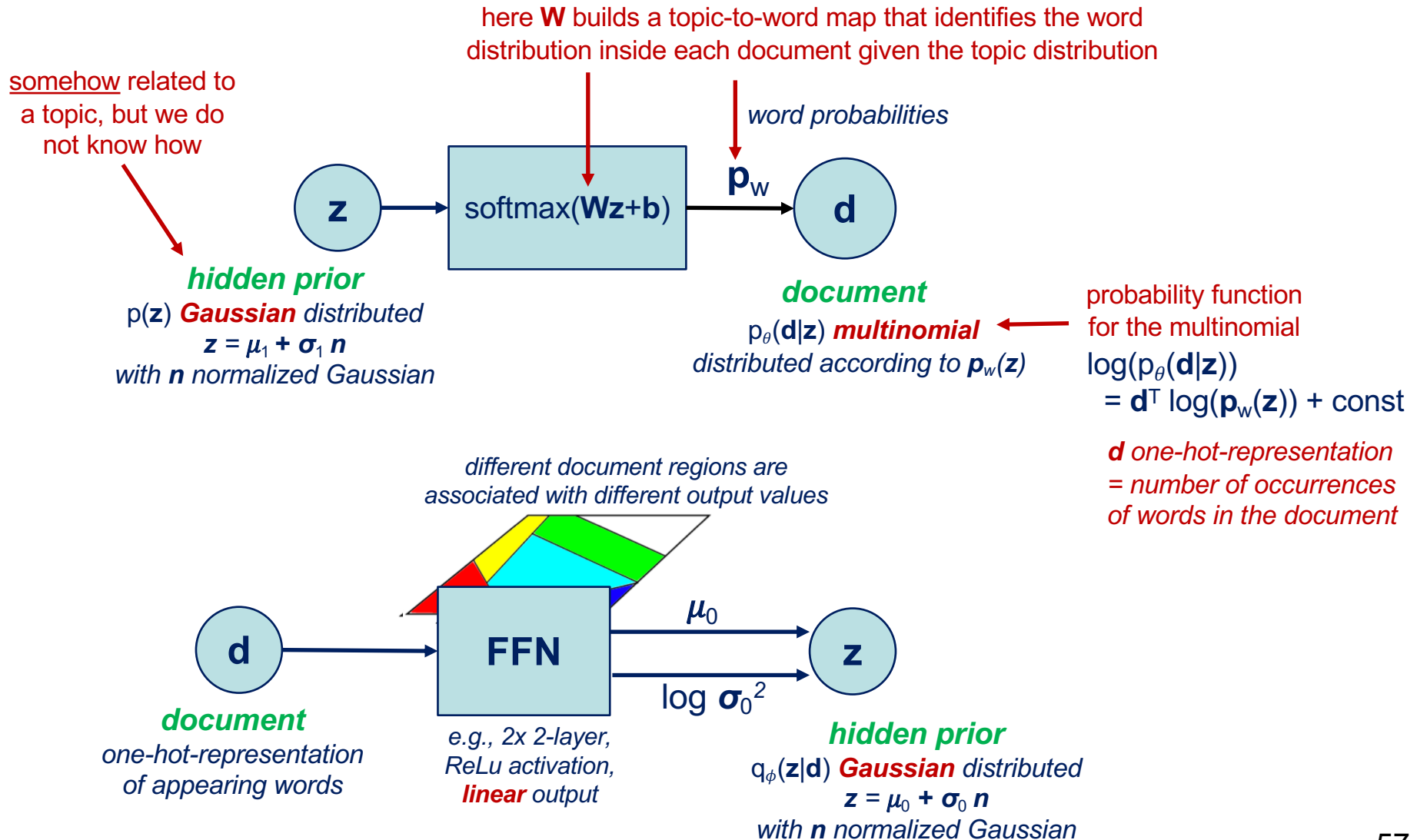
$$\mathbf{z}_\ell \sim q_\phi(\mathbf{z}|\mathbf{d})$$

e.g., the Gaussian case

$$\mathbf{z}_\ell = \boldsymbol{\mu}_\phi(\mathbf{d}) + \boldsymbol{\sigma}_\phi(\mathbf{d}) \mathbf{n}_\ell$$

need to generate these once, then use them throughout the process

$$\mathbf{n}_\ell \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$





one-hot-representation of a document = number of occurrences of words in the document

$$\mathcal{L}(\theta, \phi) = \frac{1}{L} \sum_{\ell=1}^L \sum_m \mathbf{d}_m^T \log \left( \underbrace{\text{softmax}(\mathbf{b} + \mathbf{W}(\boldsymbol{\mu}_0(\mathbf{d}_m) + \boldsymbol{\sigma}_0(\mathbf{d}_m) \mathbf{n}_{m,\ell}))}_{\text{decoder map}} \right)$$

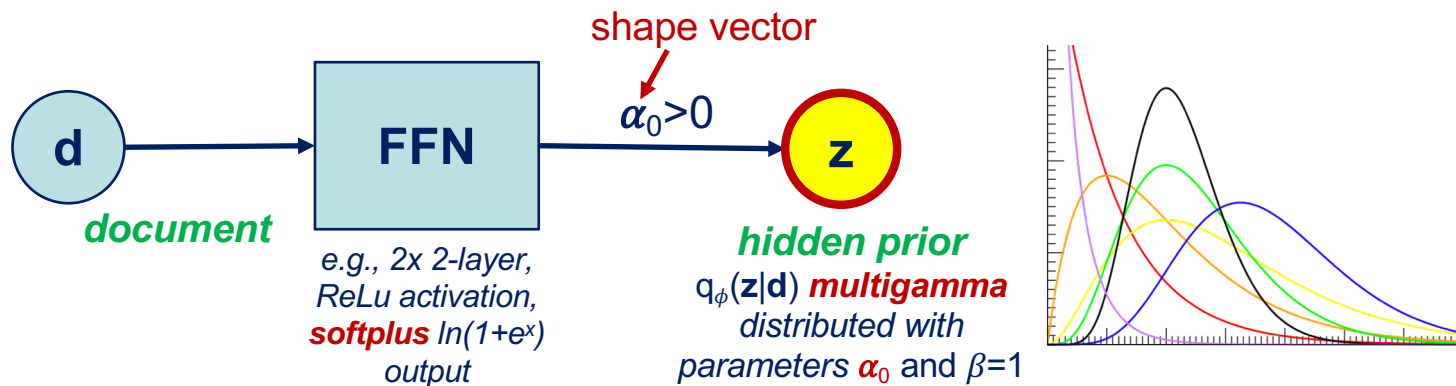
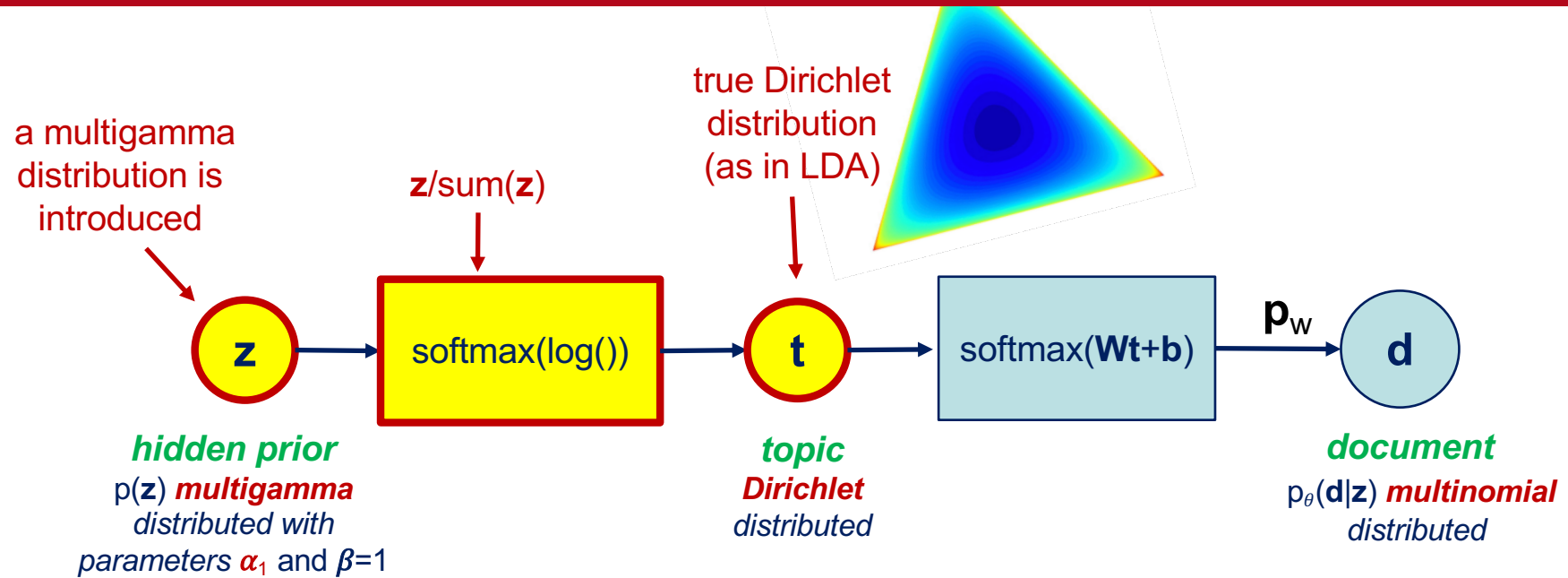
decoder model (green arrows pointing to  $\mathbf{W}$  and  $\boldsymbol{\mu}_0$ )  
encoder model (blue arrows pointing to  $\boldsymbol{\mu}_0$  and  $\boldsymbol{\sigma}_0$ )  
normalized Gaussian samples (blue arrow pointing to  $\mathbf{n}_{m,\ell}$ )

$$+ \frac{1}{2} \sum_m \sum_i 1 + \log \left( \frac{\sigma_{0,i}^2(\mathbf{d}_m)}{\sigma_{1,i}^2} \right) - \frac{\sigma_{0,i}^2(\mathbf{d}_m)}{\sigma_{1,i}^2} - \frac{(\mu_{0,i}(\mathbf{d}_m) - \mu_{1,i})^2}{\sigma_{1,i}^2}$$

a-priori model (red arrows pointing to  $\sigma_{1,i}^2$  and  $\mu_{1,i}$ )

Not very clear where the topic is, though!







one-hot-representation of a document = number of occurrences of words in the document

$$\mathcal{L}(\theta, \phi) = \frac{1}{L} \sum_{\ell=1}^L \sum_m \mathbf{d}_m^T \log \left( \underbrace{\text{softmax}(\mathbf{b} + \mathbf{W} \text{softmaxlog}(\mathbf{f}(\mathbf{u}_{m,\ell}, \boldsymbol{\alpha}_0(\mathbf{d}_m))))}_{\text{decoder map}} \right)$$

$$f(\mathbf{u}, \boldsymbol{\alpha}) = (\mathbf{u} \boldsymbol{\alpha} \Gamma(\boldsymbol{\alpha}_1))^{1/\alpha}$$

approx. uniform to multigamma map

normalized uniform samples

decoder model

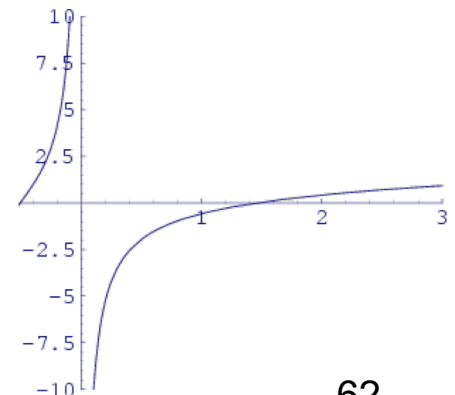
encoder model

$$+ \sum_m \sum_i \log \left( \frac{\Gamma(\alpha_{0,i}(\mathbf{d}_m))}{\Gamma(\alpha_{1,i})} \right) - (\alpha_{0,i}(\mathbf{d}_m) - \alpha_{1,i}) \psi(\alpha_{0,i}(\mathbf{d}_m))$$

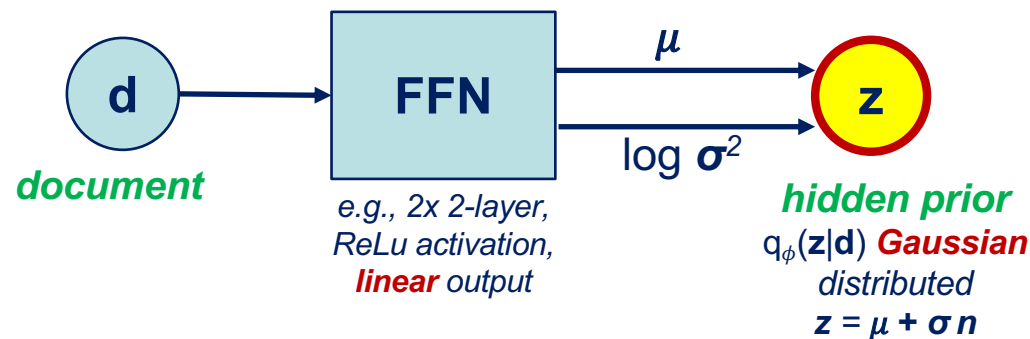
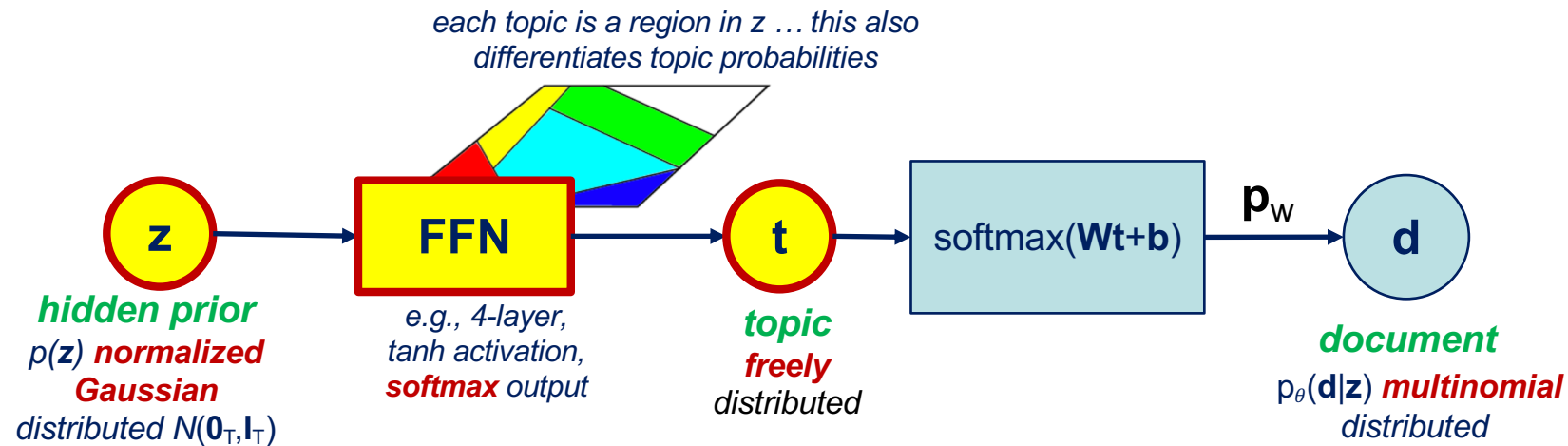
a-priori model

digamma function

$$\psi(x) = \frac{\Gamma'(x)}{\Gamma(x)}$$



Now we know where the topic is!



