

Lexical Semantic Change: Models, Data and Evaluation

LREC 2022 - Tutorial - 20 June 2022

Pierpaolo Basile¹, Annalina Caputo², Pierluigi Cassotti¹ and Rossella Varvara³
University of Bari¹, Dublin City University², Université de Fribourg³

Outline

- **Background concepts**
 - PPMI matrix factorization
 - Word2Vec Skip-gram with Negative Sampling (SGNS)
 - BERT-based models
- **Lexical Semantic Change Models**
 - Alignment models
 - Post-alignment models
 - Orthogonal Procrustes
 - Jointly Alignment Models
 - Explicit Alignment Models
 - Dynamic Word Embedding (DWE)
 - Dynamic Bernoulli Embedding (DBE)
 - Implicit Alignment Models
 - Temporal Referencing (TR)
 - Temporal Random Indexing (TRI)
 - Temporal Word Embedding with a Compass (TWEC)
 - Contextualized Models
 - TempoBERT
 - Temporal Attention
 - Deep Mistake
 - Gloss Reader
 - Other Models
 - Local Neighborhood measure
 - Word Sense Induction
 - Grammatical Features

PPMI

Factorization

Mutual Information

SYMMETRIC
NON NEGATIVE

$$I(X, Y) = \sum_x \sum_y P(x, y) \log_2 \frac{P(x, y)}{P(x)P(y)}$$

Mutual information (MI) of two random variables is a measure of the mutual dependence between the two variables.

It quantifies the "amount of information" obtained about one random variable by observing the other random variable

Pointwise Mutual Information

$$\text{PMI}(w, c) = \log_2 \frac{P(w, c)}{P(w)P(c)}$$

PMI(w,c) = 0 w and c are statistically independent

PMI(w,c) > 0 w and c co-occur more frequently than would be expected under an independence assumption

PMI(w,c) < 0 w and c co-occur less frequently than would be expected

Positive Pointwise Mutual Information

	computer	data	result	pie	sugar	count(w)
cherry	2	8	9	442	25	486
strawberry	0	0	1	60	19	80
digital	1670	1683	85	5	4	3447
information	3325	3982	378	5	13	7703
count(context)	4997	5673	473	512	61	11716

$$\text{PPMI}(w, c) = \max\left(\log_2 \frac{P(w, c)}{P(w)P(c)}, 0\right)$$

Positive Pointwise Mutual Information

	computer	data	result	pie	sugar	count(w)
cherry	2	8	9	442	25	486
strawberry	0	0	1	60	19	80
digital	1670	1683	85	5	4	3447
information	3325	3982	378	5	13	7703
count(context)	4997	5673	473	512	61	11716

$$p_{ij} = \frac{f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}, \quad p_{i*} = \frac{\sum_{j=1}^C f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}, \quad p_{*j} = \frac{\sum_{i=1}^W f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}$$

$$\text{PPMI}_{ij} = \max\left(\log_2 \frac{p_{ij}}{p_{i*} p_{*j}}, 0\right)$$

Positive Pointwise Mutual Information

PPMI_{digital,data}

	computer	data	result	pie	sugar	count(w)
cherry	2	8	9	442	25	486
strawberry	0	0	1	60	19	80
digital	1670	1683	85	5	4	3447
information	3325	3982	378	5	13	7703
count(context)	4997	5673	473	512	61	11716

$$p_{ij} = \frac{f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}, \quad p_{i*} = \frac{\sum_{j=1}^C f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}, \quad p_{*j} = \frac{\sum_{i=1}^W f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}$$

$$\text{PPMI}_{ij} = \max\left(\log_2 \frac{p_{ij}}{p_{i*} p_{*j}}, 0\right)$$

Positive Pointwise Mutual Information

PPMI_{digital,data}

	computer	data	result	pie	sugar	count(w)
cherry	2	8	9	442	25	486
strawberry	0	0	1	60	19	80
digital	1670	1683	85	5	4	3447
information	3325	3982	378	5	13	7703
count(context)	4997	5673	473	512	61	11716

$$p_{ij} = \frac{f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}, \quad p_{i*} = \frac{\sum_{j=1}^C f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}, \quad p_{*j} = \frac{\sum_{i=1}^W f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}$$

