

Measuring online sentiment with BERT

Pieter Delobelle

Nov 24, 2023

Outline

- **Language models:** BERT, GPT, ...
 - What?
 - How?
 - Where?
- **Measuring sentiment:** illustration with COVID measures
- **Topic modelling**
- **Practical session**

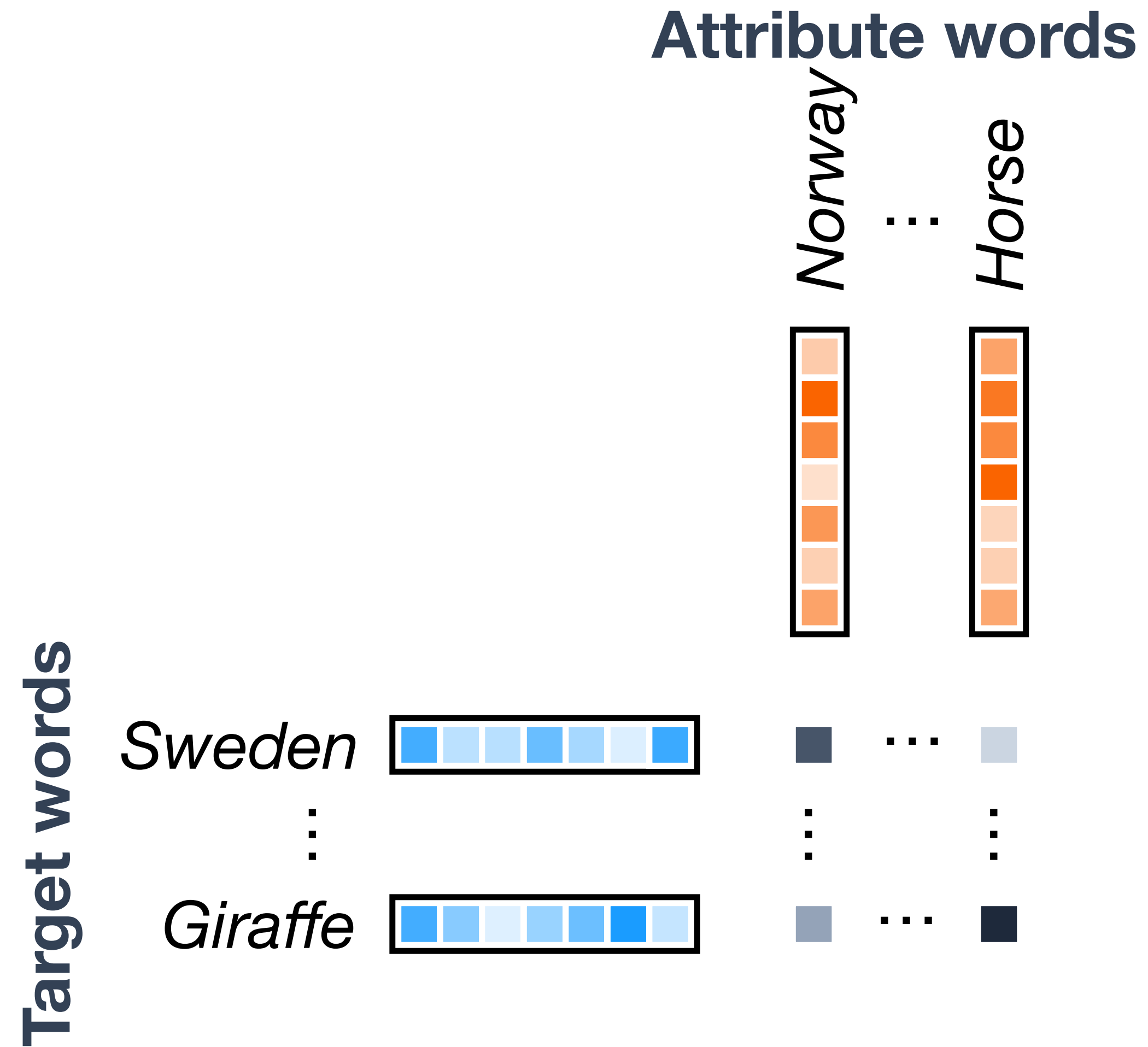
Language models

Word embeddings

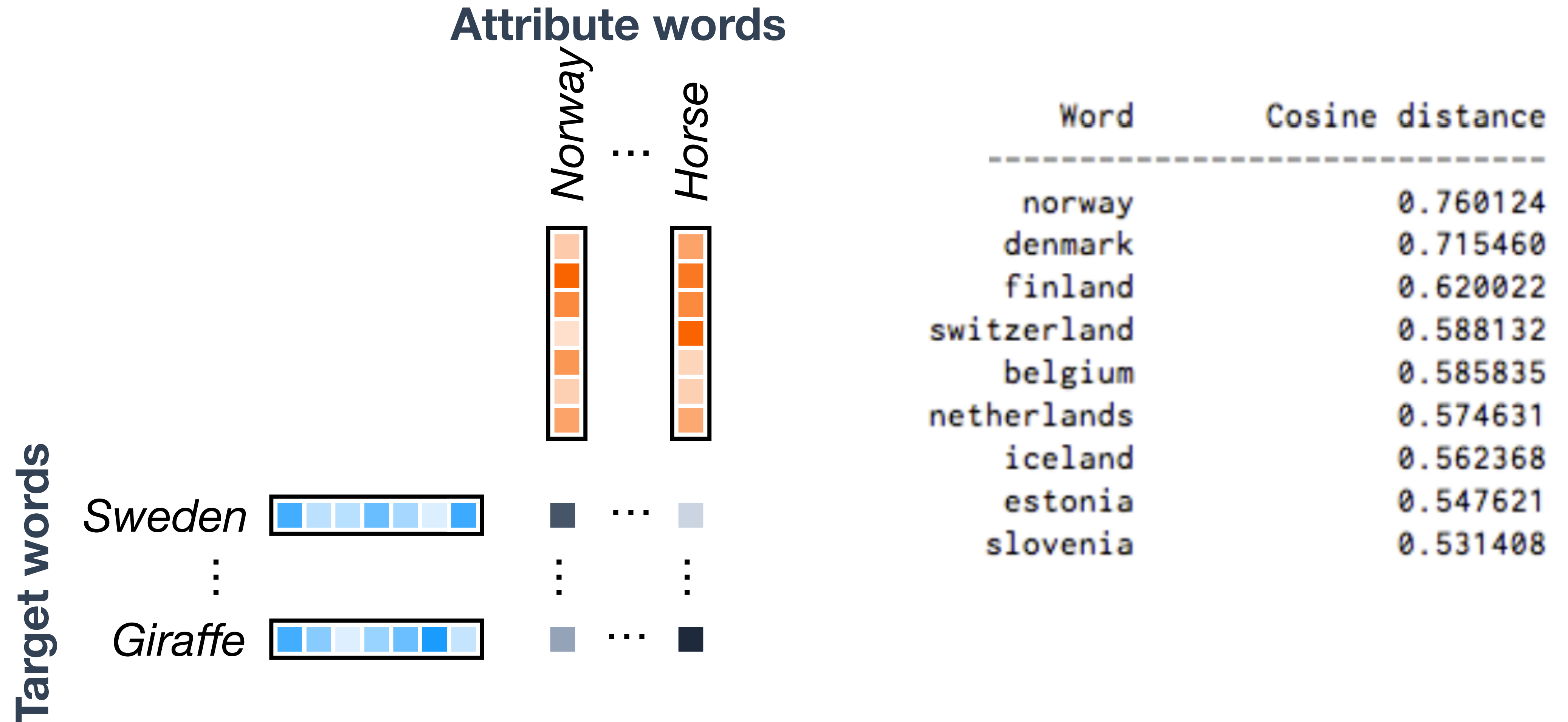