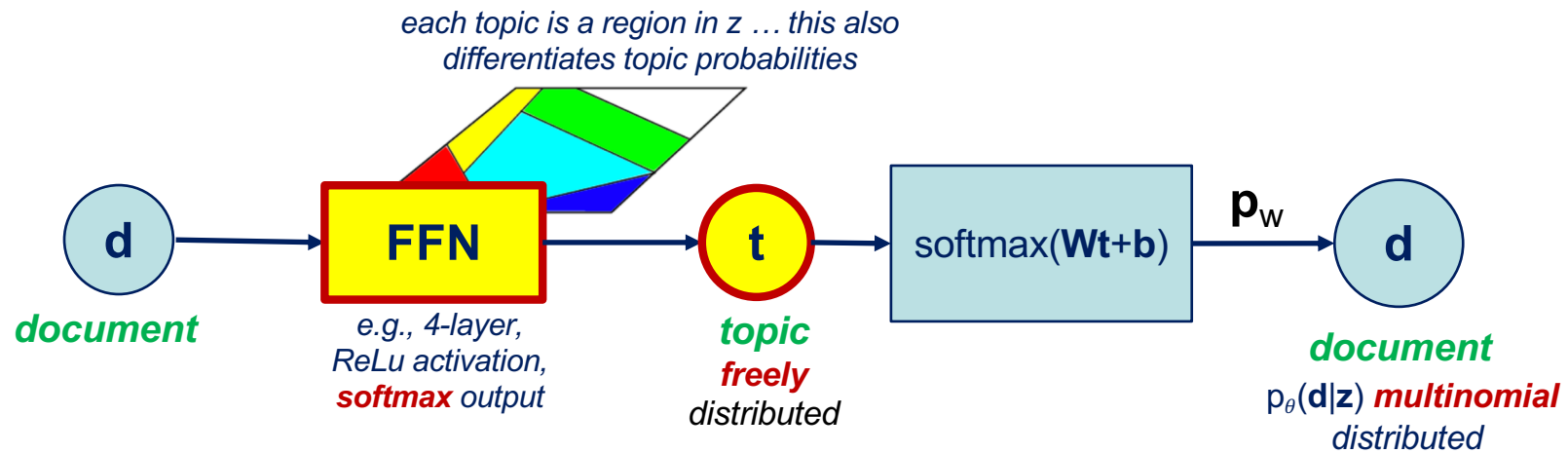


one-hot-representation of a document = number of occurrences of words in the document

$$\mathcal{L}(\theta, \phi) = \frac{1}{L} \sum_{\ell=1}^L \sum_m \mathbf{d}_m^T \log \left( \underbrace{\text{softmax}(\mathbf{b} + \mathbf{W} \text{FFN}_1)}_{\text{decoder map}}(\underbrace{\mu_0(\mathbf{d}_m) + \sigma_0(\mathbf{d}_m) \mathbf{n}_{m,\ell}}_{\text{normalized Gaussian samples}})) \right) + \frac{1}{2} \sum_m \sum_i 1 + \log(\sigma_{0,i}^2(\mathbf{d}_m)) - \sigma_{0,i}^2(\mathbf{d}_m) - (\mu_{0,i}(\mathbf{d}_m))^2$$

Our estimate of the **topic distribution** for the  $m$ th document!

$$\mathbf{c}_m = \frac{1}{L} \sum_{\ell=1}^L \text{FFN}_1(\mu_0(\mathbf{d}_m) + \sigma_0(\mathbf{d}_m) \mathbf{n}_{m,\ell})$$



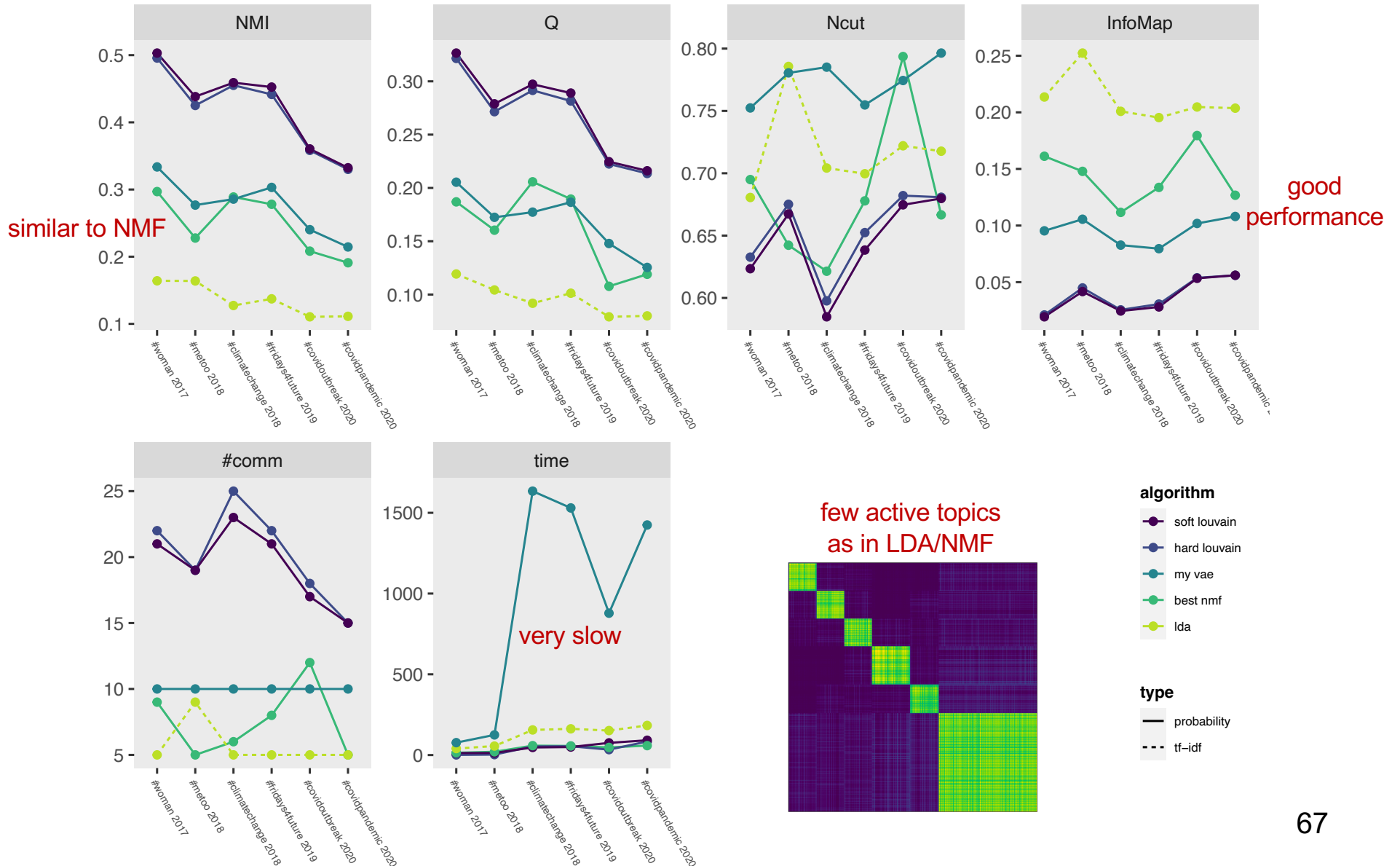
one-hot-representation of a document = number of occurrences of words in the document

decoder model      encoder model

$$\mathcal{L}(\theta, \phi) = \sum_m \mathbf{d}_m^T \log(\text{softmax}(\mathbf{b} + \mathbf{W} \text{FFN}(\mathbf{d}_m)))$$



# A comparison with NMF, LDA, and Louvain





- ❑ Naturally provides a soft topic assignment
- ❑ VAE – interesting approach
  - more flexible model than NMF or LDA
  - gives improvements
- ❑ Comparison – with Louvain
  - still far away
  - would be nice to see other **Deep Learning** approaches
  - ... your task! 😊

# Transformer Architecture

with application to BERT, RoBERTa, OpenAI GPT





## ☰ Attention (machine learning)

---

Article [Talk](#)

---

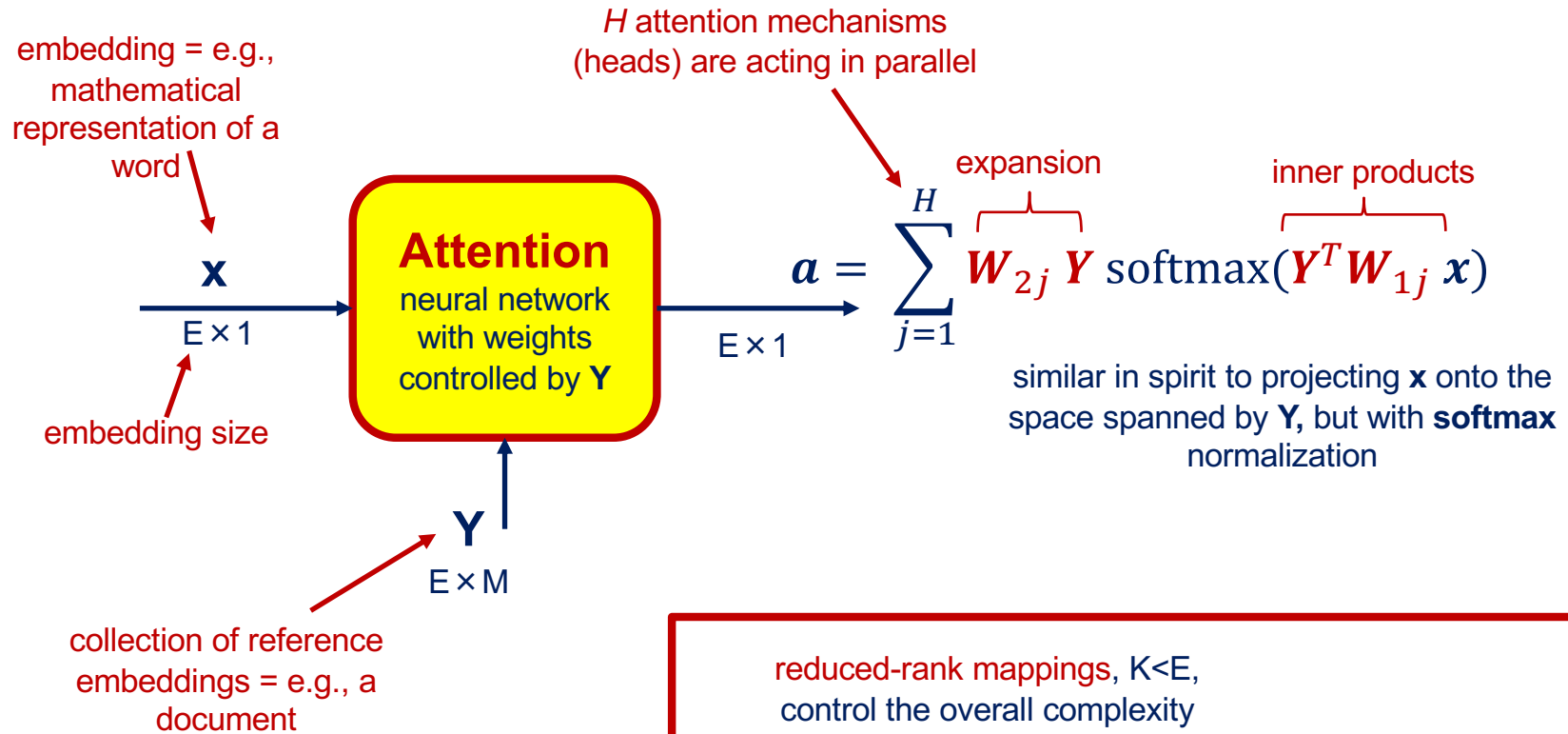
From Wikipedia, the free encyclopedia

In [artificial neural networks](#), **attention** is a technique that is meant to mimic [cognitive attention](#). This effect **enhances** some **parts of the input data while diminishing other parts** — the motivation being that the network should devote more focus to the important parts of the data, even though they may be small portion of an image or sentence. Learning which part of the data is more important than another depends on the context, and this is trained by [gradient descent](#).



# The Attention Module

Vaswani, Ashish, et al. "Attention is all you need" (2017)



reduced-rank mappings,  $K < E$ , control the overall complexity

$$\mathbf{W}_{ij} = \mathbf{V}_{ij} \mathbf{U}_{ij}$$

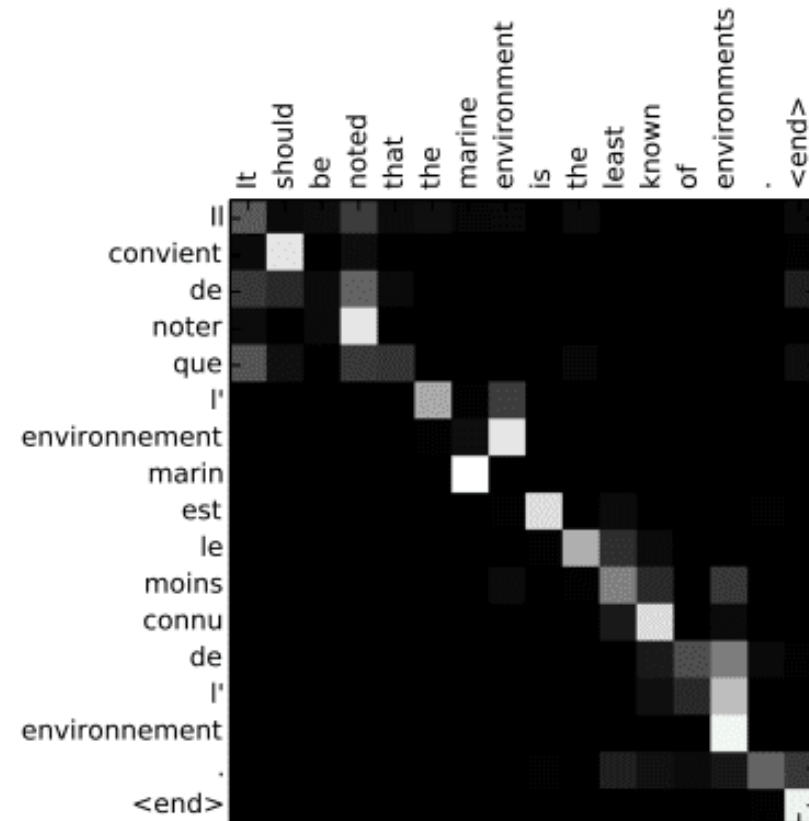
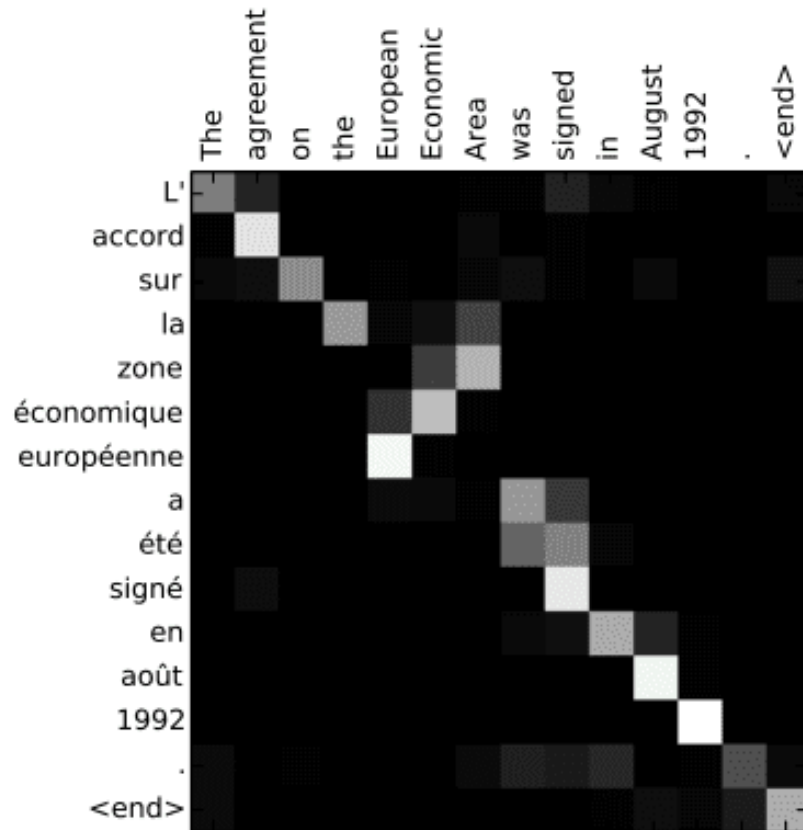
$E \times E$        $E \times K$        $K \times E$