$$\text{PPMI}_{ij} = \max\left(\log_2 \frac{p_{ij}}{p_{i*} p_{*j}}, 0\right)$$

Positive Pointwise Mutual Information

PPMI_{digital,data}

	computer	data	result	pie	sugar	count(w)
cherry	2	8	9	442	25	486
strawberry	0	0	1	60	19	80
digital	1670	1683	85	5	4	3447
information	3325	3982	378	5	13	7703
count(context)	4997	5673	473	512	61	11716

$$p_{ij} = \frac{f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}, \quad p_{i*} = \frac{\sum_{j=1}^C f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}, \quad p_{*j} = \frac{\sum_{i=1}^W f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}$$

$$\text{PPMI}_{ij} = \max\left(\log_2 \frac{p_{ij}}{p_{i*} p_{*j}}, 0\right)$$

Positive Pointwise Mutual Information

PPMI_{digital,data}

	computer	data	result	pie	sugar	count(w)
cherry	2	8	9	442	25	486
strawberry	0	0	1	60	19	80
digital	1670	1683	85	5	4	3447
information	3325	3982	378	5	13	7703
count(context)	4997	5673	473	512	61	11716

$$p_{ij} = \frac{f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}, \quad p_{i*} = \frac{\sum_{j=1}^C f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}, \quad p_{*j} = \frac{\sum_{i=1}^W f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}$$

$$\text{PPMI}_{ij} = \max\left(\log_2 \frac{p_{ij}}{p_{i*} p_{*j}}, 0\right)$$

Positive Pointwise Mutual Information

	computer	data	result	pie	sugar
cherry	0	0	0	4.38	3.30
strawberry	0	0	0	4.10	5.51
digital	0.18	0.01	0	0	0
information	0.02	0.09	0.28	0	0

PPMI Factorization

	computer	data	result	pie	sugar
cherry	0	0	0	4.38	3.30
strawberry	0	0	0	4.10	5.51
digital	0.18	0.01	0	0	0
information	0.02	0.09	0.28	0	0