stick

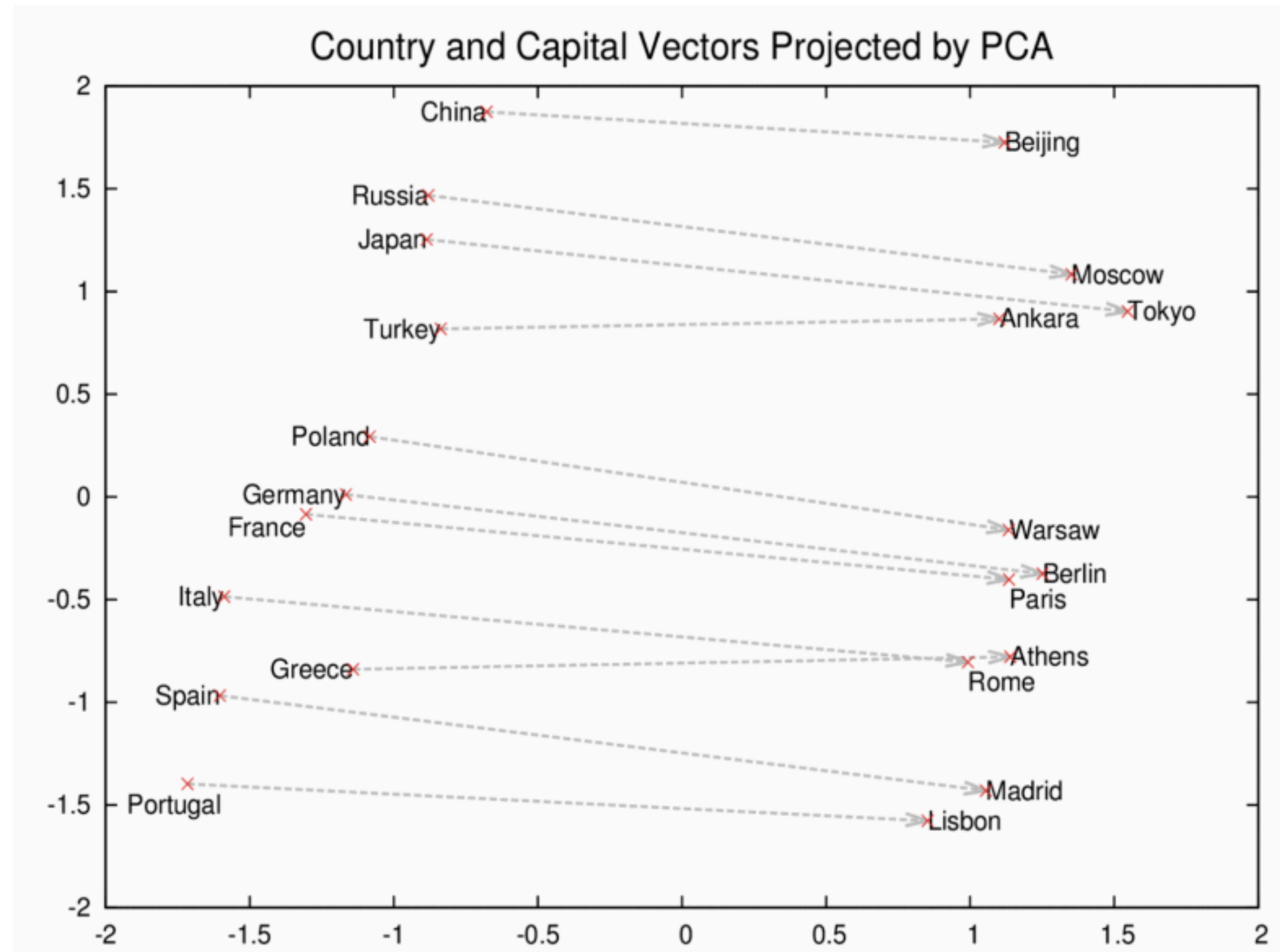
Word embeddings



Word embeddings



Embeddings have meaningful principal components



Word embeddings



stick

Word embeddings don't understand polysemy



Bank



Bank

Word embeddings don't understand polysemy



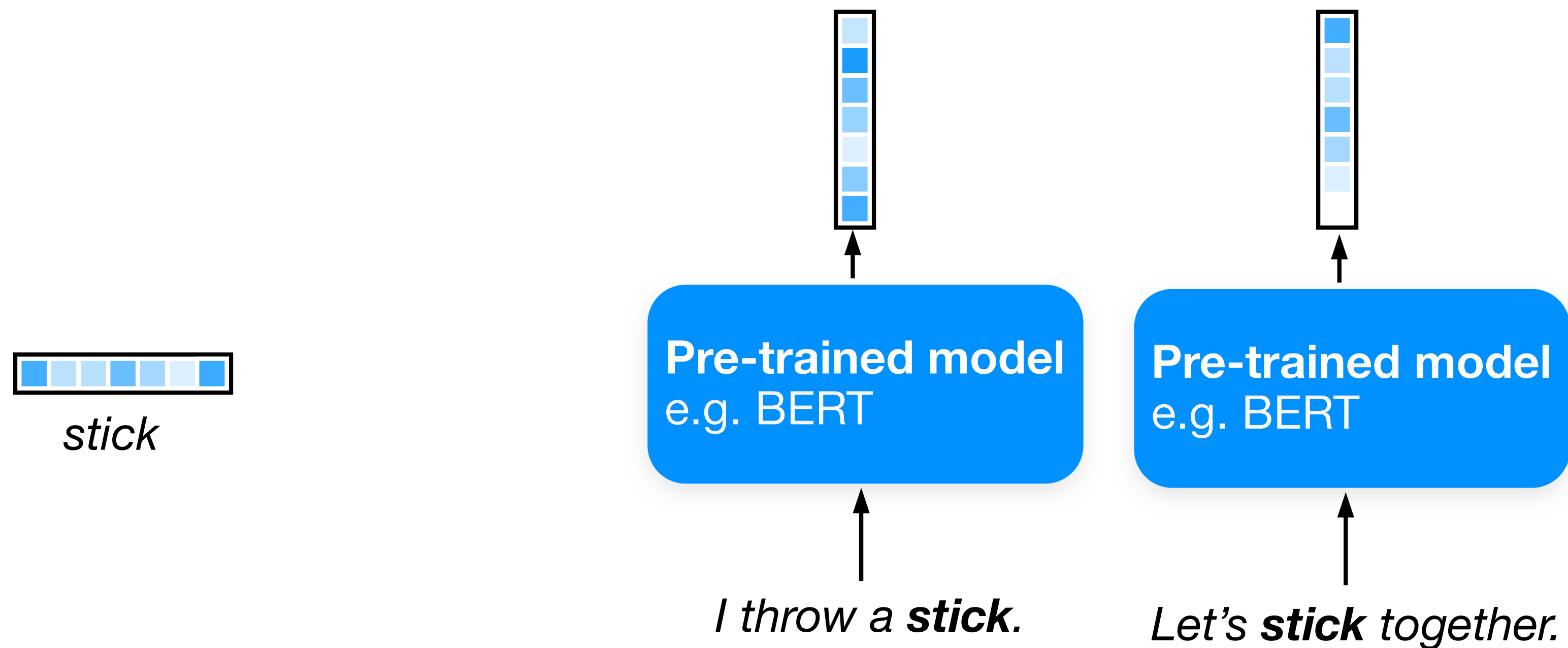
Bank



Bank

→ How to incorporate context?