overall complexity  $4EKH$

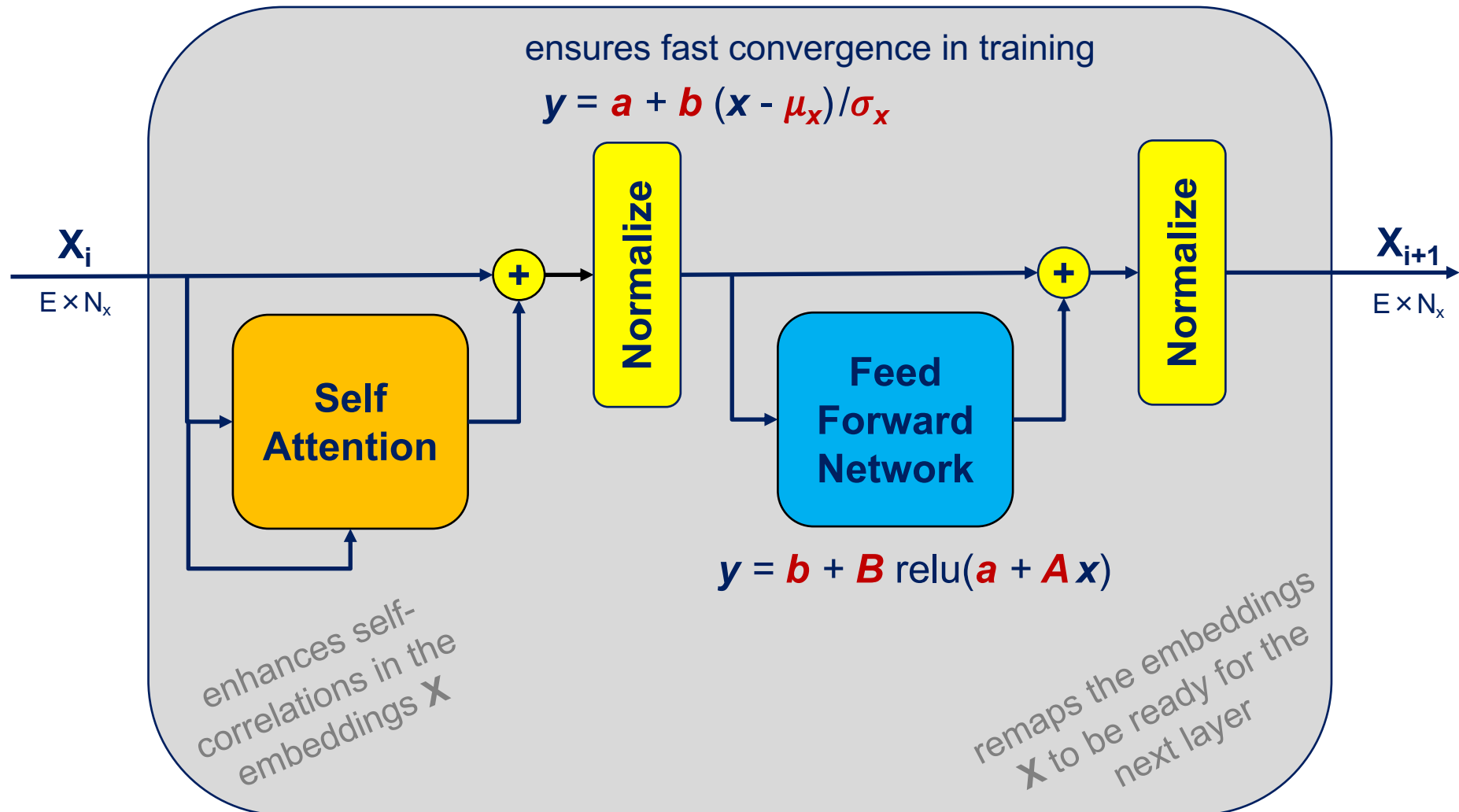


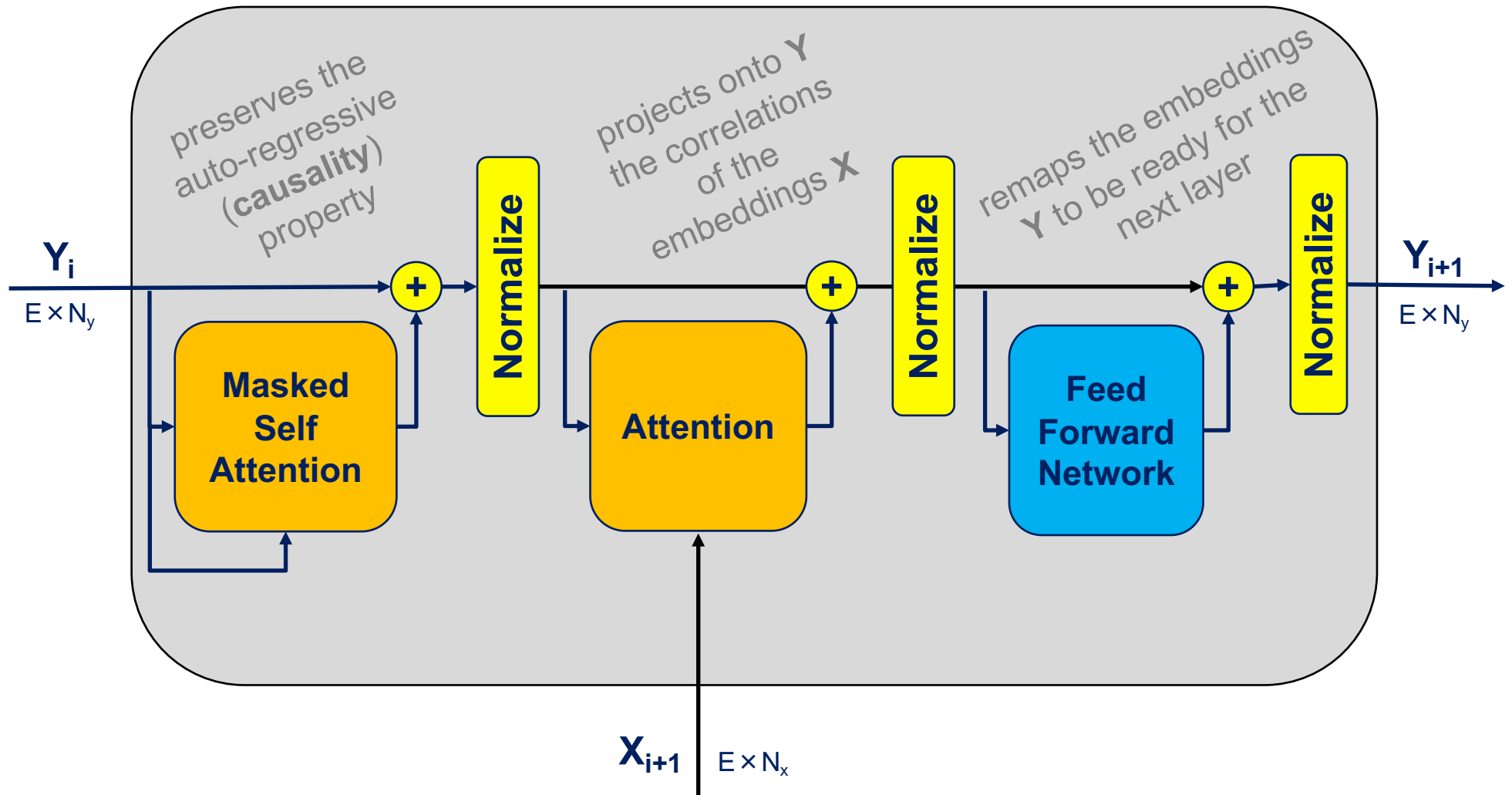
# Visualizing Attention

in a translation experiment (X English, Y French)



$$\text{softmax}(Y^T W_{1j} X)$$

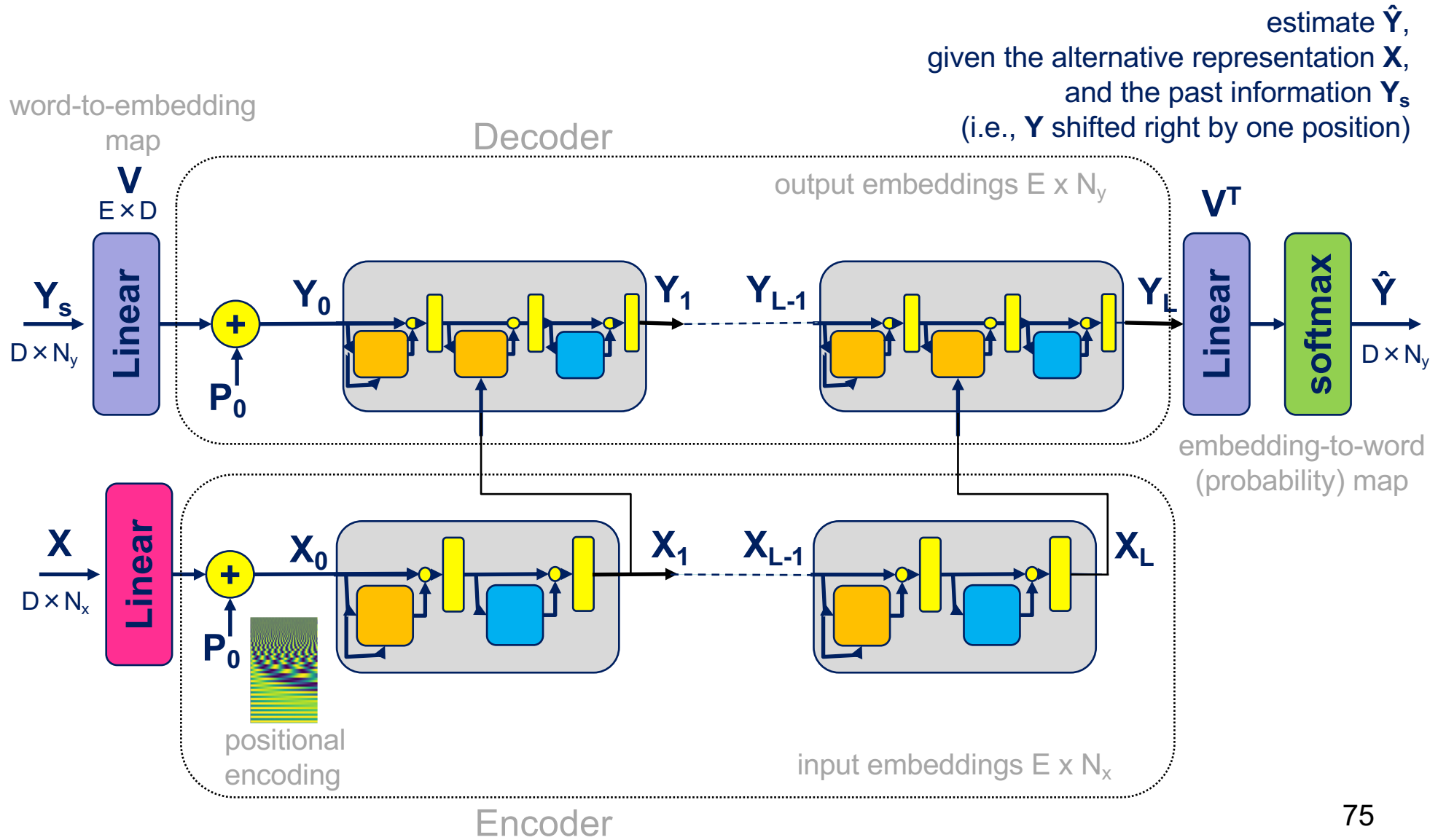






# Transformer Architecture

Vaswani, Ashish, et al. "Attention is all you need" (2017)  
Google's patent <https://patents.google.com/patent/US10452978B2/en>





---

## The Annotated Transformer

---

Apr 3, 2018

-----

There is now a [new version](#) of this blog post updated for modern PyTorch.

-----

```
from IPython.display import Image
Image(filename='images/aiayn.png')
```

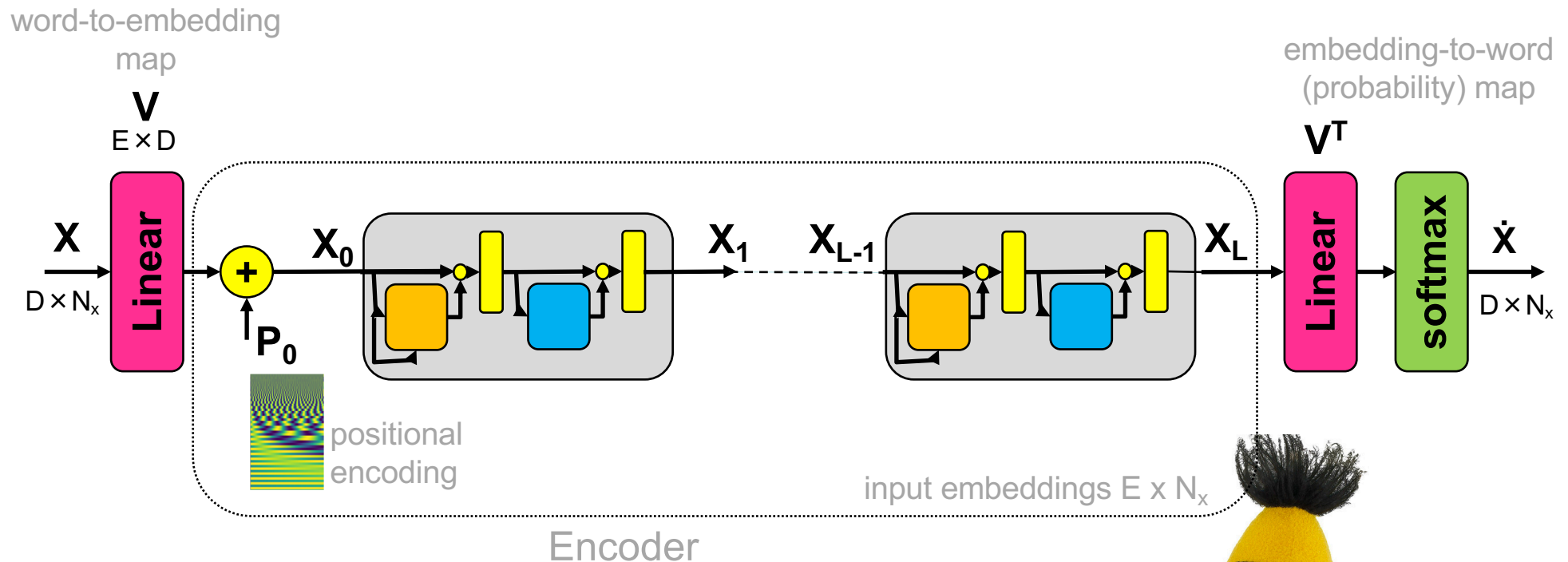
---

### Attention Is All You Need

---



<https://github.com/google-research/bert>







# BERT parameters

	Embeddings size E	Self-attention heads H	Head dimension K = E/H	FFN inner size I = 4E	Parameters per layer $12E^2+9E$	Layers L	Dictionary size D	Total parameters
BERT base	768	12	64	3072	7.1M	12	30.5K	110M
BERT large	1024	16	64	4096	12.6M	24	30.5K	340M



max tokens  $N_x = 512$

Created by researchers at Google AI Language



# BERT pre-training procedure

BooksCorpus (800M words) + English Wikipedia (2,500M words)

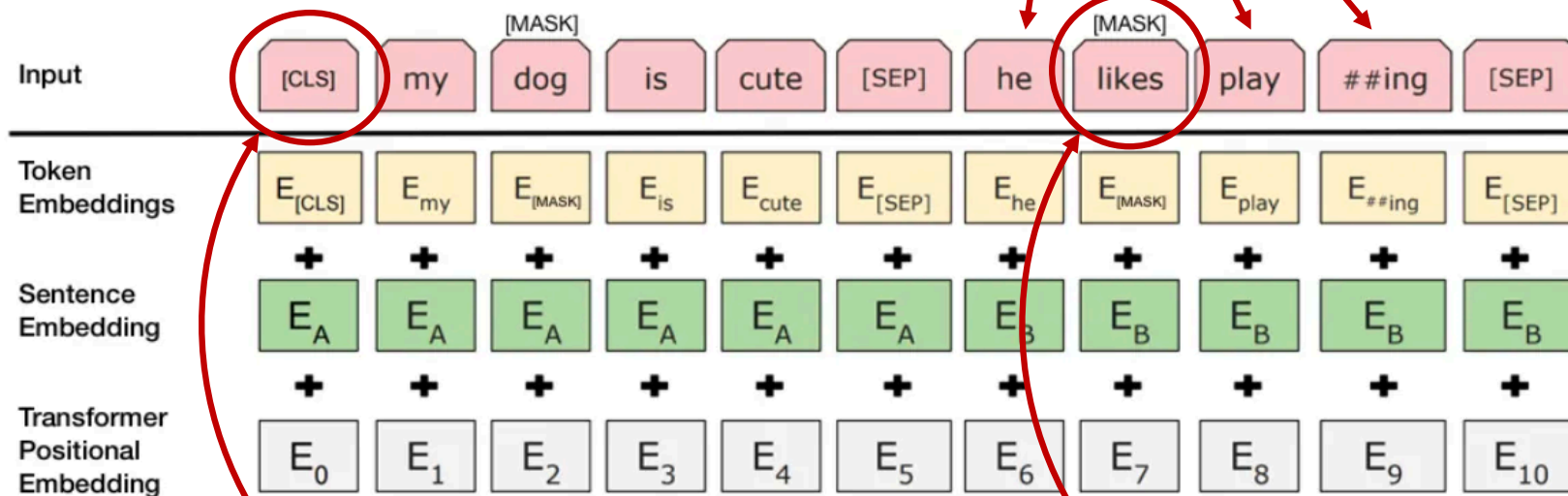
## Masked Language Model

15% **masked tokens** replaced with:

- [MASK] token (80% of the times)
- Original token (10%)
- Random token (10%)

## Next Sequence Prediction

- Next sequence (50% of the times)
- Random sequence (50%)