$m \times n$

=

U

$m \times m$

	0	0
0		0
0	0	
0	0	0

Σ

$m \times n$

V*

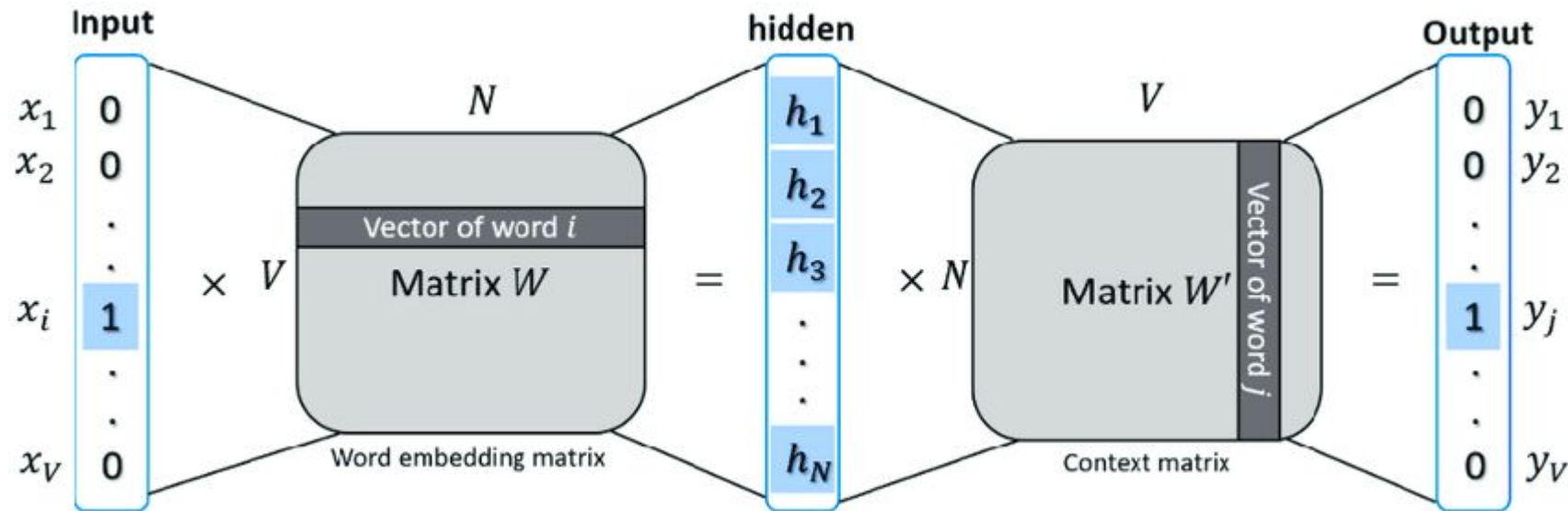
$n \times n$

Word2Vec

Skip-gram with

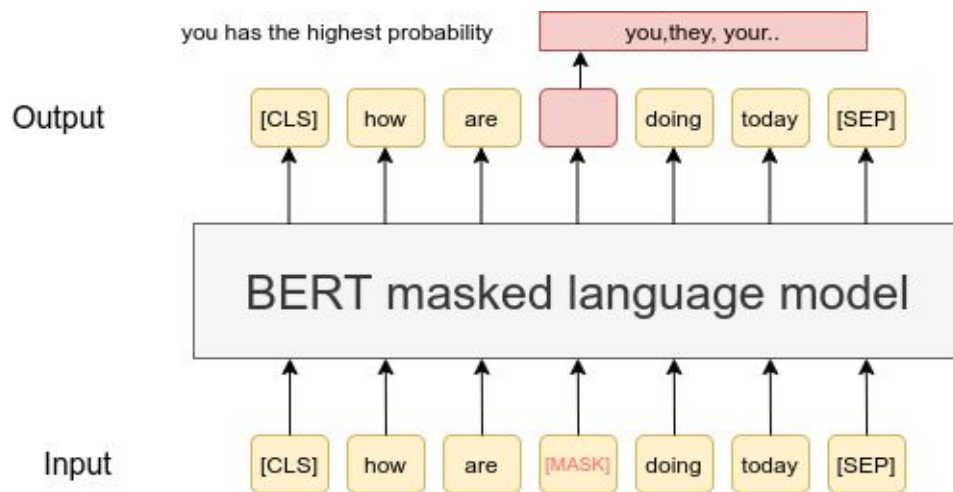
Negative sampling

Word2Vec Skip-gram with Negative Sampling (SGNS)



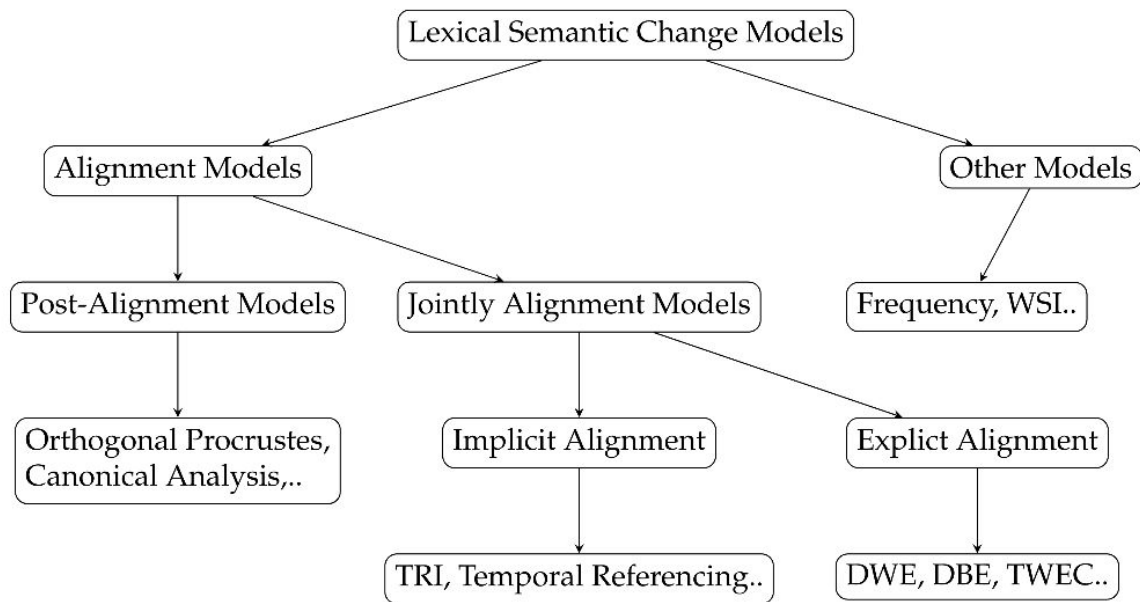
BERT-based models

BERT-based models



Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Lexical Semantic Change Models