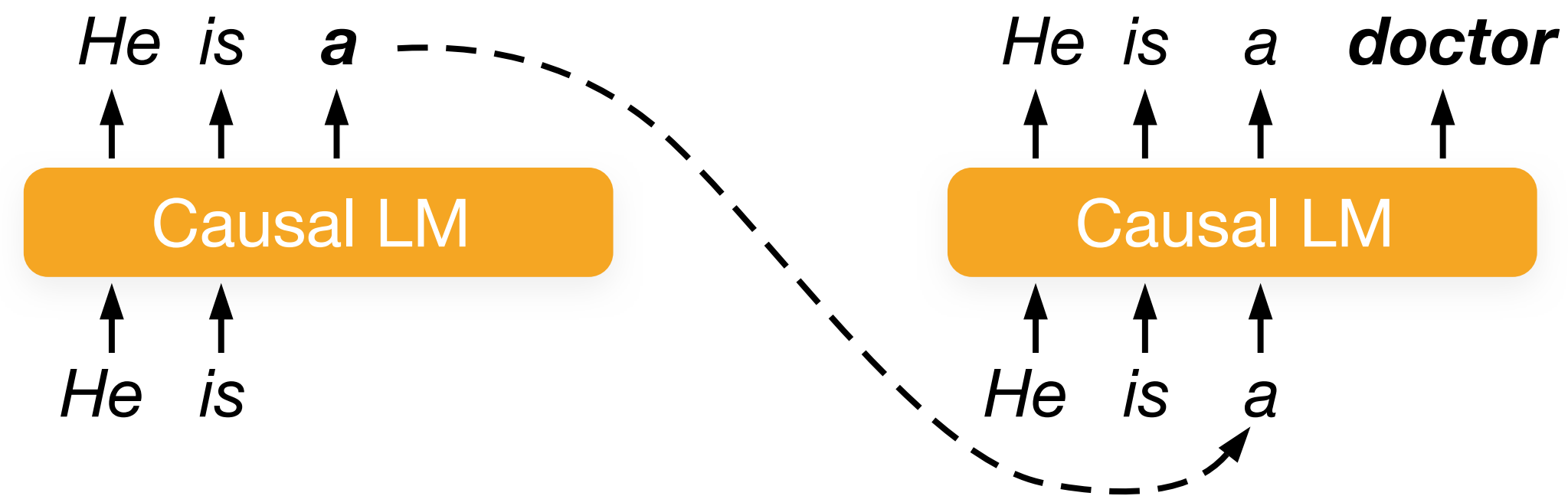
Language models address polysemy



Language modeling



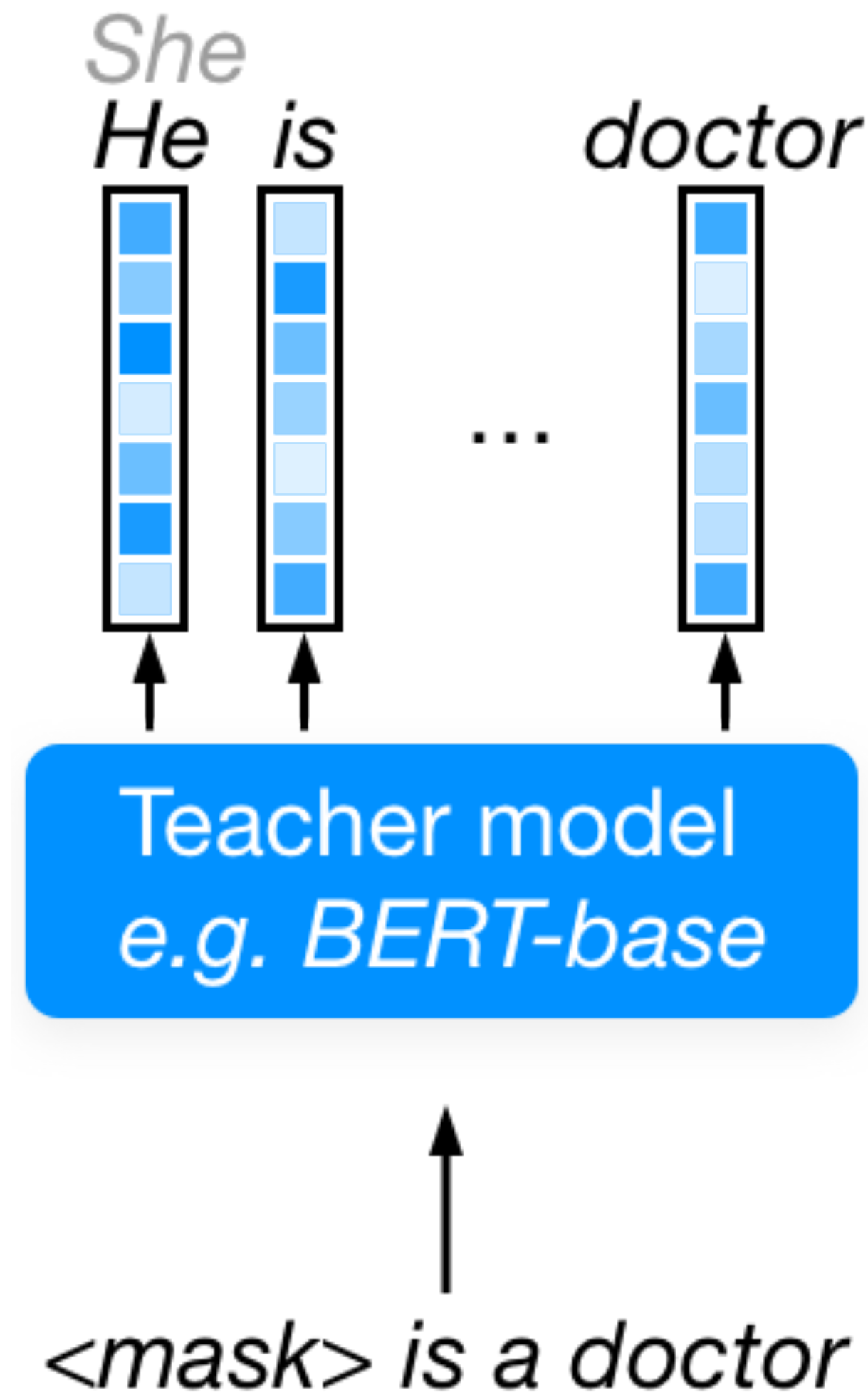
1. Causal language modeling (CLM)



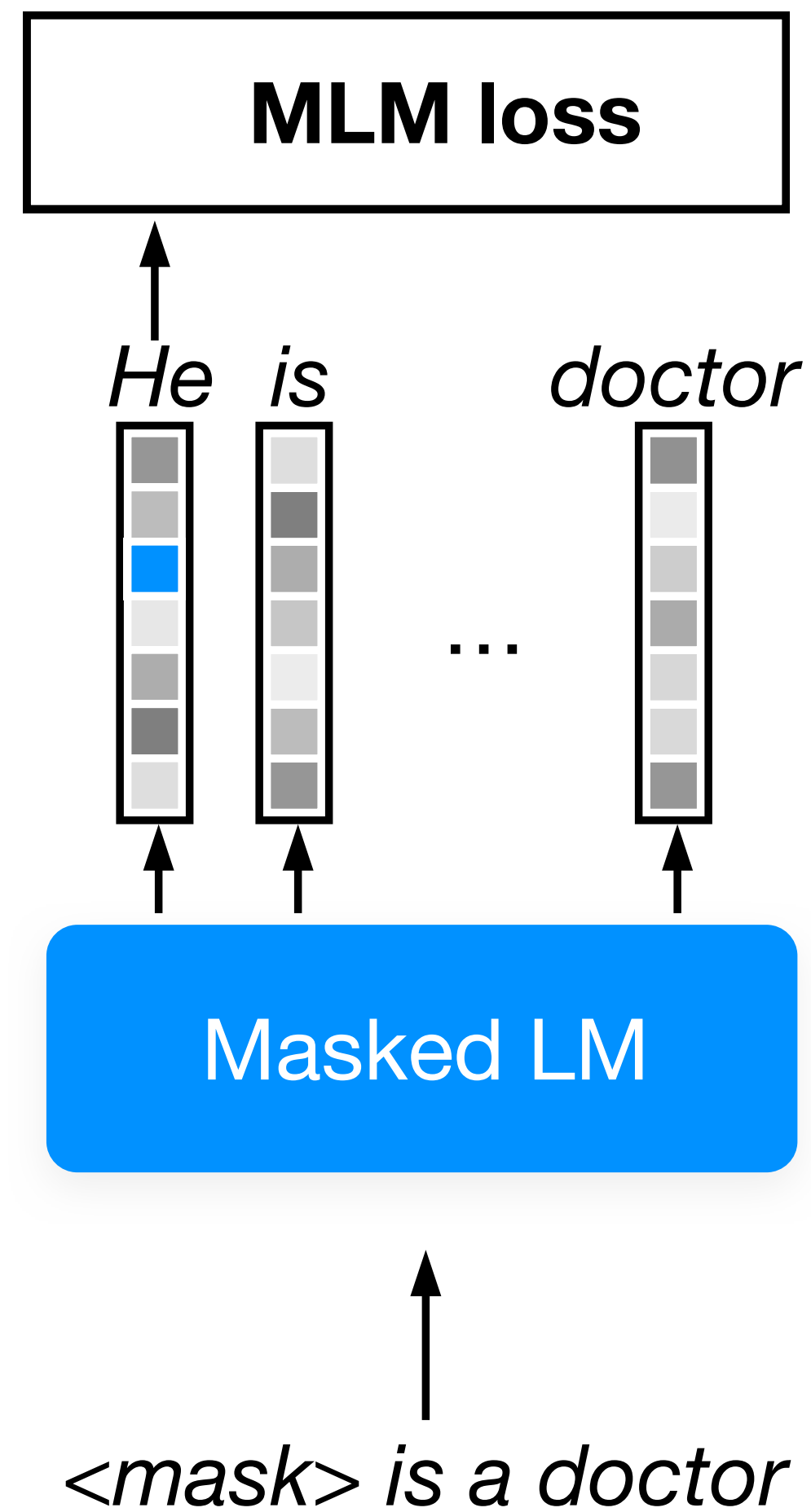
2. Masked language modeling (MLM)