Output [CLS] fed into an additional output layer for softmax classification (of correct/wrong next sequence)

Output **masked tokens** fed into the output layer  $V^T$  and evaluated for probability of correct estimate



## **Larger training corpora** (10x larger)

*training on BookCorpus + Wikipedia and also CC-News, OpenWebText, Stories*

## **Dynamic masking**

*training data was duplicated 10 times so that each sequence is masked in 10 different ways over the 40 epochs of training*

## **Full-sentences without NSP loss**

*full sentences sampled contiguously from one or more documents, such that the total length is at most 512 tokens*

## **Large mini-batches**

## **A larger byte-level BPE (byte pair encoding)** of 50K subword units

*a hybrid between character- and word-level representations that allows handling the large vocabularies common in natural language corpora*

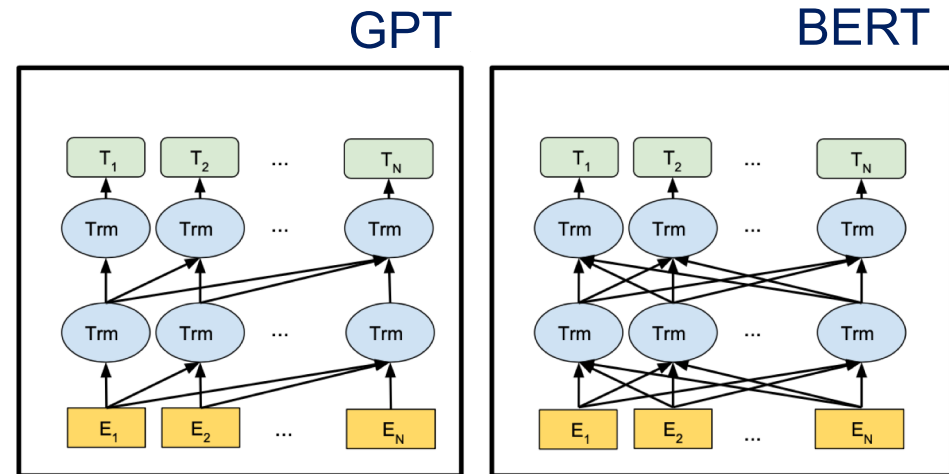
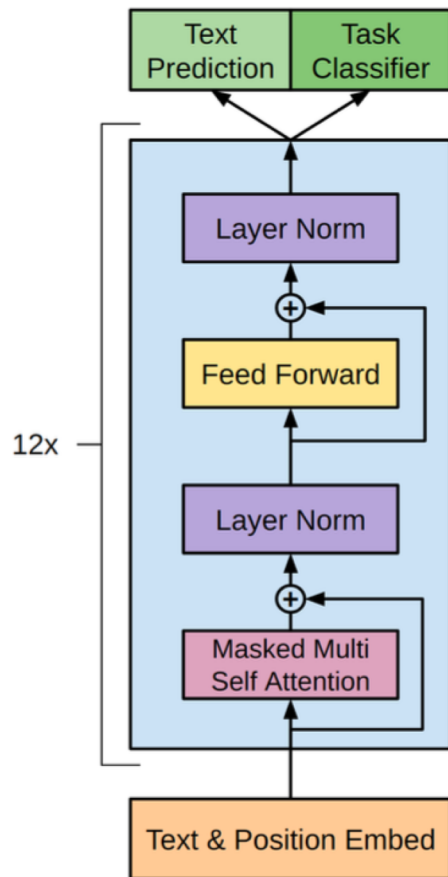


# Generative Pre-Training (GPT)

Radford, Alec, et al. "Improving language understanding by generative pre-training." (2018)

(unsupervised) pre-training on **Language Modelling (no mask)**

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta)$$



same parameters of BERT-base, but with **Masked Attention** trained on BookCorpus only



## McCann et al. (2018)

language provides a flexible way to specify tasks, inputs, and outputs all as a sequence of symbols... it is therefore possible to **train a single model** with **sufficient capacity** to infer and perform many **different tasks**

model gets complex!



data gets larger!

Parameters	Layers	$d_{model}$	
117M	12	768	GPT, BERT-base
345M	24	1024	BERT-large
762M	36	1280	
1542M	48	1600	GPT-2

## WebText

scraping all outbound links (45M links) from Reddit, a social media platform, which received at least 3 karma – exclude Wikipedia



increasingly larger data and model!

Model Name	$n_{\text{params}}$	$n_{\text{layers}}$	$d_{\text{model}}$	$n_{\text{heads}}$	$d_{\text{head}}$	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	$6.0 \times 10^{-4}$
GPT-3 Medium	350M	24	1024	16	64	0.5M	$3.0 \times 10^{-4}$
GPT-3 Large	760M	24	1536	16	96	0.5M	$2.5 \times 10^{-4}$
GPT-3 XL	1.3B	24	2048	24	128	1M	$2.0 \times 10^{-4}$
GPT-3 2.7B	2.7B	32	2560	32	80	1M	$1.6 \times 10^{-4}$
GPT-3 6.7B	6.7B	32	4096	32	128	2M	$1.2 \times 10^{-4}$
GPT-3 13B	13.0B	40	5140	40	128	2M	$1.0 \times 10^{-4}$
GPT-3 175B or "GPT-3"	175.0B	96	12288	96	128	3.2M	$0.6 \times 10^{-4}$

**Layer normalization at the input** (plus one at the output)

## **Sparse attention patterns**

*alternating dense and locally banded sparse attention patterns in the layers*

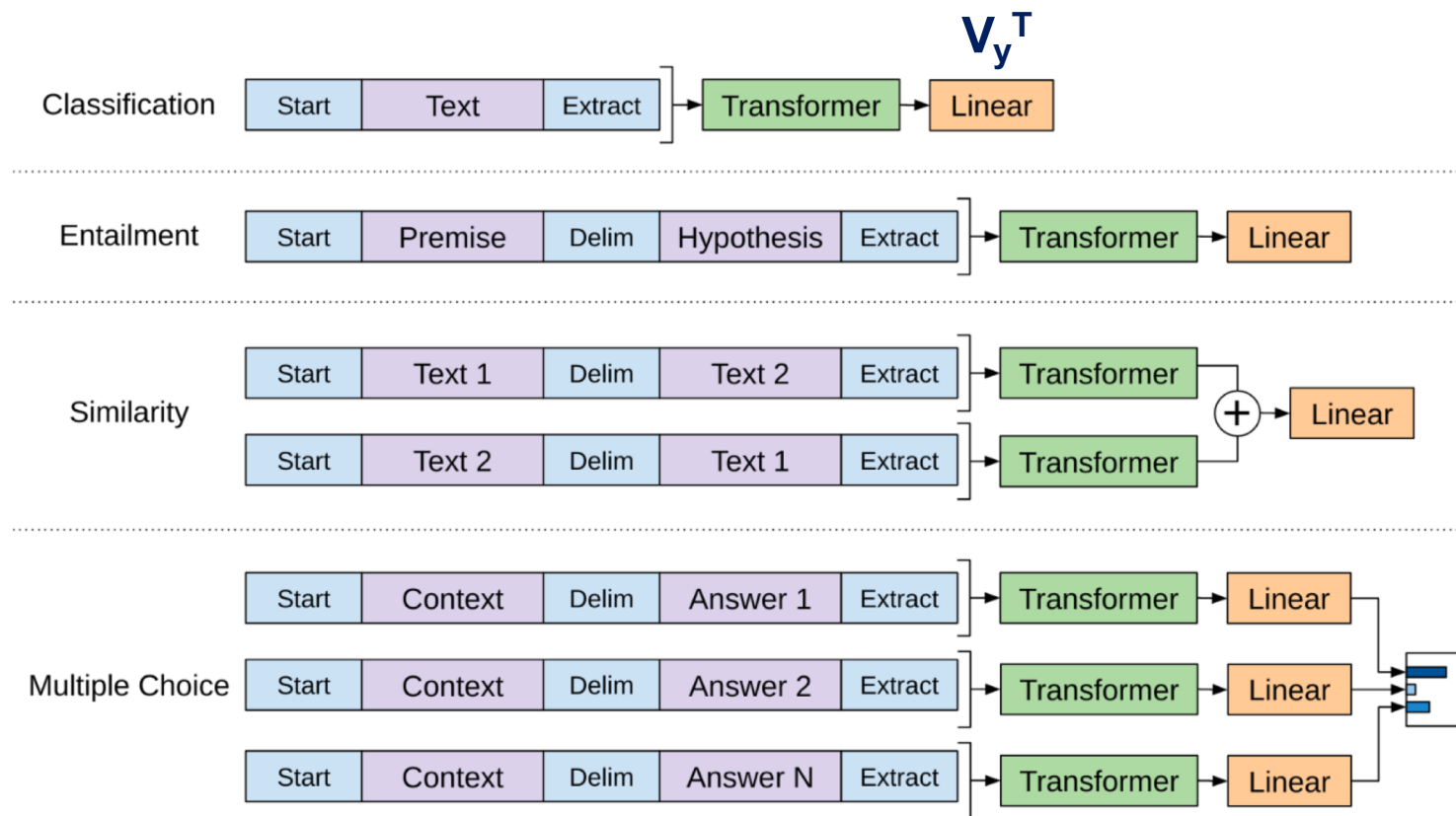
## **Byte-level BPE (byte pair encoding)** of 50K subword units

*also prevent BPE from merging across character categories (to avoid dog, dog!, dog?)*

## **Modified initialization**



$$\log \text{softmax}(\mathbf{V}_y^T \mathbf{X}_L) \quad \xrightarrow{\quad} \quad L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C}) \quad \xleftarrow{\quad} \quad \text{Language Modelling loss}$$





Task	Description	Possible approach
Masked language prediction	predict masked words in a text	This is what BERT model is pre-trained for
Text classification or Sentiment analysis	assign a label to a given sequence of text	Apply <b>linear transform+softmax on K classes</b> , and train the model for the specific classification task
Text translation	translate a text	Need to pre-train a full Transformer Architecture for this task
Summarization	generate a summary of a document	GPT example: context given by a document; then generate 100 tokens by <b>top-2 random sampling</b> (Fan et al., 2018), i.e., take at each step the most likely next word at random among the top-2 candidates; finally select first 3 sentences as abstract
Question answering	answer a question	GPT example: the context of the language model is seeded with example question answer pairs which helps the model infer the <b>short answer style</b> of the dataset
Document question answering	answer a question on a given text	GPT example: context seeded by a text; then as for question answering
Conversational	ChatBot	InstructGPT/ChatGPT: Fine-tuned models using <b>reinforcement learning</b> from human feedback





## Hugging Face

<https://huggingface.co/docs/transformers/v4.29.1/en/index>

State-of-the-art Machine Learning  
for PyTorch, TensorFlow, and JAX



PyTorch



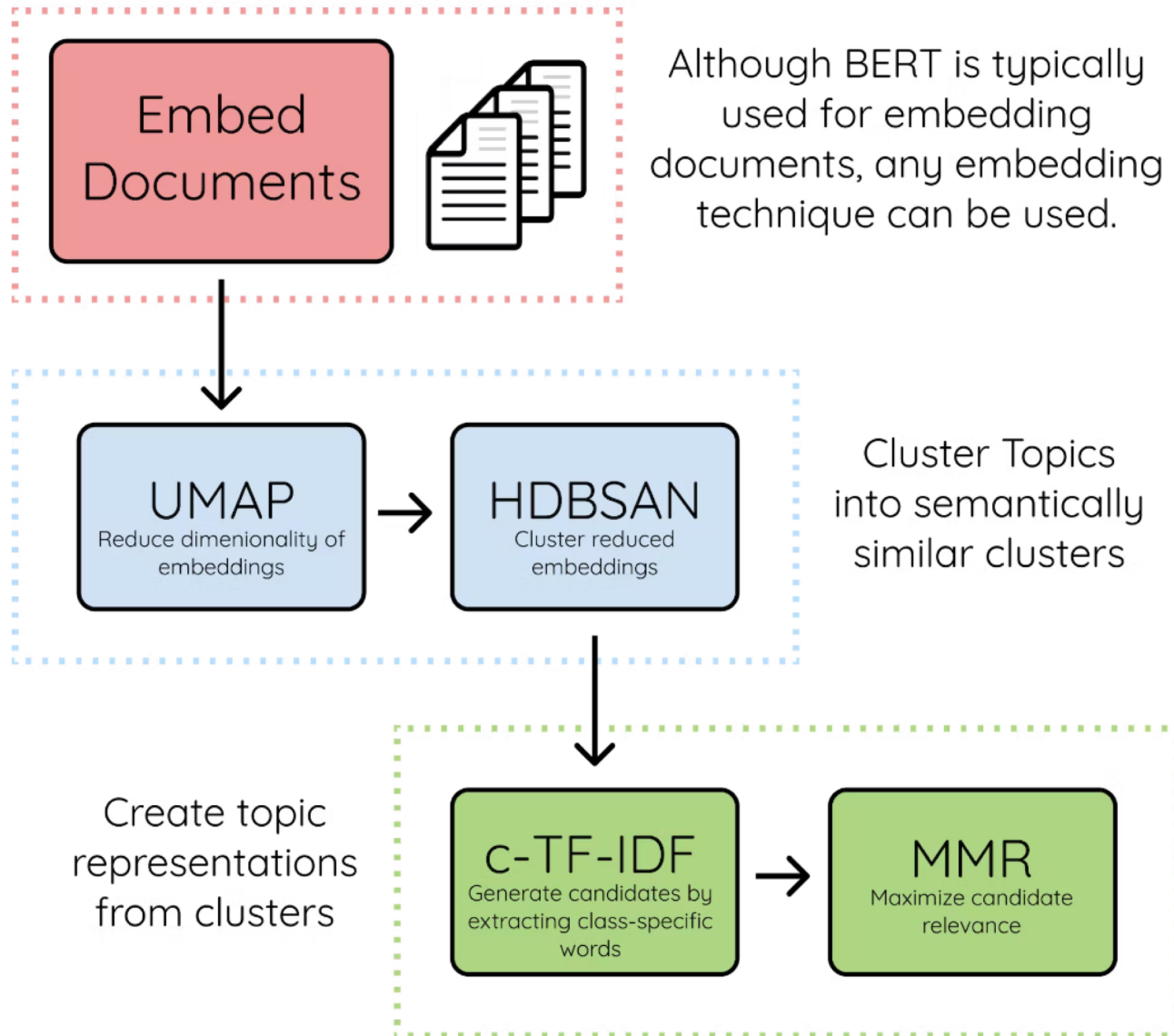
TensorFlow



ALBERT, BART, **BERT**, BigBird, BigBird-Pegasus, BioGpt, BLOOM, CamemBERT, CANINE, ConvBERT, CTRL, Data2VecText, DeBERTa, DeBERTa-v2, DistilBERT, ELECTRA, ERNIE, ErnieM, ESM, FlauBERT, FNet, Funnel Transformer, GPT-Sw3, **OpenAI GPT-2**, GPTBigCode, GPT Neo, GPT NeoX, GPT-J, I-BERT, LayoutLM, LayoutLMv2, LayoutLMv3, LED, LiLT, LLaMA, Longformer, LUKE, MarkupLM, mBART, MEGA, Megatron-BERT, MobileBERT, MPNet, MVP, Nezha, Nyströmformer, OpenLlama, **OpenAI GPT**, OPT, Perceiver, PLBart, QDQBert, Reformer, RemBERT, **RoBERTa**, RoBERTa-PreLayerNorm, RoCBert, RoFormer, SqueezeBERT, TAPAS, Transformer-XL, XLM, XLM-RoBERTa, XLM-RoBERTa-XL, XLNet, X-MOD, YOSO