Alignment Models

Alignment approach

Post-alignment

- Post-alignment models first train static word embeddings for each time slice and then align them

Jointly alignment

- Jointly Alignment models train word embeddings and jointly align vectors across all time slices
- Jointly Alignment models can be distinguished in Explicit alignment models and Implicit alignment models.

Alignment approach

Post-alignment

- Post-alignment models first train static word embeddings for each time slice and then align them

Jointly alignment

- Jointly Alignment models train word embeddings and jointly align vectors across all time slices
- Jointly Alignment models can be distinguished in Explicit alignment models and Implicit alignment models.



Post-alignment and Explicit alignment rely on the assumption that only few words change their meaning

Orthogonal Procrustes (OP)

$$R = \arg \min_{Q^T Q = I} \|QW^t - W^{t+1}\|_F$$

The diagram illustrates the Orthogonal Procrustes problem. The equation $R = \arg \min_{Q^T Q = I} \|QW^t - W^{t+1}\|_F$ is shown. Three orange arrows point from parts of the equation to labels: one from R to "Rotation matrix", one from Q to "Embedding matrix time t", and one from W^{t+1} to "Embedding matrix time t+1".

Jointly Alignment - Alignment constraint

Explicit alignment

- The objective function of explicit alignment models involves constraints on word vectors
- Typically those constraints require that the distance of two-word vectors in two consecutive periods is the smallest possible

Implicit alignment

- In the implicit alignment, the alignment is automatically performed by sharing the same word context vectors across all the time spans

Dynamic Word Embedding (DWE)

$$\min_{U(t)} \frac{1}{2} \|Y(t) - U(t)U(t)^T\|_F^2 + \frac{\lambda}{2} \|U(t)\|_F^2 + \frac{\tau}{2} \left(\|U(t-1) - U(t)\|_F^2 + \|U(t) - U(t+1)\|_F^2 \right)$$

PMI factorization

Regularization term

Explicit temporal constraints

Dynamic Bernoulli Embedding (DBE)

$$\mathcal{L}(\boldsymbol{\rho}, \boldsymbol{\alpha}) = \mathcal{L}_{\text{pos}} + \mathcal{L}_{\text{neg}} + \mathcal{L}_{\text{prior}}$$

$$\mathcal{L}_{\text{pos}} = \sum_{i=1}^N \sum_{v=1}^V x_{iv} \log \sigma(\eta_{iv})$$

$$\mathcal{L}_{\text{neg}} = \sum_{i=1}^N \sum_{v=1}^V (1 - x_{iv}) \log(1 - \sigma(\eta_{iv}))$$

$$\eta_{iv} = \rho_v^{(t_i)\top} \left(\sum_{j \in c_j} \sum_{v'} \alpha_{v'} x_{jv'} \right)$$

embedding vector $\rho_v \in \mathbb{R}^K$

context vector $\alpha_v \in \mathbb{R}^K$

$$\mathcal{L}_{\text{prior}} = \log p(\boldsymbol{\alpha}) + \log p(\boldsymbol{\rho})$$

$$\log p(\boldsymbol{\alpha}) = -\frac{\lambda_0}{2} \sum_v \|\alpha_v\|^2$$

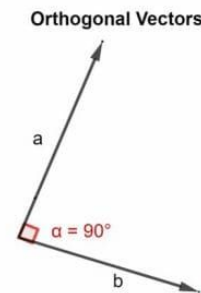
$$\log p(\boldsymbol{\rho}) = -\frac{\lambda_0}{2} \sum_v \|\rho_v^{(0)}\|^2 - \frac{\lambda}{2} \sum_{v,t} \|\rho_v^{(t)} - \rho_v^{(t-1)}\|^2$$

Regularization
term

Explicit temporal
constraints

Temporal Random Indexing (TRI)

- Produce aligned word embeddings in a single step.
- Count-based method.
- TRI is based on Random Indexing: near-orthogonality random index vectors shared across all time slices so that word spaces are comparable.