MLMs learn a probability for each word

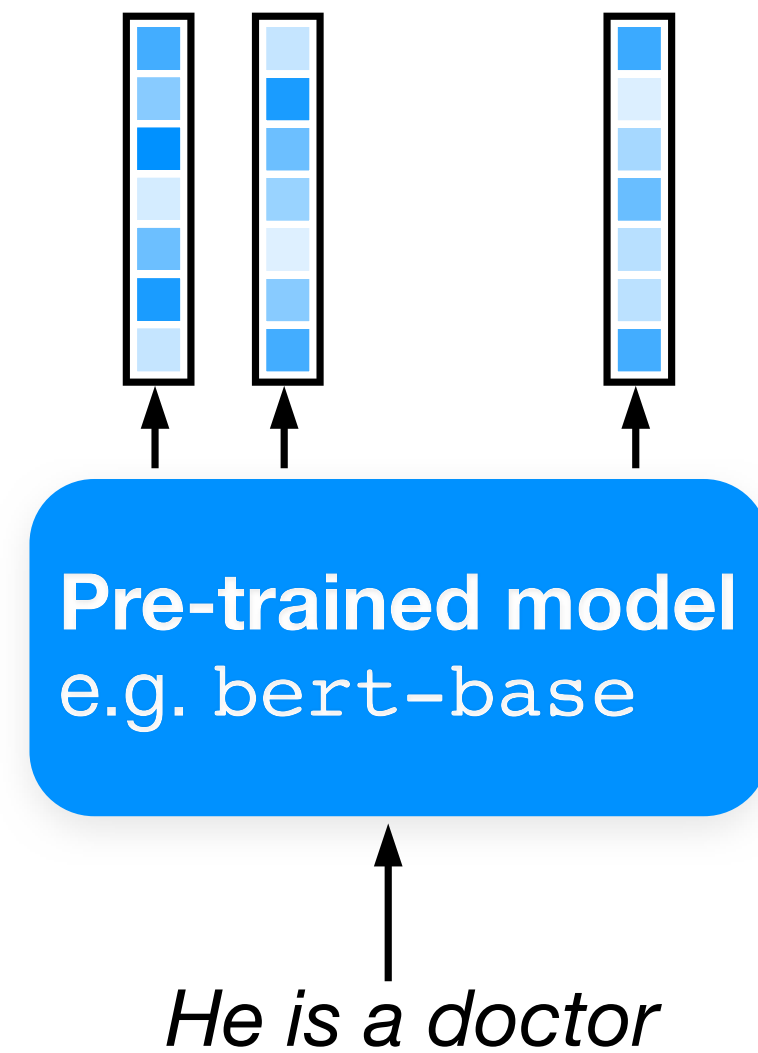


MLMs learn a probability for each word



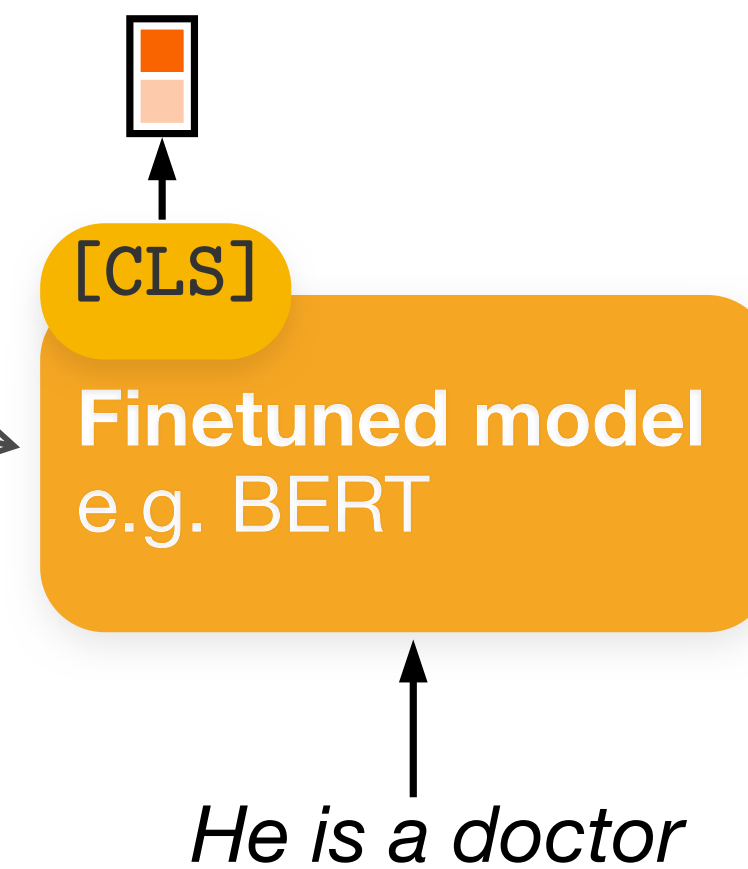
MLMs are trained twice

1. Pretraining step
e.g. OSCAR, Wikipedia, ...













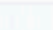
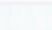
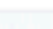
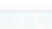


Transfer
learning

2. Finetuning step
e.g. sentiment analysis,
named entity recognition



Where? On the HuggingFace hub

Models 408,066 new Full-text search Sort: Trending

 stabilityai/stable-video-diffusion-img2vid-xt Updated 2 days ago • ❤️ 628	 openai/whisper-large-v3 Automatic Speech Recognition • Updated 2 days ago • ↓ 99k • ❤️ 882
 microsoft/Orca-2-13b Text Generation • Updated 1 day ago • ↓ 2.77k • ❤️ 289	 stabilityai/stable-video-diffusion-img2vid Updated about 16 hours ago • ❤️ 181
 coqui/XTTS-v2 Text-to-Speech • Updated 6 days ago • ↓ 66.4k • ❤️ 292	 Intel/neural-chat-7b-v3-1 Text Generation • Updated about 7 hours ago • ↓ 4.84k • ❤️ 152
 openchat/openchat_3.5 Text Generation • Updated 1 day ago • ↓ 24.9k • ❤️ 674	 latent-consistency/lcm-lora-sdxl Text-to-Image • Updated 8 days ago • ↓ 51.5k • ❤️ 360
 google/switch-c-2048 Text2Text Generation • Updated 4 days ago • ↓ 595 • ❤️ 159	 01-ai/Yi-34B Text Generation • Updated about 5 hours ago • ↓ 58.7k • ❤️ 993
 mistralai/Mistral-7B-v0.1 Text Generation • Updated Oct 12 • ↓ 436k • ❤️ 1.93k	 microsoft/Orca-2-7b Text Generation • Updated 1 day ago • ↓ 2.22k • ❤️ 97
 HuggingFaceH4/zephyr-7b-beta Text Generation • Updated 1 day ago • ↓ 169k • ❤️ 883	 meta-llama/Llama-2-7b-chat-hf Text Generation • Updated 11 days ago • ↓ 950k • ❤️ 1.89k
 stabilityai/stable-diffusion-xl-base-1.0 Text-to-Image • Updated 25 days ago • ↓ 9.97M • ❤️ 3.59k	 teknium/OpenHermes-2.5-Mistral-7B Text Generation • Updated 20 days ago • ↓ 16.7k • ❤️ 200