# BERT Topic

exploiting embeddings for topic detection





```
!pip install bertopic
from bertopic import BERTopic
from sentence_transformers import SentenceTransformer
```

```
sentence_model = SentenceTransformer("all-MiniLM-L6-v2")
bert_model = BERTopic(embedding_model=sentence_model,
                      min_topic_size=20, nr_topics='auto')
```

initialise model

```
docs = list(df2["text_sup_clean"])
topics, probabilities = bert_model.fit_transform(docs)
```

fit model

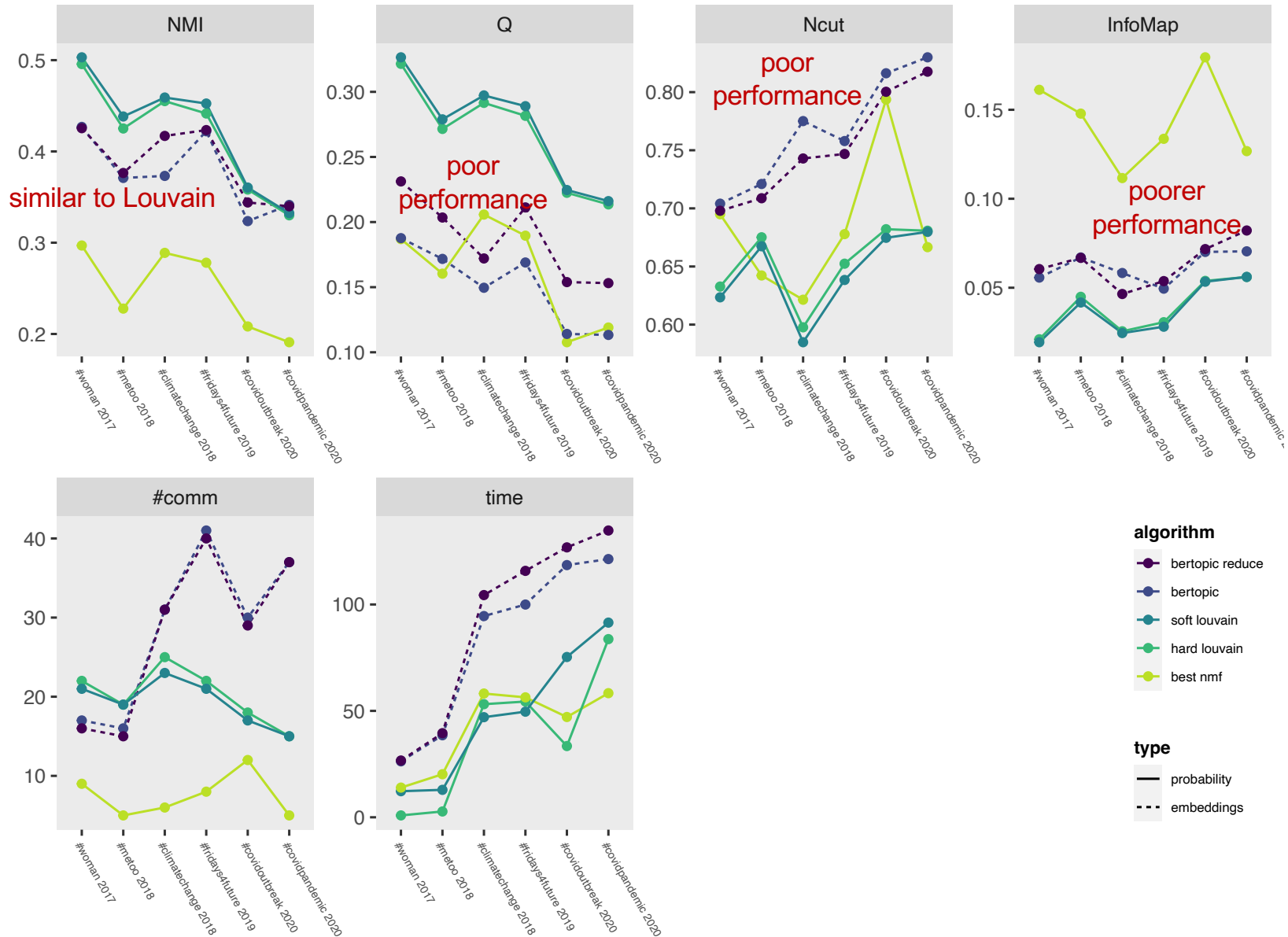
```
topics = bert_model.reduce_outliers(docs, topics)
```

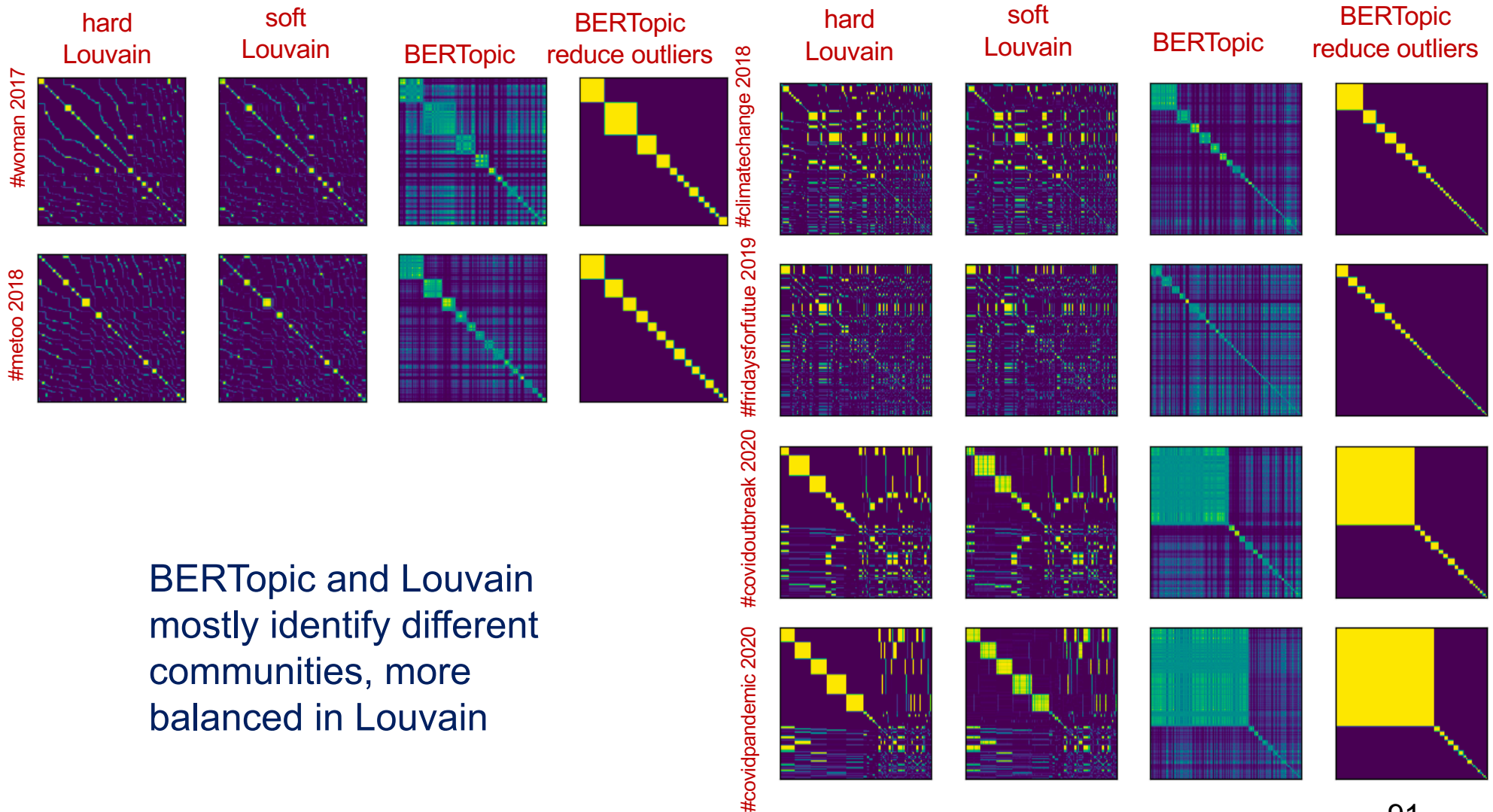
reduce outliers

```
# extract community assignments
C = sps.csr_matrix((len(topics), max(topics)+2))
for i in range(C.shape[1]):
    C[np.array(topics)==(i-1), i] = 1

# remove zero assignments
C = C[:, np.unique(scipy.sparse.find(C)[1])]
```

extract C from topic  
assignment



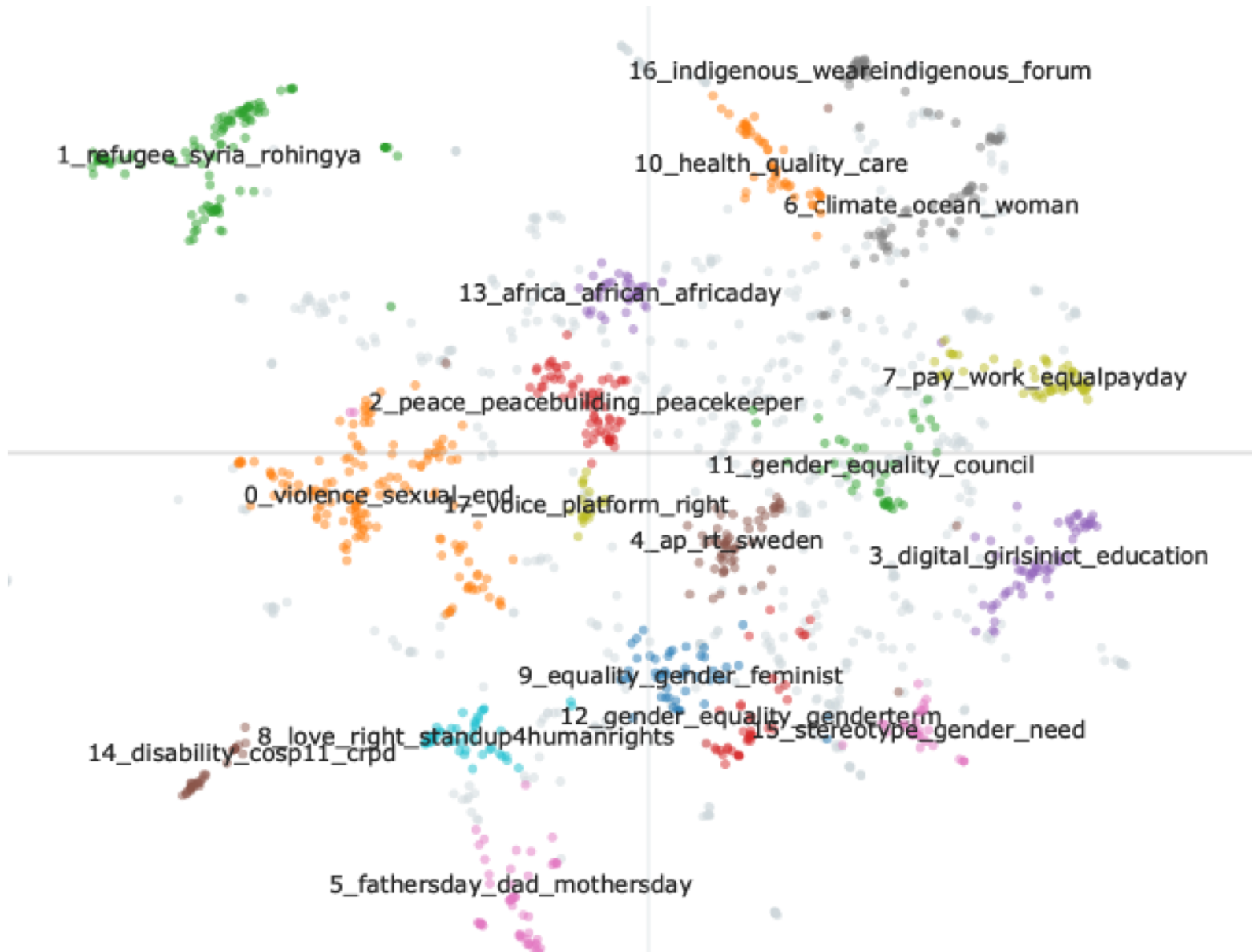


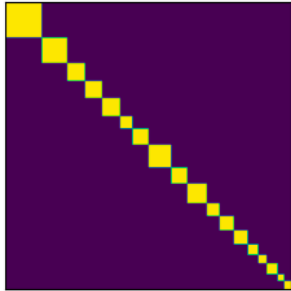
BERTopic and Louvain  
mostly identify different  
communities, more  
balanced in Louvain



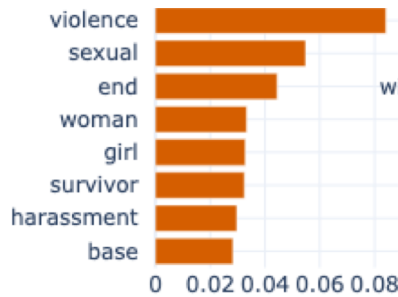
# bert\_model.visualize\_documents(docs)

#metoo2018





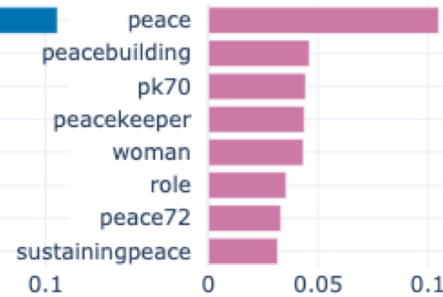
### sexual violence



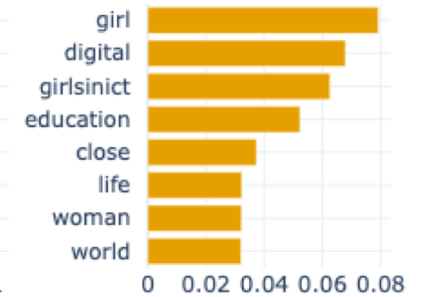
### refugees



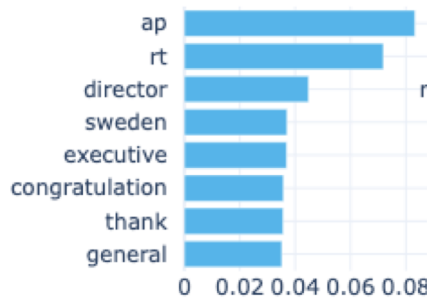
### peace



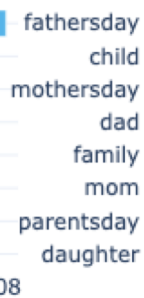
### girlsinct



### executive director



### mothersday



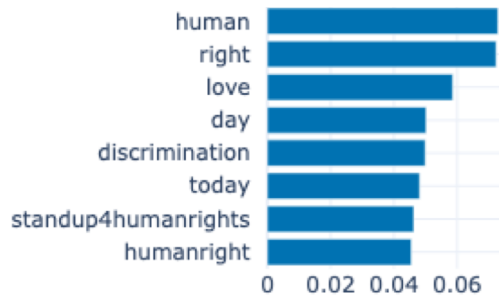
### sustainability



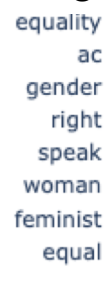
### equal pay



### discrimination



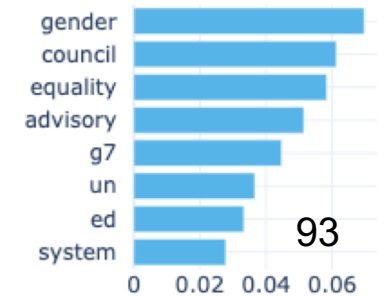
### gender equality



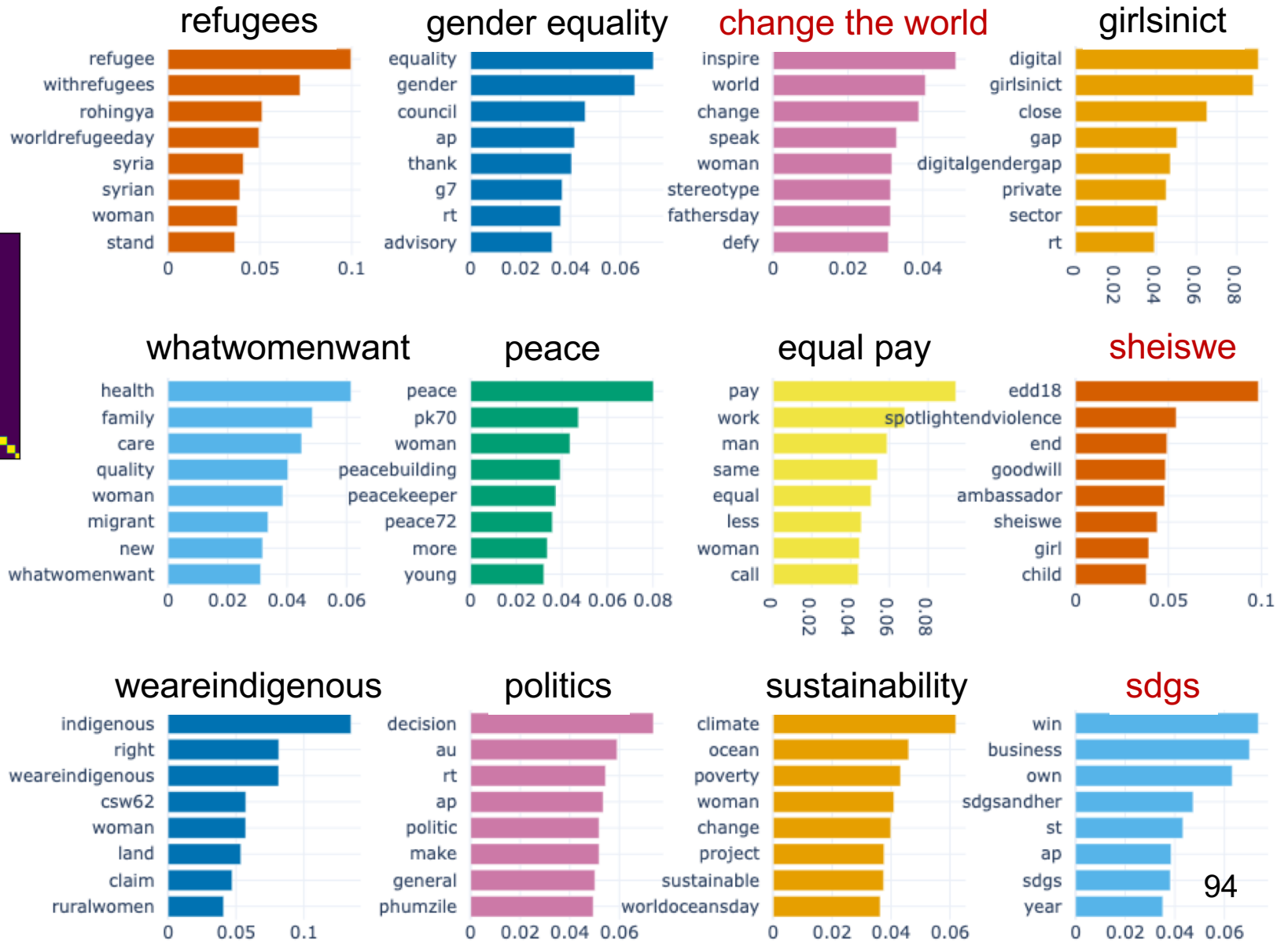
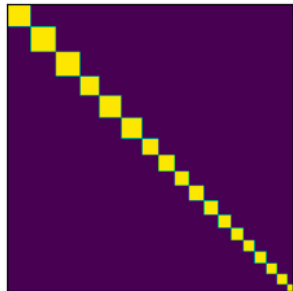
### whatwomenwant