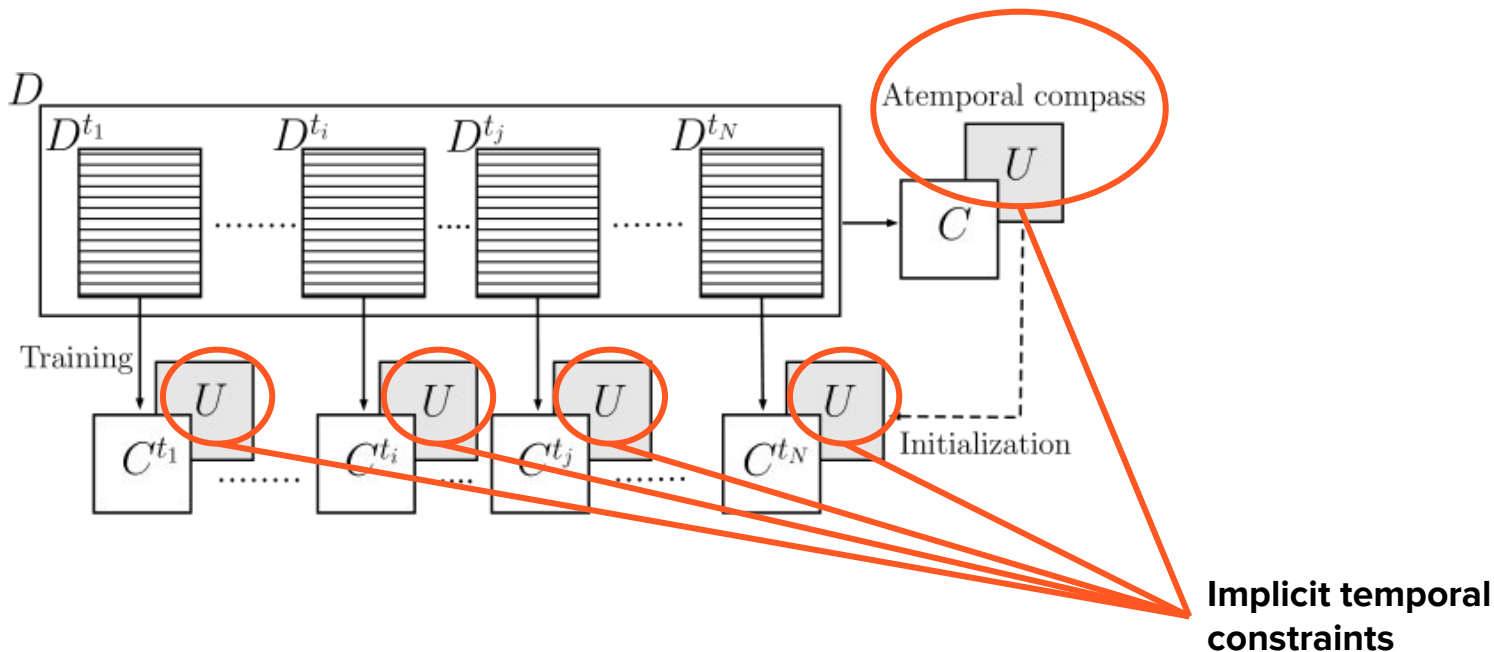


$$sv_i = \sum_{d \in C} \sum_{-m < i < +m} c_i$$

Temporal Referencing (TR)

- Replace a subset of words in the dictionary (target words) with time-specific tokens
- Temporal Referencing is not performed when the word is considered a context word
- Since TR is a generic framework, it can be applied to both low-dimensional embeddings learned with SGNS and high-dimensional sparse PPMI vectors


Temporal Word Embedding with a Compass (TWEC)



Contextualized Models

TempoBERT

- Use time as additional context
- Exploit time masking



YEAR: 1800 → "<1800> The mountains have an awful majesty."
YEAR: 2020 → "<2020> You look awful today."

(a) TempoBERT is trained on temporal corpora, where each sequence is prepended with temporal context information.