Measuring sentiment on Twitter

Context

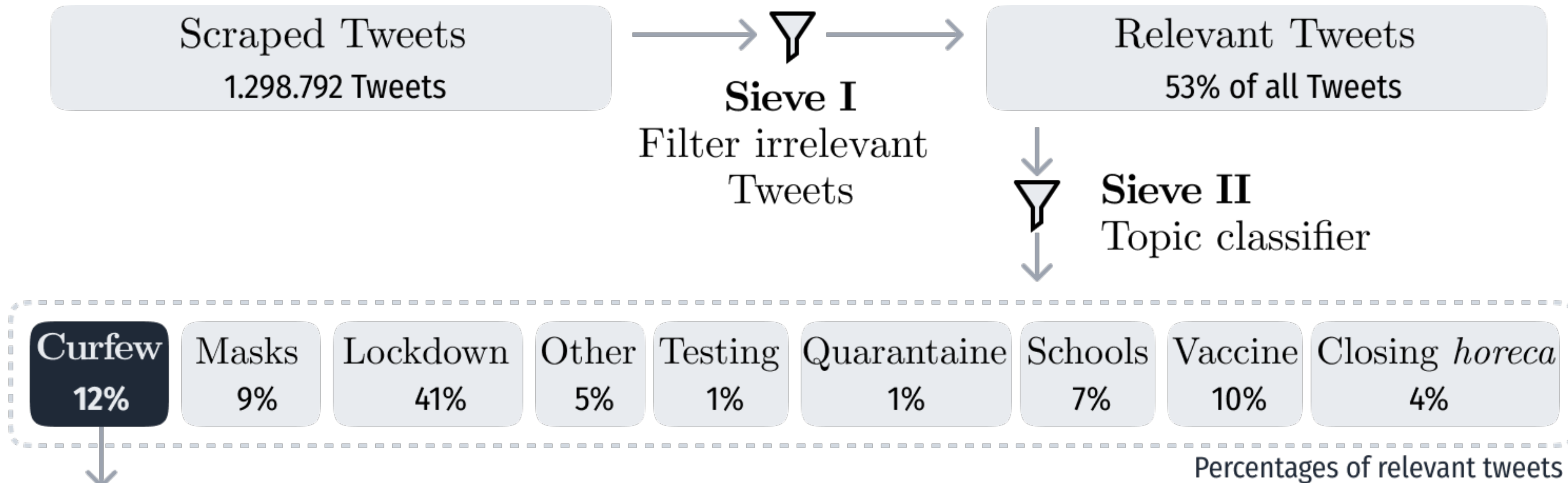
- Evaluation of 1.3M collected Tweets on COVID-19 measures
- Focused on discussion of COVID-19 policy in Belgium
- Additional focus on support for curfews
 - Belgium had multiple curfews (starting at midnight)

Labeling: Doccano

The screenshot shows the Doccano web interface. The browser address bar displays 'doccano - doccano'. The page title is 'Corona Tweets - Curfew run'. The interface includes a search bar with a filter icon and a 'Labels' section containing three active labels: 'topic:curfew (2)', 'measure:too-strict (q)', and 'government:unsupportive (s)'. The main content area displays a tweet in Dutch: 'Waarom spreken wij niet over hoelang we al met die avondklok zitten en blijven wij ons daar als brave Belgen maar gewoon aan houden? Al vier fucking maanden en nergens wordt er gesproken over versoepelingen van die avondklok'. To the right of the tweet is a metadata table with the following data:

Key	Value
language	nl
id_str	136284665404767
created_at	Fri Feb 19 19:28:57 +0000 2021
followers_count	4766

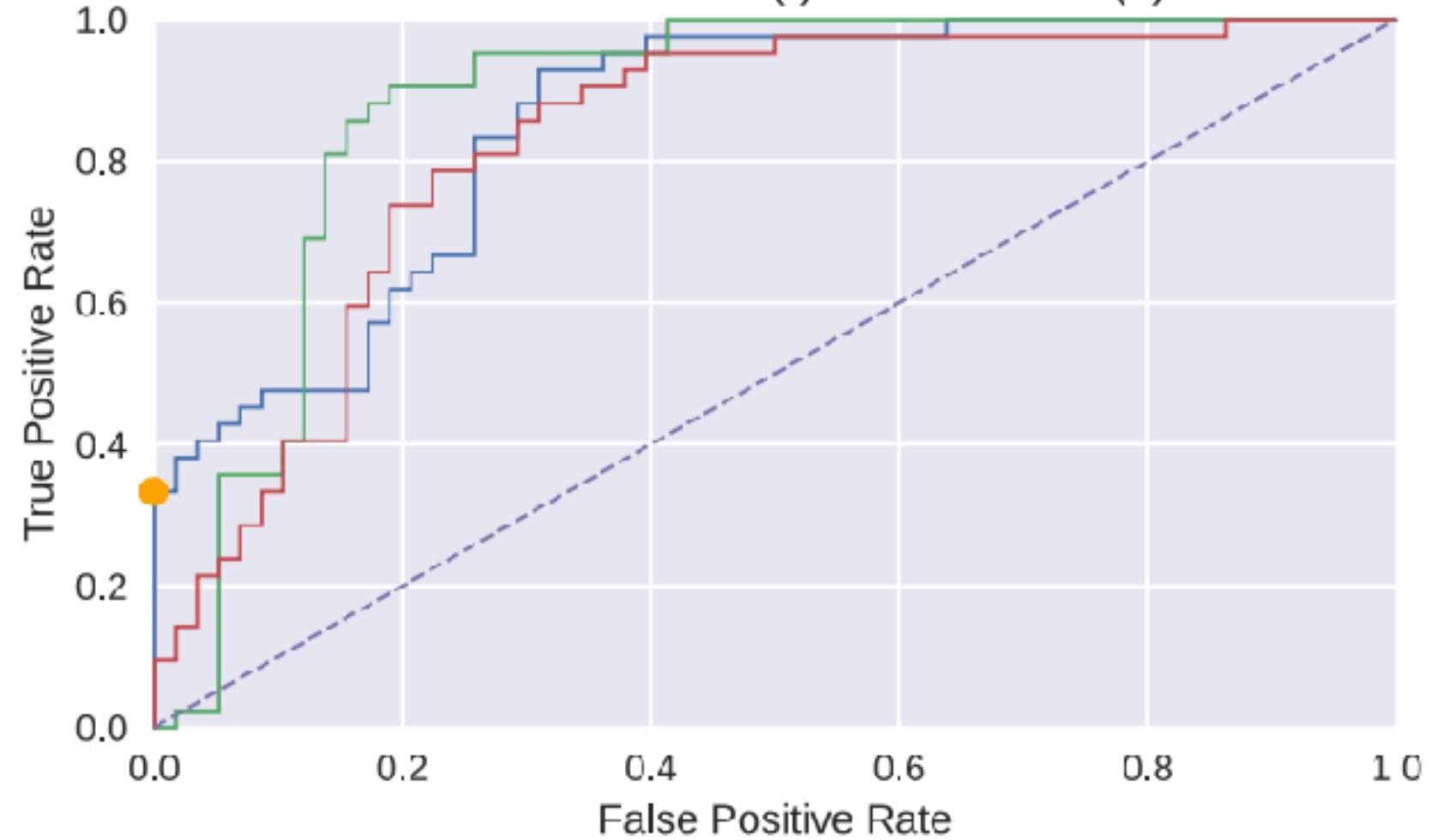
Our pipeline



- **Measure support**

"Is the curfew measure too strict, too loose, or appropriate ('ok')