### politics









# How to use BERTopic barchart

with your own topic assignment **C** in python

```
def bertopic_overwrite(bert_model_in, docs, C):
    bert_model = copy.deepcopy(bert_model_in)

    # build the documents dataframe: 'Document' + "Topic"
    documents = pd.DataFrame(docs, columns=['Document'])
    tmp = np.array([C[i].argmax() for i in range(C.shape[0])])
    documents["Topic"] = tmp

    # update topic assignment
    bert_model.topics_ = tmp.tolist()

    # build cf-idf values
    documents_per_topic = documents.groupby(['Topic'],
                                           as_index=False).agg({'Document': ' '.join})
    c_tf_idf_, words = bert_model._c_tf_idf(documents_per_topic)
    bert_model.c_tf_idf_ = c_tf_idf_

    # extract words representations
    topic_representations_ = bert_model._extract_words_per_topic(words, documents)
    bert_model.topic_representations_ = topic_representations_
    bert_model.topic_labels_ = {key: f"{key}_" + "_".join([word[0] for word in values[:4]])
                               for key, values in
                               bert_model.topic_representations_.items()}

    # exit
    return bert_model
```



- ❑ Naturally provides a hard topic assignment
- ❑ **Useful** tool
- ❑ More readable output with deep cleaned text but same performance
- ❑ Comparison – with Louvain
  - weaker in general, especially in modularity
  - equivalent NMI** = relevant topics
  - lower modularity** = the documents that identify the topics are less distinguishable
  - higher complexity involved
  - less balanced topics, but generally meaningful
  - topics correlated with Louvain

# Sentiment analysis

adding useful insights to your data



- ❑ Sentiment – e.g., positive, negative, neutral  
enduring cognitive content that defines the affective state
- ❑ Emotion – e.g., anger, disgust, fear, joy, sadness  
intense affective state of short duration with a precise cause
- ❑ Ingroup bias – e.g., use of pronouns I, we, us  
tendency to favor one's own group over other groups
- ❑ Outgroup bias – e.g., use of pronoun they  
tendency to dislike members of groups we don't identify with
- ❑ Agency – e.g., use of action verbs do, take, make  
perception that an individual is able to contribute to/a group  
can collectively reach a social change

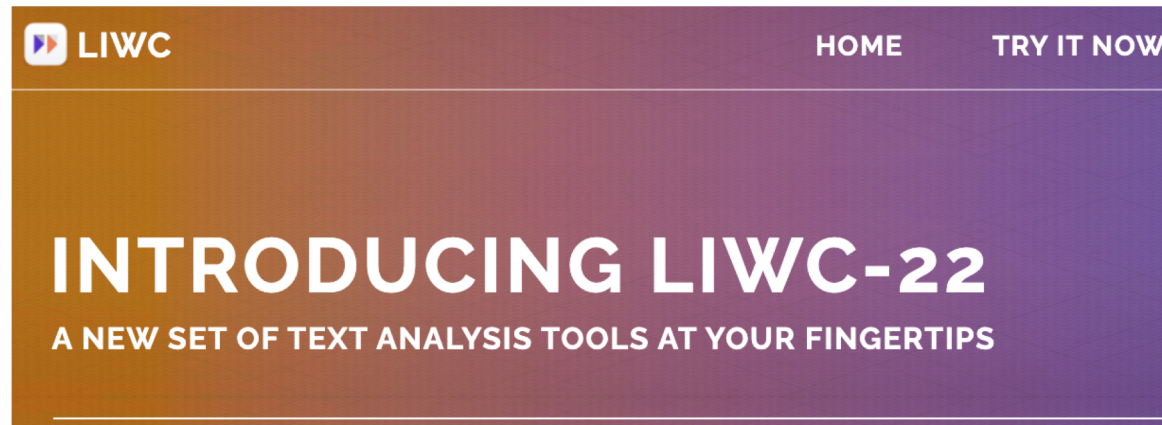


UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA

# LIWC linguistic inquiry and word count

Tausczik, Pennebaker. "The psychological meaning of words:  
LIWC and computerized text analysis methods." (2010)

<https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=79d2494cc10a9633c42115df84bb74ed447080f6>



<https://www.liwc.app/>

- ❑ word count (or **dictionary**) methodology
- ❑ over 60 dictionaries coded and **validated** for their accuracy in reflecting psychological content
- ❑ **simplicity** of implementation and usage
- ❑ state-of-the-art in psychology
- ❑ **one licence** available in the instructor's PC 😊



Category	Examples	Words in Category	Psychological Correlates
<i>Linguistic processes</i>			
Word count			Talkativeness, verbal fluency
Words/sentence			Verbal fluency, cognitive complexity
Dictionary words	(Percentage of all words captured by the program)		Informal, nontechnical language
Words >6 letters	(Percentage of all words longer than 6 letters)		Education, social class
Total function words		464	
Total pronouns	I, them, itself	116	Informal, personal
Personal pronouns	I, them, her	70	Personal, social
First-person singular	I, me, mine	12	Honest, depressed, low status, personal, emotional, informal
First-person plural	We, us, our	12	Detached, high status, socially connected to group (sometimes)
Second person	You, your, thou	20	Social, elevated status
Third-person singular	She, her, him	17	Social interests, social support
Third-person plural	They, their, they'd	10	Social interests, out-group awareness (sometimes)

ingroup



outgroup





# LIWC categories

goal orientation, aggression, social concern, emotionality

Category	Examples	Words in Category	Psychological Correlates
Indefinite pronouns	It, it's, those	46	Use of concrete nouns, interest in objects and things
Articles	A, an, the	3	
Common verbs	Walk, went, see	383	Informal, passive voice Focus on the past
Auxiliary verbs	Am, will, have	144	
Past tense	Went, ran, had	145	
Present tense	Is, does, hear	169	Living in the here and now
Future tense	Will, gonna	48	Future and goal oriented
Adverbs	Very, really, quickly	69	Education, concern with precision
Prepositions	To, with, above	60	
Conjunctions	And, but, whereas	28	Inhibition
Negations	No, not, never	57	
Quantifiers	Few, many, much	89	
Numbers	Second, thousand	34	Informal, aggression,
Swear words	Damn, piss, fuck	53	
<i>Psychological processes</i>			
Social processes	Mate, talk, they, child	455	Social concerns, social support
Family	Daughter, husband	64	Emotionality
Friends	Buddy, friend, neighbor	37	
Humans	Adult, baby, boy	61	
Affective processes	Happy, cried, abandon	915	

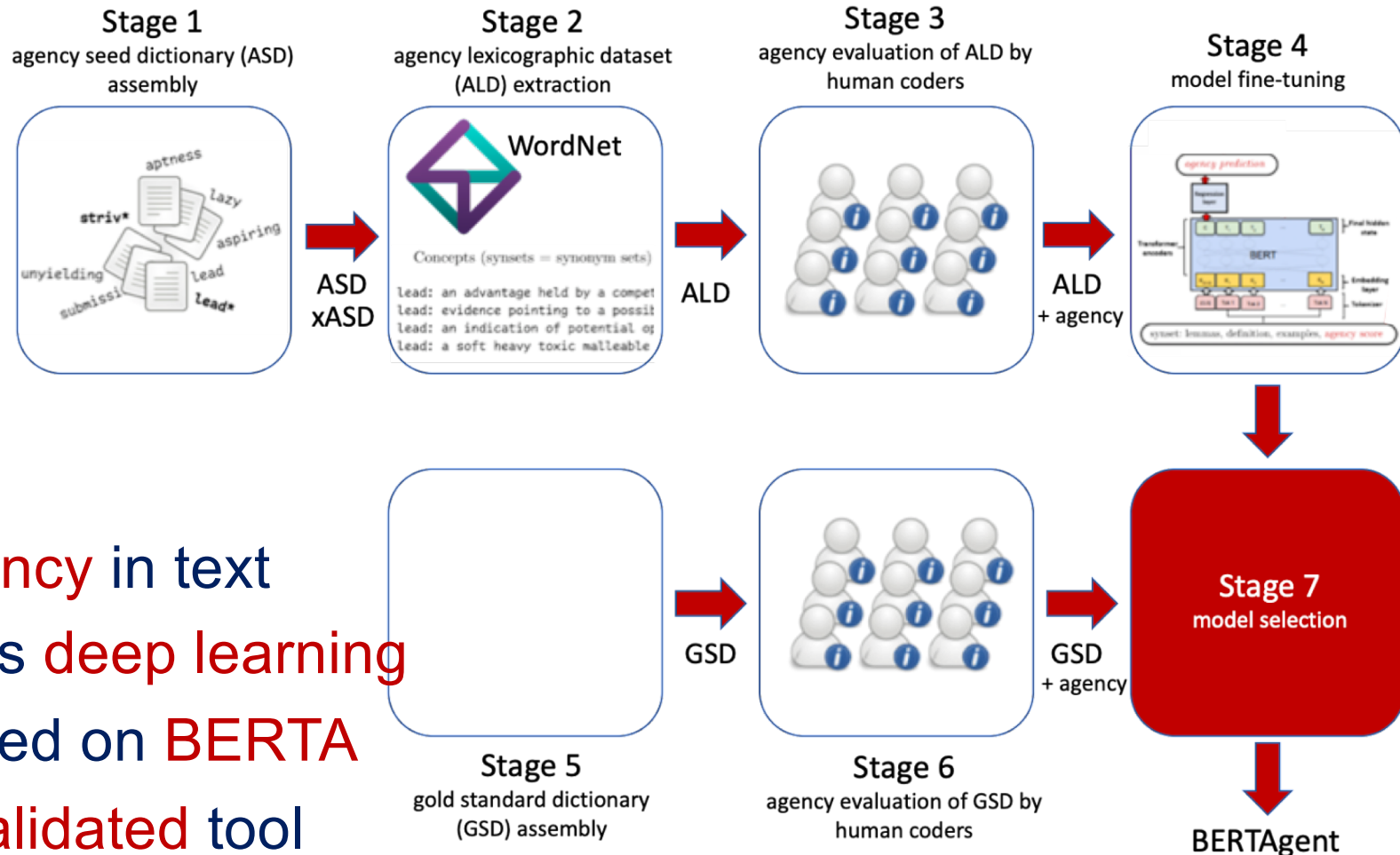
focus on  
past, present  
or future





WC	Analytic	Clout	Authentic	Tone	WPS	Sixltr	Dic	function	pronoun
ppron	<b>i</b>	<b>we</b>	<b>you</b>	<b>shehe</b>	<b>they</b>	ipron	article	prep	auxverb
adverb	conj	negate	verb	adj	compare	interrog	number	quant	affect
<b>posemo</b>	<b>negemo</b>	<b>anx</b>	<b>anger</b>	<b>sad</b>	social	family	friend	female	male
insight	cause	discrep	tentat	certain	differ	percept	see	hear	feel
bio	<b>body</b>	<b>health</b>	<b>sexual</b>	ingest	drives	<b>affiliation</b>	<b>achieve</b>	<b>power</b>	<b>reward</b>
risk	<b>focus past</b>	<b>focus present</b>	<b>focus future</b>	relativ	motion	space	time	work	leisure
home	money	<b>relig</b>	death	<b>informal</b>	<b>swear</b>	netspeak	assent	nonflu	filler
AllPunc	Period	Comma	Colon	SemiC	QMark	Exclam	Dash	Quote	Apostro
Parenth	cogproc								

Choose the ones of **interest** to your project!



- ❑ agency in text
- ❑ uses deep learning
- ❑ based on BERT
- ❑ a validated tool
- ❑ available on Python

BERTAgent  
<https://pypi.org/project/bertagent/>



# Validation of BERTAgent

deep learning wins versus DWC = dictionary word count

Variable	<i>M</i>	<i>SD</i>	1	2	3	4	5	6	7	8
1. HumEval	0.12	1.54								
2. PietA	0.05	0.05	.17** [.06, .28]		-1.25	0.28	0.05	5.35**	-1.78	-10.95**
3. PietB	0.02	0.03	.25** [.14, .35]	.40** [.30, .49]		1.27	1.16	6.58**	-0.70	-10.00**
4. PietC	0.05	0.05	.17** [.06, .28]	.99** [.99, 1.00]	.40** [.30, .49]		0.03	5.34**	-1.80	-10.93**
5. NicoPos	0.03	0.04	.17** [.05, .27]	.18** [.07, .29]	.23** [.12, .34]	.17** [.06, .28]		5.49**	-3.81**	-11.08**
6. NicoNeg	0.01	0.03	-.28** [-.38, -.17]	-.10 [-.21, .01]	-.01 [-.12, .11]	-.10 [-.21, .02]	-.03 [-.14, .09]		-5.73**	-13.40**
7. NicoCom	0.02	0.05	.30** [.19, .40]	.20** [.09, .31]	.19** [.08, .30]	.19** [.08, .30]	.82** [.78, .85]	-.60** [-.67, -.52]		-10.38**
8. BATot	0.09	0.35	.78** [.73, .82]	.21** [.10, .31]	.24** [.13, .34]	.20** [.09, .31]	.22** [.11, .33]	-.42** [-.51, -.33]	.42** [.33, .51]	