Time prediction: "[MASK] Today's weather is awful." → <2020>

Time-dependent MLM:
<1800> He has an awful [MASK]. → presence
<2020> He has an awful [MASK]. → temper

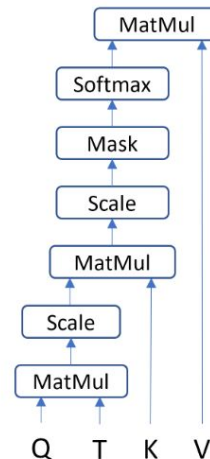
(b) TempoBERT can be used for inference in two modes: (1) time prediction; (2) time-dependent mask filling.

Temporal Attention

- Extends self-attention to include time dimension

$$\text{TemporalAttention}(Q, K, V, T) = \text{softmax} \left(\frac{Q \frac{T^T T}{\|T\|} K^T}{\sqrt{d_k}} \right) V$$

Time-specific weight matrix



XLM-RoBERTa

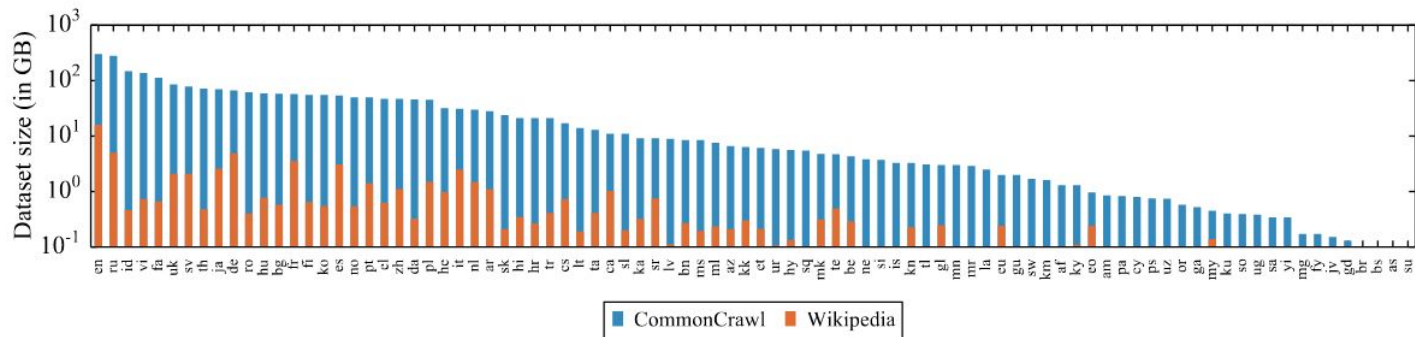


Figure 1: Amount of data in GiB (log-scale) for the 88 languages that appear in both the Wiki-100 corpus used for mBERT and XLM-100, and the CC-100 used for XLM-R. CC-100 increases the amount of data by several orders of magnitude, in particular for low-resource languages.

Gloss Reader

- Rely on XLM-RoBERTa and trained on an English Word Sense Disambiguation (WSD) dataset (SemCor)
- Zero-shot ability on other languages such as Russian

The **bank** can guarantee deposits will eventually cover future tuition costs because it invests in adjustable-rate mortgage securities.



Context Encoder

bank ¹	Gloss:	a financial institution that accepts deposits and channels the money into lending activities
-------------------	--------	--



Gloss Encoder

bank ²	Gloss:	sloping land (especially the slope beside a body of water)
-------------------	--------	--

Deep Mistake

- Pretrained XLM-R finetuned on MCL-WiC task
- Not depends on fixed sense inventories

Lang	Target	Context-1	Context-2	Label
EN	Beat	We <u>beat</u> the competition	Agassi <u>beat</u> Becker in the tennis championship.	True
DA	Tro	Jeg <u>tror</u> p ^o a det, min mor fortalte.	Maria <u>troede</u> ikke sine egne øjne.	True
ET	Ruum	Uhel hetkel olin v ^o aljaspool aega ja <u>ruumi</u> .	Umberringi oli l ^o oputu t ^o uhi <u>ruum</u> .	True
FR	Causticité	Sa <u>causticité</u> lui a fait bien des ennemis.	La <u>causticité</u> des acides.	False
KO	틀림	틀림이 있는지 없는지 세어 보시오.	그 아이 하는 짓에 틀림이 있다면 모두 이 어미 죄이지요.	False
ZH	發	建築師希望發大火燒掉城市的三分之一。	如果南美洲氣壓偏低，則印度可能發乾旱	True
FA	صرف	صرف غذا نیم ساعت طول کشید	معلم صرف افعال ماضی عربی را آموزش داد	False