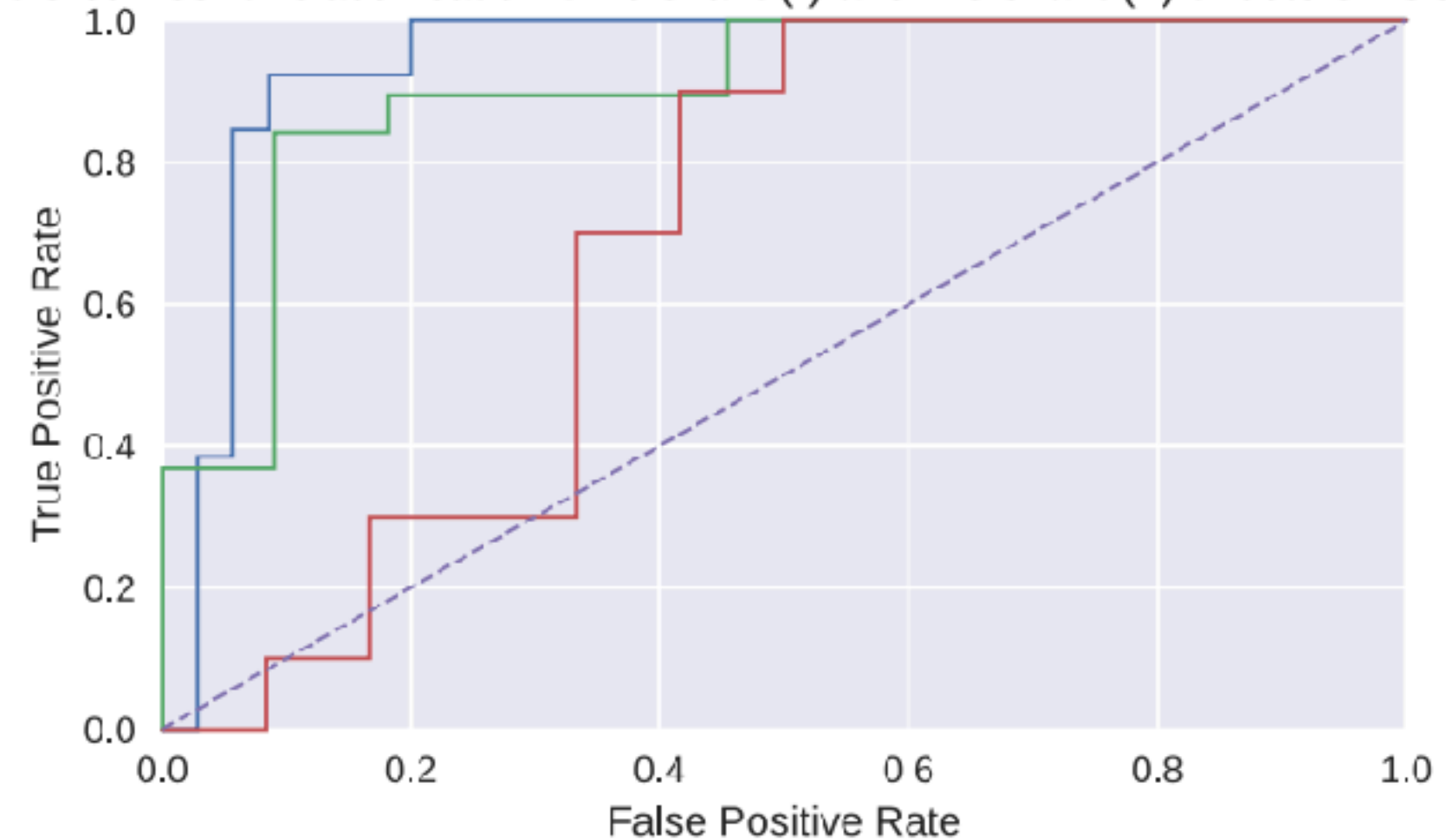
ROC curves for classification of relevant (-) and irrelevant (+) tweets on COVID-19



- 400 tweets mBERT (AUC = 0.85, model version = 2020-12-03)
- 2k tweets mBERT (AUC = 0.88, model version = 2021-01-05)
- 400 + 2k tweets mBERT (AUC = 0.83, model version = 2021-01-05)
- Labeling threshold

(a) ROC curves for different model versions, including the threshold set on the first (400 tweets) model used as Sieve 1.

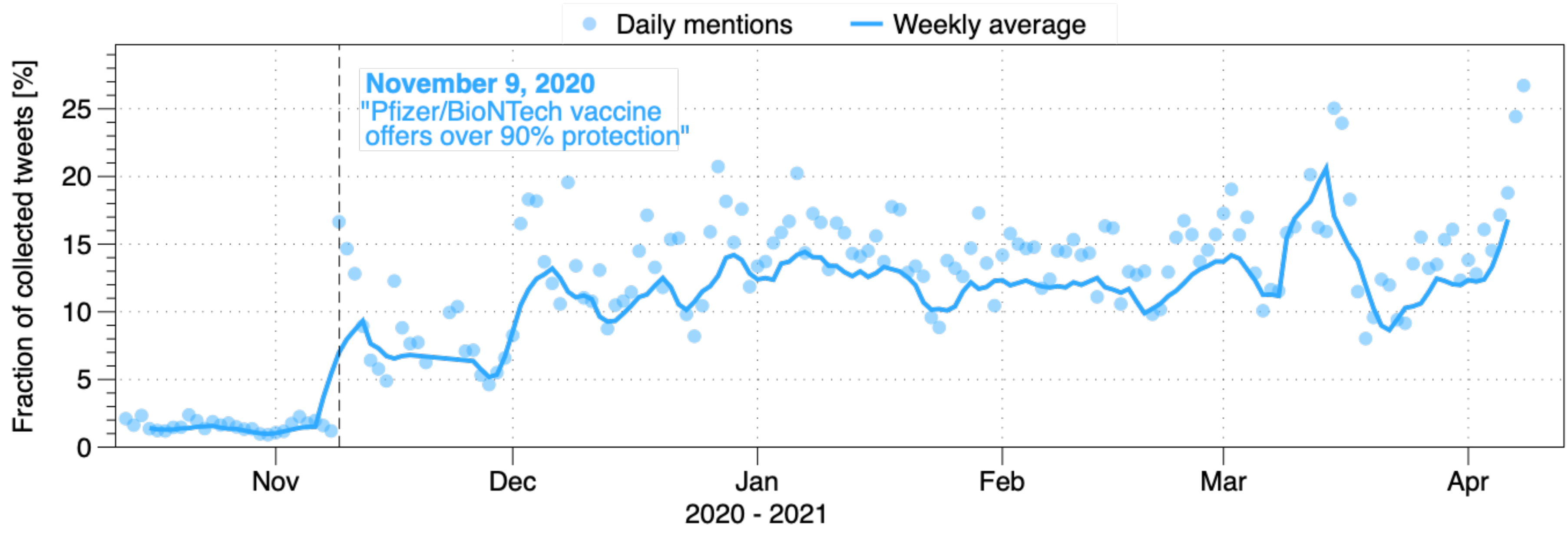
ROC curves for classification of relevant (-) and irrelevant (+) tweets on COVID-19



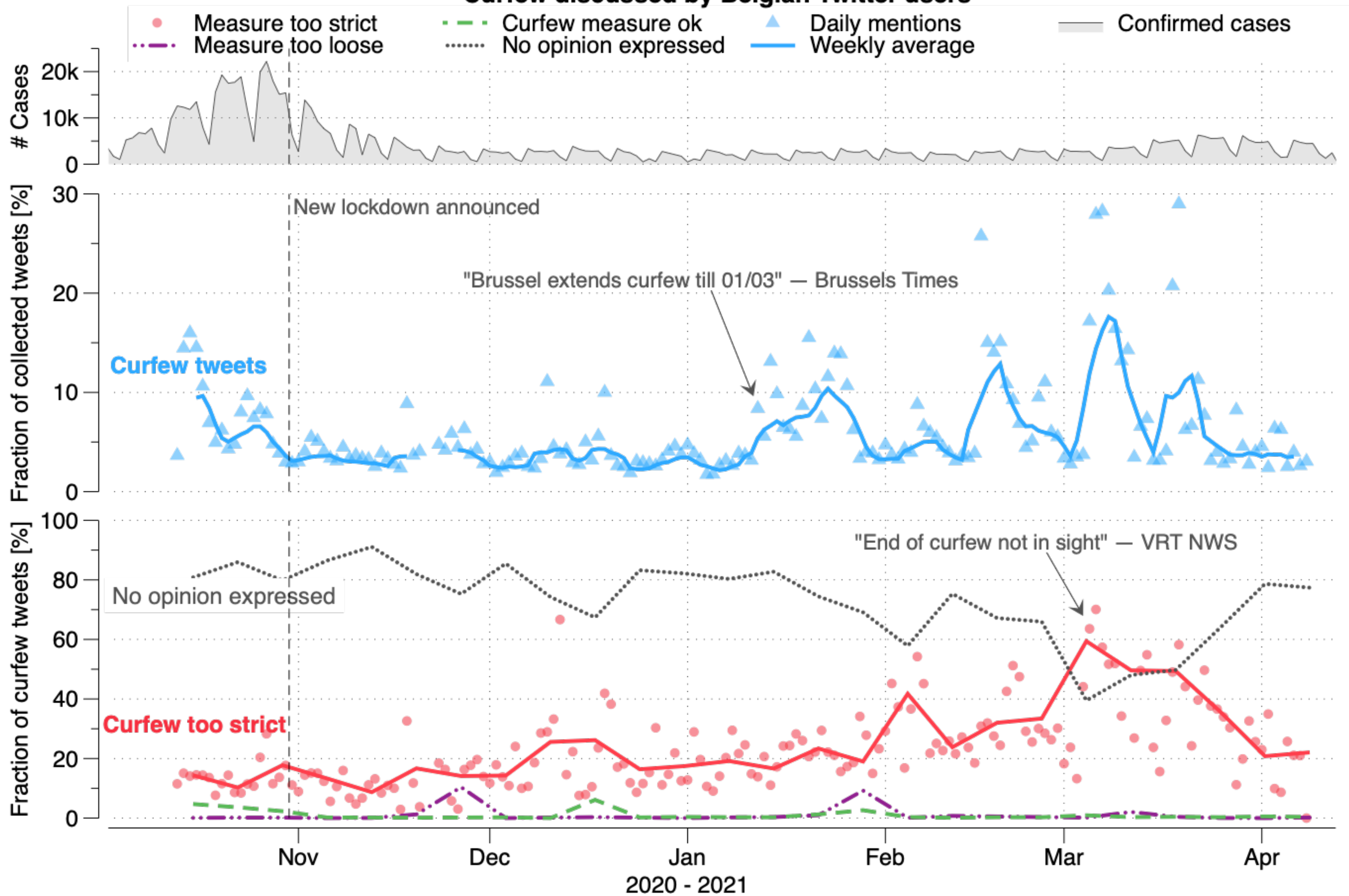
- NL (AUC = 0.94, model version = 2020-01-05)
- EN (AUC = 0.90, model version = 2020-01-05)
- FR (AUC = 0.69, model version = 2020-01-05)

(b) ROC curves conditioned on language (English, Dutch and French) for the best-performing model: mBERT trained on 2k tweets.