Human  
evaluation

BERTAgent

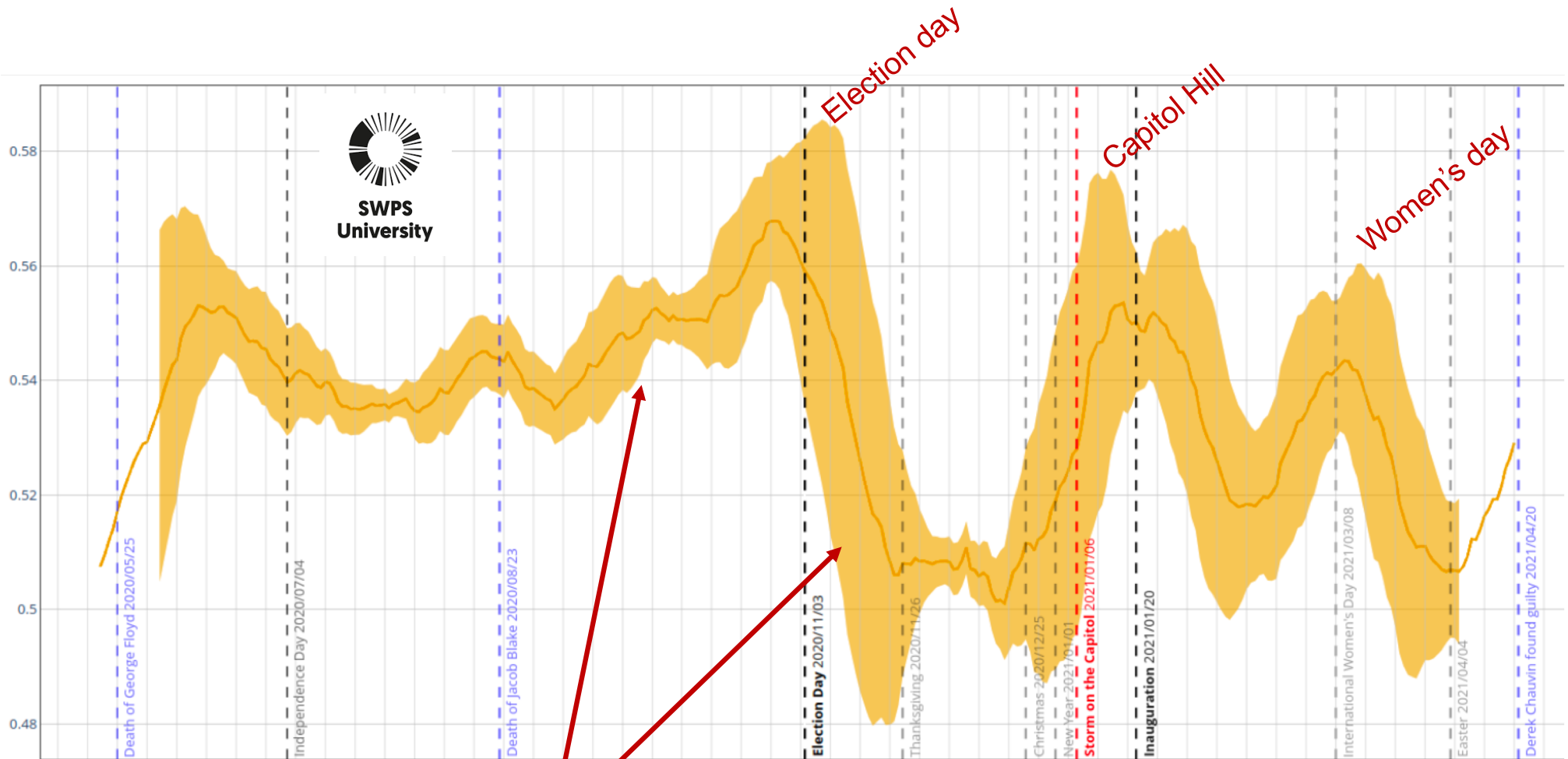
best correlation with  
Human evaluation

Z-statistics:  
correlation is statistically  
more relevant than DWC



# Agency in US elections

Twitter, 2020-2021  
by Jan Nikadon @ swps



agency raises before elections than drops





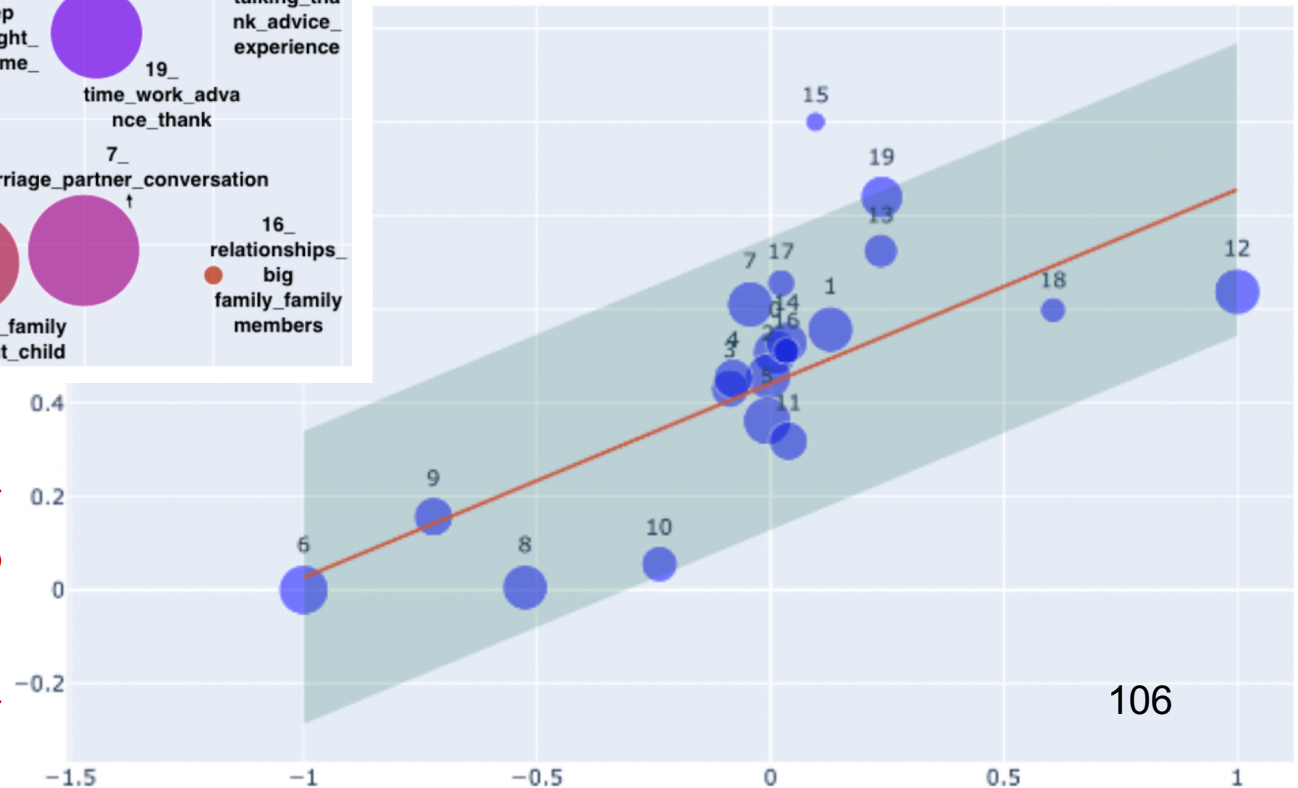
# Agency in postpartum depression

Reddit Posts 2021  
by Selen Arslan @ unipd/swps



composite emotion score  
(LIWC, EmoPos - EmoNeg)

semantic agency  
(BERTAgent)





```
[1] !pip install bertagent
```

install, import,  
set instance

```
▶ from bertagent import BERTAgent  
ba0 = BERTAgent()
```

```
▶ # provide example sentences
```

input sentences must  
be superficially cleaned

```
sents = ["hardly wo  
"hard work  
"striving  
"strugglin  
"strugglin  
"unable to  
"this car  
"this car  
"this poli  
]
```

'hardly working individual' : -0.57	BERTAgent output
'hard working individual' : 0.44	
'striving to achieve my goals' : 0.73	
'struggling to achieve my goals' : -0.67	
'struggling to survive' : -0.52	
'unable to survive' : -0.57	
'this car runs on gasoline with lead' : -0.03	
'this car runs on gasoline and it will lead us' : 0.09	
'this politician runs for office and he will lead us' : 0.58	

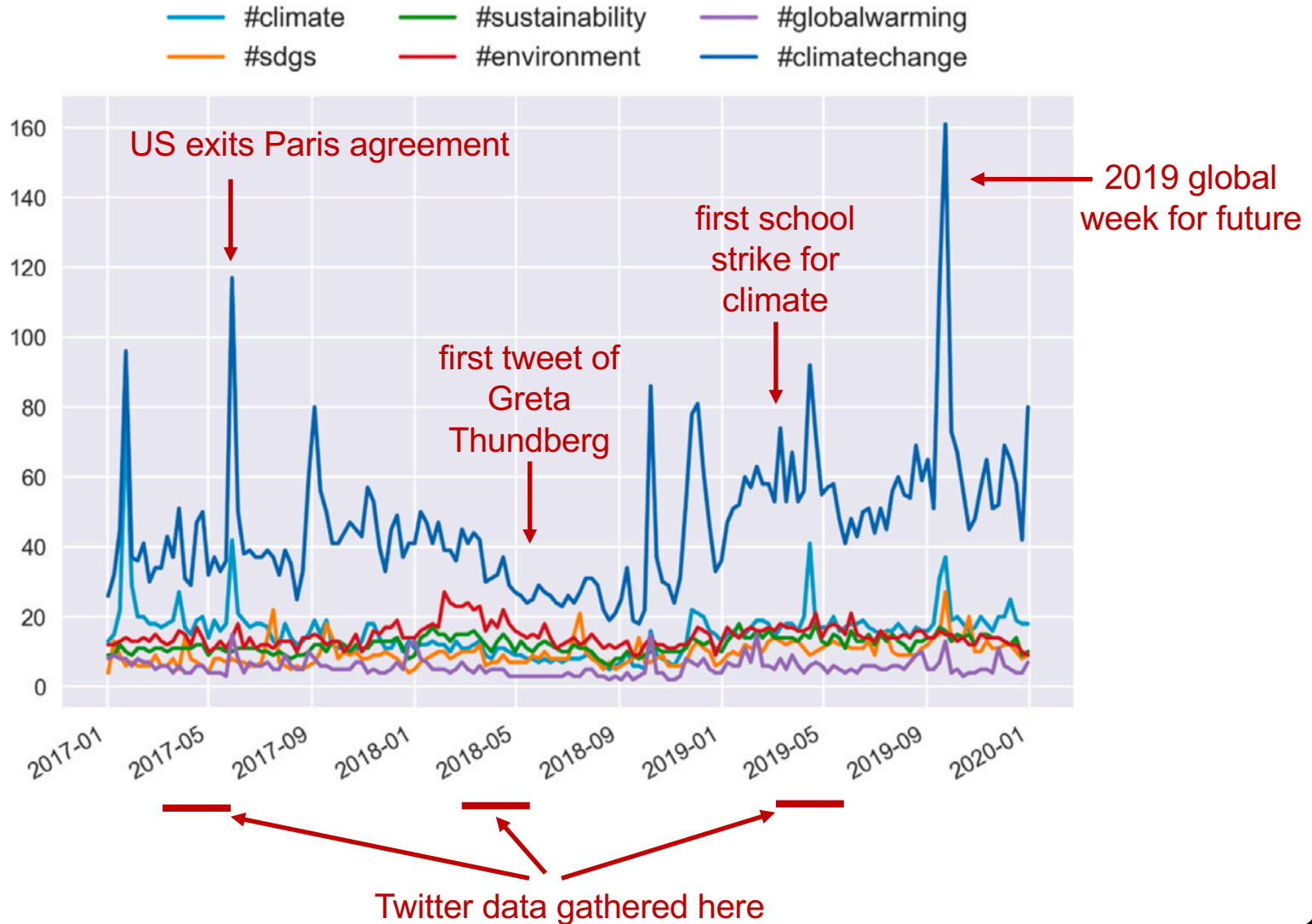
```
# assign agency  
vals = ba0.predict(sents)  
# print results  
for item in zip(sents, vals):  
    print(f" {item[0]!r} : {item[1]:.2f}")
```

run BERTAgent



# Using sentiment analysis

an overview on how it can be useful in your projects





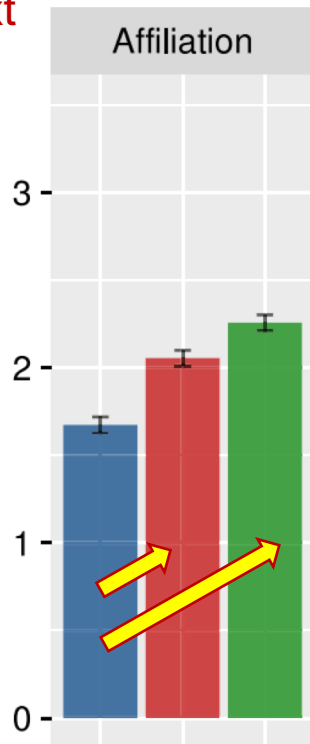


# Socio-psychological linguistic markers

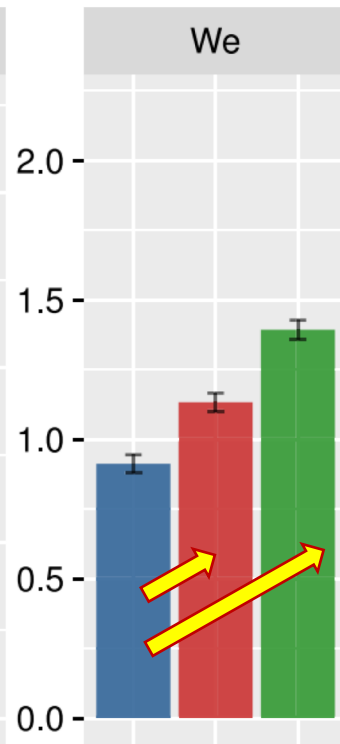
a view on the entire tweets corpus

■ 2017 ■ 2018 ■ 2019

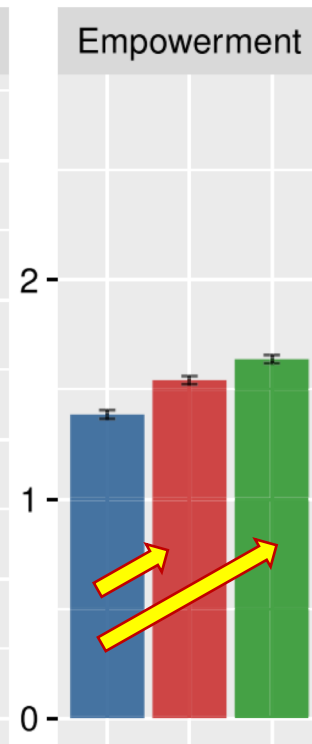
ingroup community  
orientation within  
the text



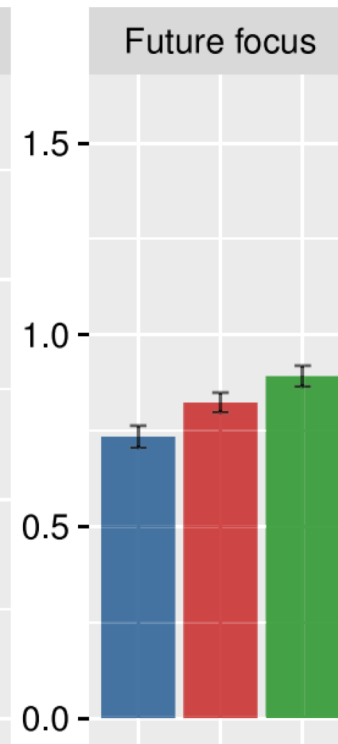
salience of  
group  
membership,  
sense of  
belonging



a person's  
striving to be  
independent to  
assert, protect  
and expand  
one's self



orientation of tweets to the  
past or future



only a few  
statistically  
relevant  
changes



## ☰ Student's *t*-test

[Article](#) [Talk](#)

From Wikipedia, the free encyclopedia

A ***t*-test** is a type of statistical analysis used to **compare the averages of two groups** and determine whether the differences between them are more likely to arise from random chance. It is any **statistical hypothesis test** in which the **test statistic** follows a **Student's *t*-distribution** under the **null hypothesis**. It is most commonly applied when the test statistic would follow a **normal distribution** if the value of a **scaling term** in the test statistic were known (typically, the scaling term is unknown and is therefore a **nuisance parameter**). When the scaling term is estimated based on the **data**, the test statistic—under certain conditions—follows a Student's *t* distribution. **The *t*-test's most common application is to test whether the means of two populations are different.**



## Assumption:

the **average** values of the two populations being compared should follow a **normal distribution** (e.g., with many samples)

## Hypothesis (to be tested):

H0 - the two sets have **equal** average value,  $m_x = m_y$

H1 - the two sets have different average value,  $m_x \neq m_y$

Test statistic:  $t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s_x^2}{N_x} + \frac{s_y^2}{N_y}}} \sim t_\nu$

sample averages

unbiased estimator of variance

samples sizes

Student's t distribution with  $\nu$  degrees of freedom (an approximation under H0)

estimate of  $\nu$  (degrees of freedom) used for calculations

$$\hat{\nu} = \frac{\left(\frac{s_x^2}{N_x} + \frac{s_y^2}{N_y}\right)^2}{\frac{(s_x^2/N_x)^2}{N_x - 1} + \frac{(s_y^2/N_y)^2}{N_y - 1}}$$