Other models

Local Neighborhood measure

- Global measure

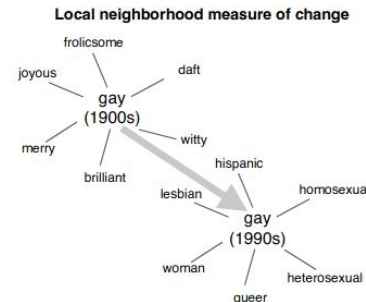
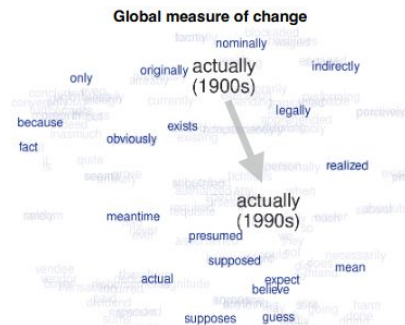
$$d^G(w_i^{(t)}, w_i^{(t+1)}) = \text{cos-dist}(\mathbf{w}_i^{(t)}, \mathbf{w}_i^{(t+1)})$$

- Local Neighborhood measure

$$\mathbf{s}^{(t)}(j) = \text{cos-sim}(\mathbf{w}_i^{(t)}, \mathbf{w}_j^{(t)})$$

$$\forall w_j \in \mathcal{N}_k(w_i^{(t)}) \cup \mathcal{N}_k(w_i^{(t+1)})$$

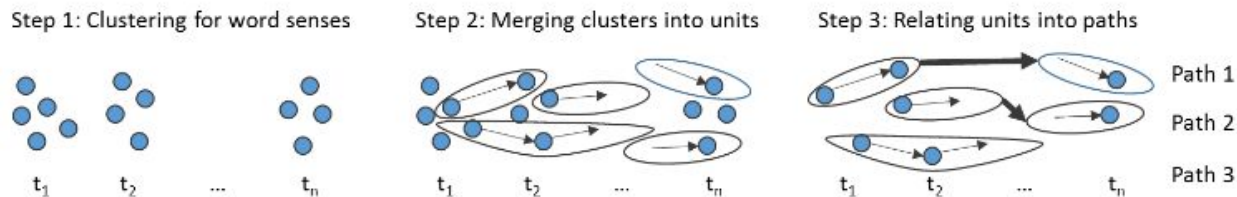
$$d^L(w_i^{(t)}, w_i^{(t+1)}) = \text{cos-dist}(\mathbf{s}_i^{(t)}, \mathbf{s}_i^{(t+1)})$$



***k* nearest-neighbors**

Word Sense Induction

- Curvature clustering
- *lin* measure (based on the WordNet synset similarity)



Grammatical Features

- Grammatical features such as PoS tags, dependency labels, number, case, tense
- Grammatical features are language-dependent
- Interpretable results

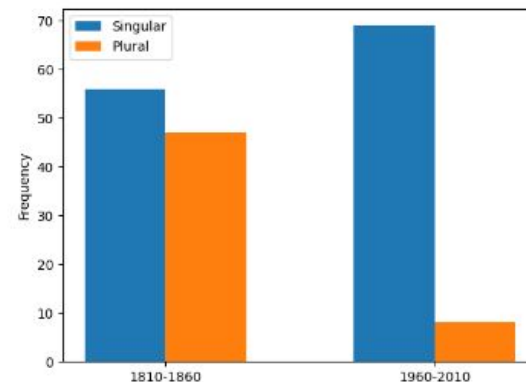


Figure 1: Changes in the number category distribution for the English noun 'lass' over time, calculated on the English corpora of the SemEval 2020 shared task 1 (Schlechtweg et al., 2020). 'Lass' is annotated as semantically changed in the SemEval dataset.