Vaccines discussed by Belgian Twitter-users



Curfew discussed by Belgian Twitter users



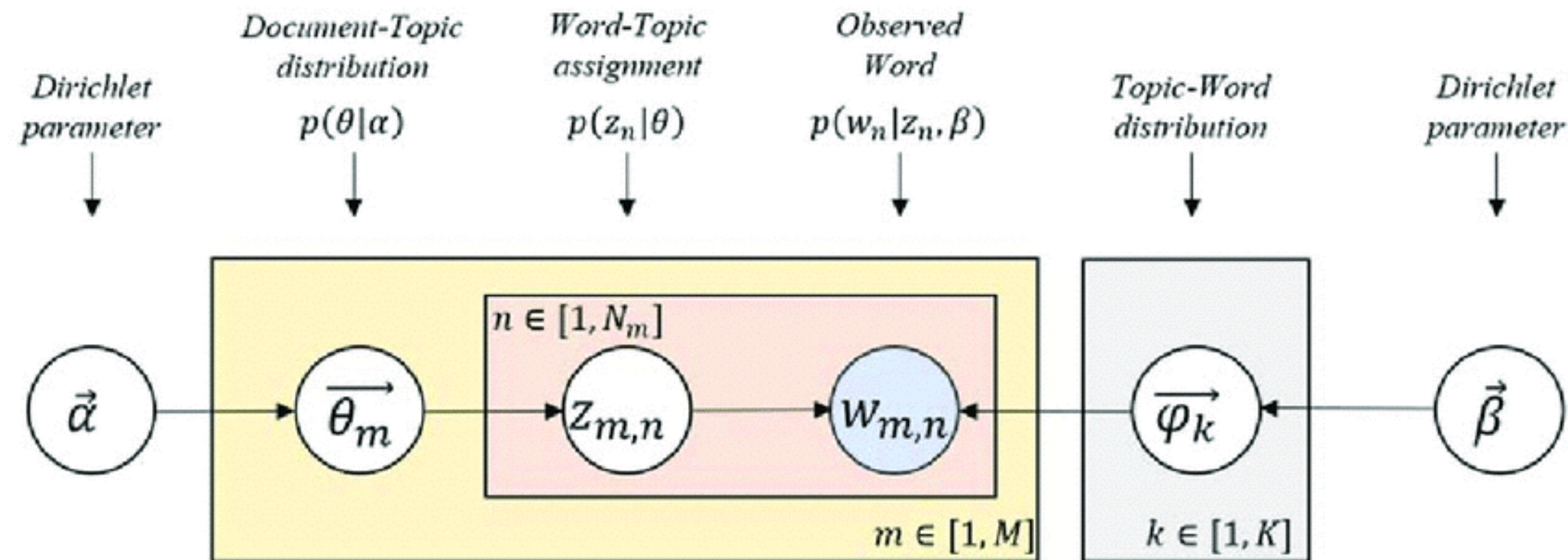
Topic modelling

Topic modelling

- Unsupervised machine learning
- Finding topics in a corpus
- Several approaches:
 - Latent Dirichlet Allocation (LDA)
 - BERTopic
 - ...

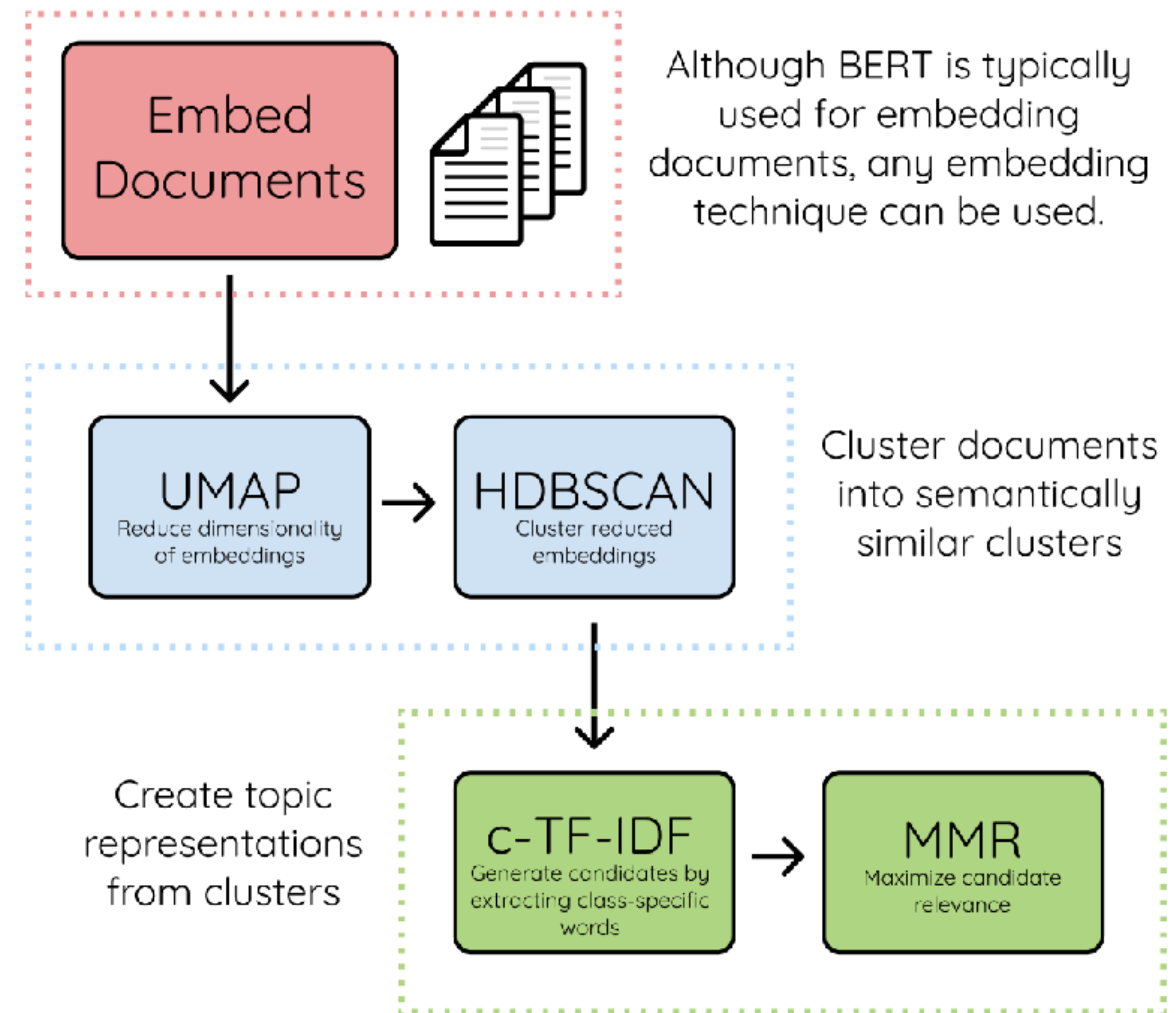
Latent Dirichlet Allocation

- Bayesian approach to model topics
 - words are generated by topics



BERTopic

- Topic modelling as clustering task
 1. Embedding with BERT
 2. Dim. reduction
 3. Clustering
 4. Get topic words