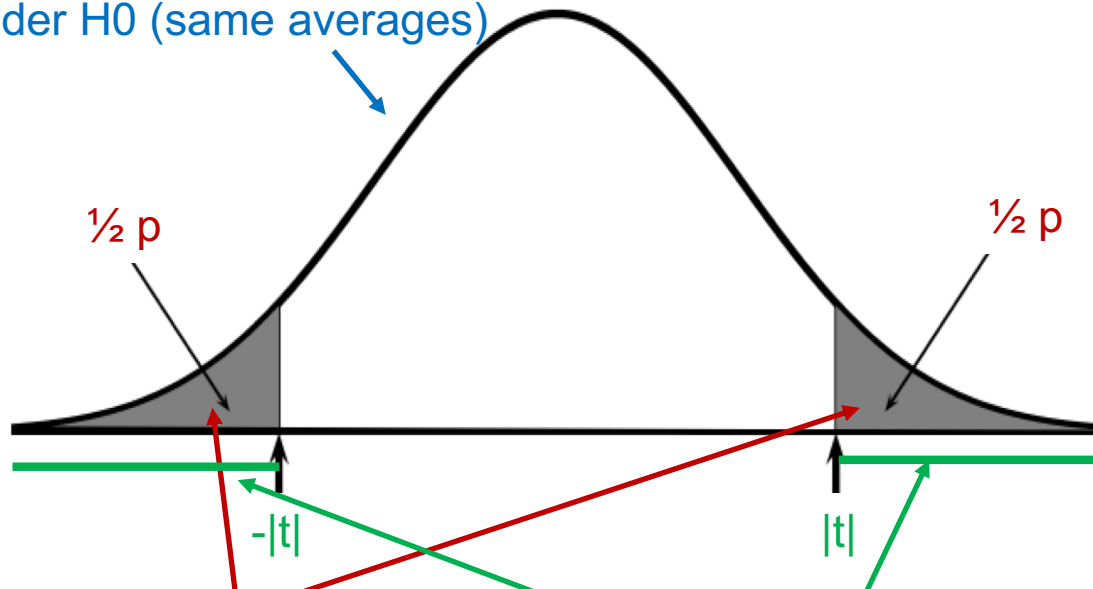
# Student/Welch t-test

p value = probability of error under  $H_0$

we declare different average values  $H_1$

p-value	evidence
$< .01$	very strong evidence against $H_0$
$.01 - .05$	strong evidence against $H_0$
$.05 - .10$	weak evidence against $H_0$
$> .1$	little or no evidence against $H_0$

Student-t PDF with  $\nu$  degrees of freedom under  $H_0$  (same averages)



probability  $p$  that, by confirming  $H_1$ , we are making errors on  $H_0$

(conservative) region choice to declare  $H_1$  from the  $t$  value

the  $p$  value informs on whether a difference exists (statistically)

it is a false positive rate



# Cohen's d-value

evaluating effect sizes

$$\text{Cohen's } d = \frac{\bar{x} - \bar{y}}{\sqrt{a s_x^2 + (1-a) s_y^2}}$$

$$a = \frac{N_x - 1}{N_x + N_y - 2}$$

Relative size	Effect size
	0.0
Small	0.2
Medium	0.5
Large	0.8
	1.4

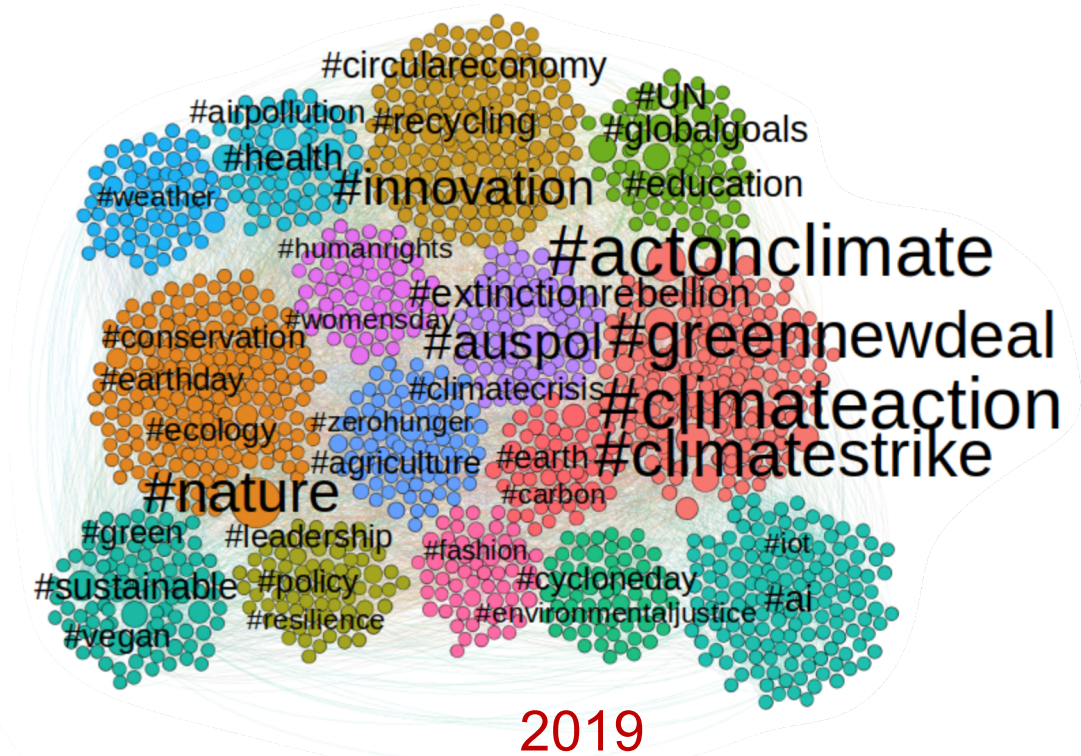
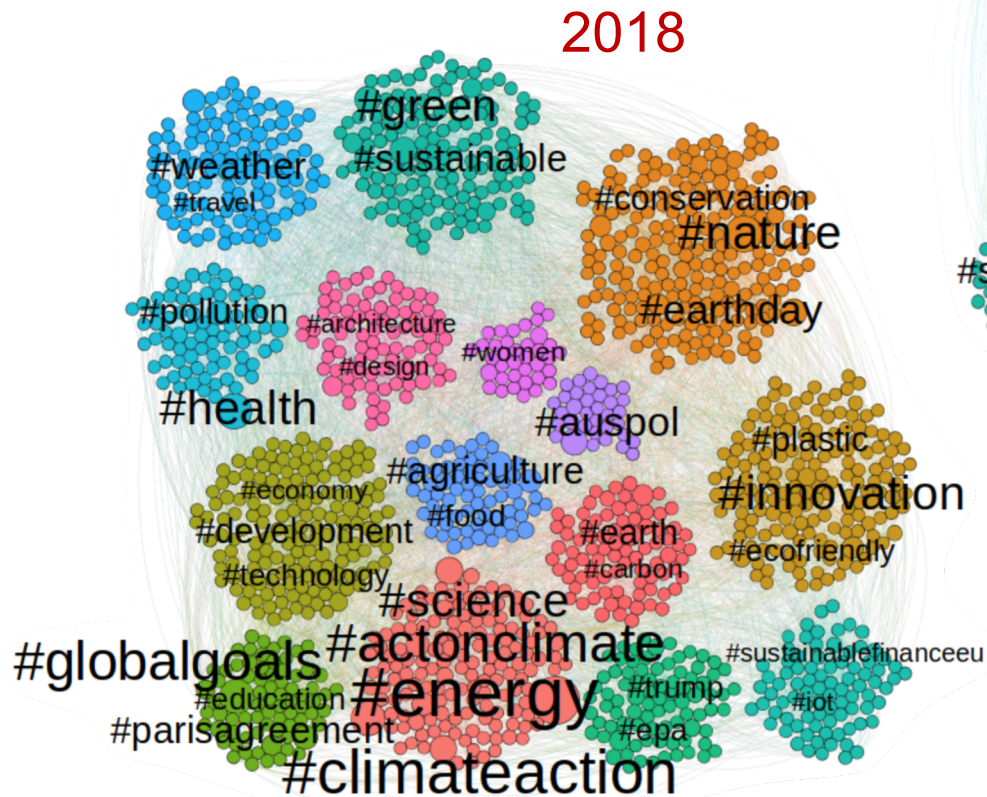
we confirm H1



the *d* value informs on the size of the effect

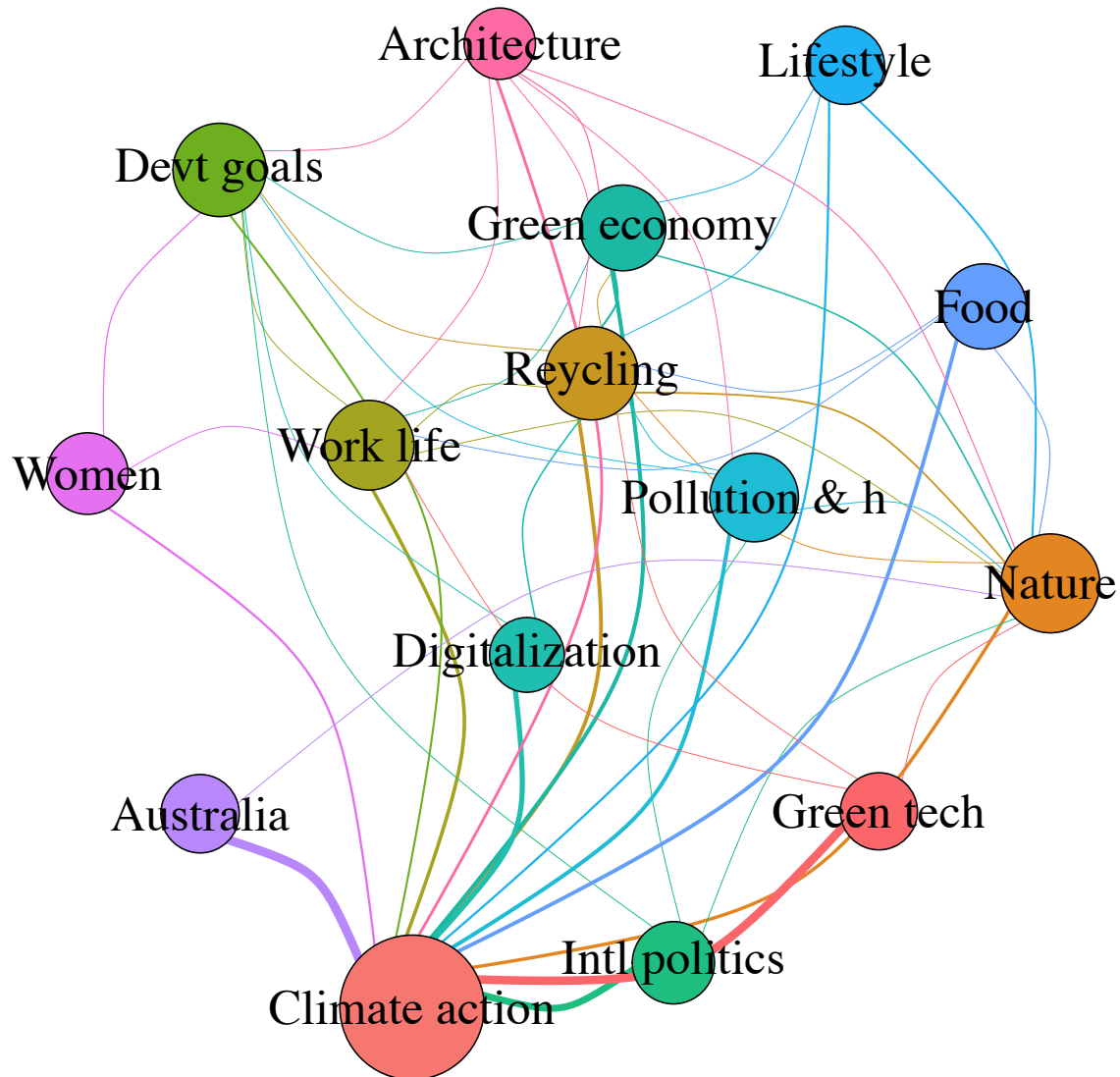


# Topics in #climateaction on Twitter in 2017, 2018, 2019





# Topics interdependencies



projecting the  
adjacency matrix on  
topics

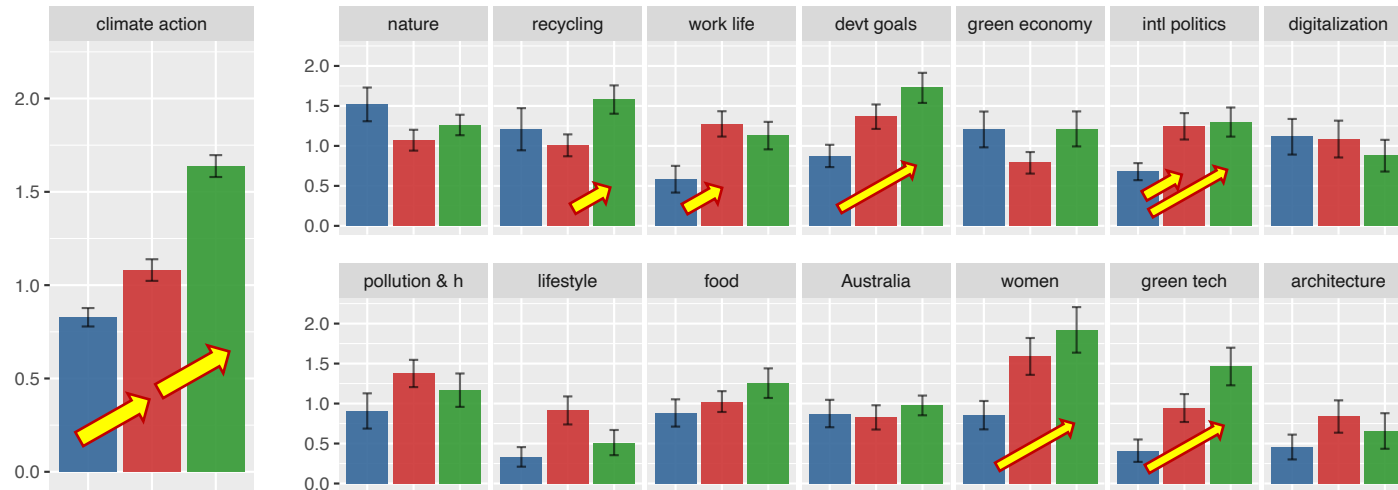
$P_{11}$	$P_{12}$	$P_{13}$
$P_{21}$	$P_{22}$	$P_{23}$
$P_{31}$	$P_{32}$	$P_{33}$



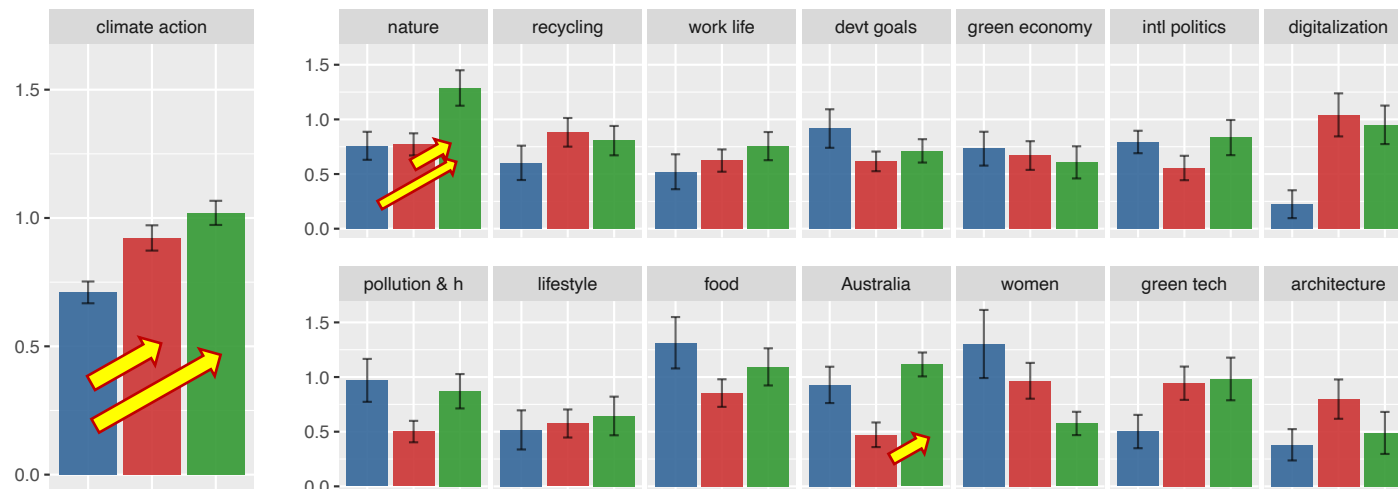
# Socio-psychological linguistic markers a view inside topics

2017 2018 2019

(b) We



(d) Future focus



relevant  
statistically  
changes of  
we-future  
only in the  
climate  
action  
community



# Wrap-up

on topic detection



- ❑ What available tools should be used
  - Louvain & BERTopic
  - compare their performance through NMI, modularity, etc.
  - LIWC & BERTAgent
  - to enrich your analysis under a socio-psychological lens
  
- ❑ What available tools should **NOT** be used
  - InfoMap, NMF & LDA
  - they show poor performance
  
- ❑ What would be nice to see implemented
  - soft Louvain made fast
  - performance of BigCLAM and SMBs
  - NFTM VAE and its performance