UMAP

- Dimensionality reduction
- Graph constructed from nearest neighbours

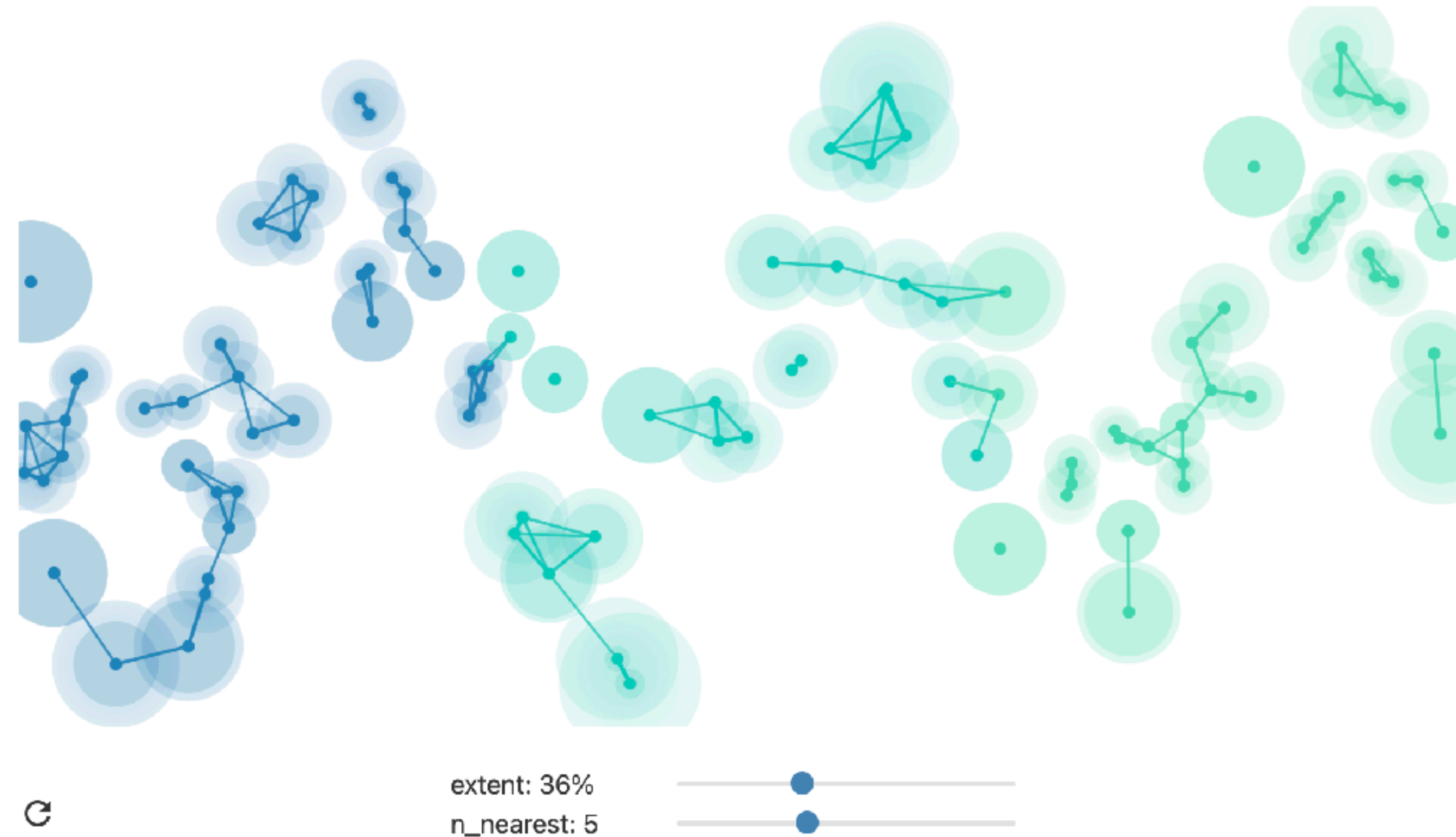


Figure 3: Adjust the slider to extend a radius outwards from each point, computed by the distance to its n th nearest neighbor. Notice that past the intersection with the first neighbor, the radius begins to get fuzzy, with subsequent connections appearing with less weight;

UMAP

- Dimensionality reduction
- Graph constructed from nearest neighbours
- Projected in 2D space by mapping the graph to 2D
- High-dim. neighbours are closer together in 2D

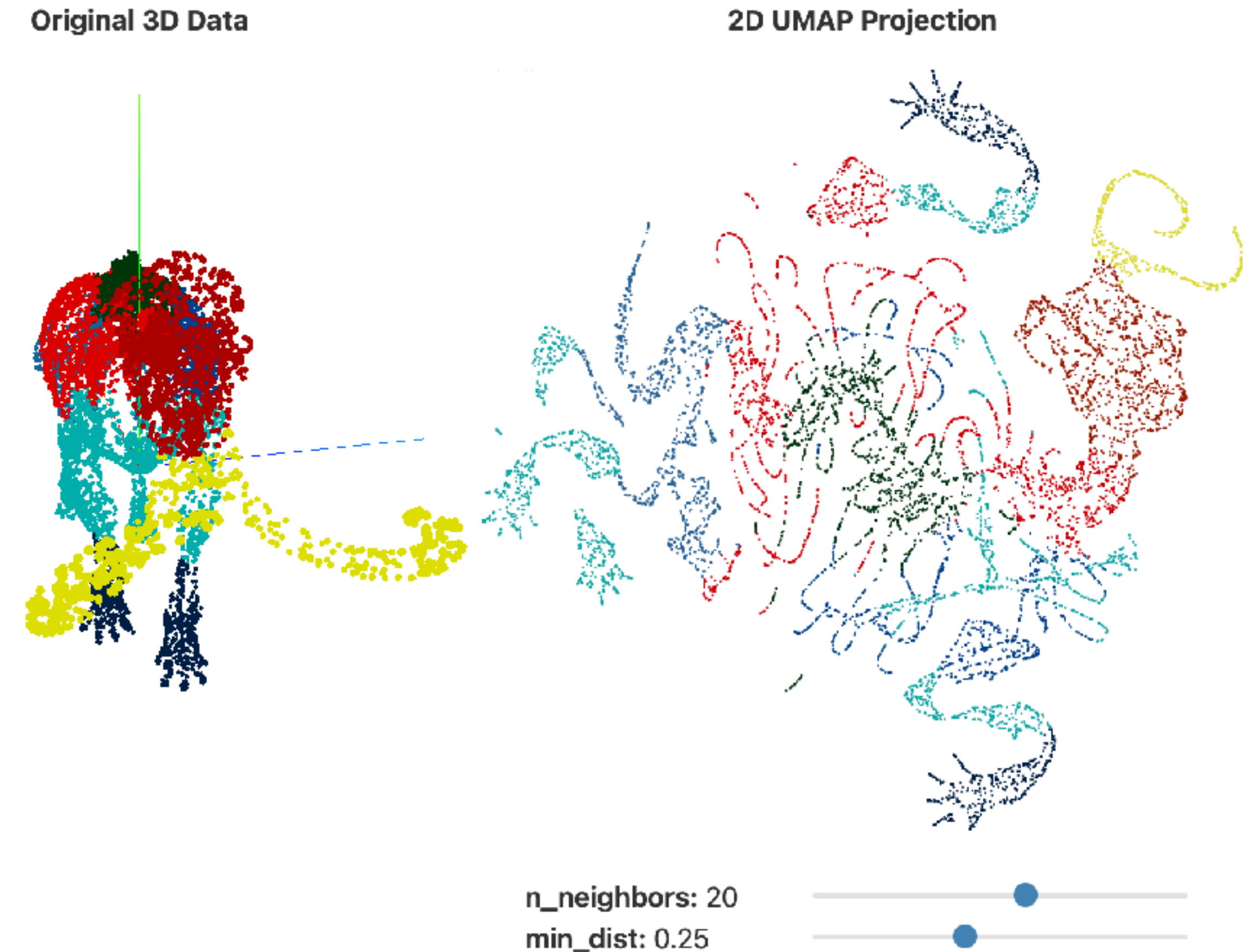


Figure 5: UMAP projections of a 3D woolly mammoth skeleton (50k points, 10k shown) into 2 dimensions, with various settings for the `n_neighbors` and `min_dist` parameters.

TF-IDF

- Measure of words that are unique for a document
- But corrected for words that are in every document (e.g. 'the')

$$tf(t, d) = \frac{f_{t,d}}{\sum_{t' \in d} f_{t',d}}$$

$$idf(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

BERTopic

Modularity

- Minimizing Assumptions -

Relatively few assumptions with respect to the dependencies of one algorithm on another

