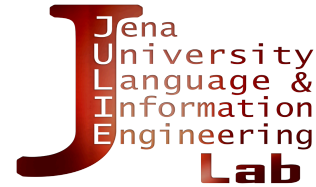




FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA



MODELL
ROMANTIK



Don't Get Fooled by Word Embeddings— Better Watch Their Neighborhood



Johannes Hellrich^{1,2}

1: Graduate School 'The Romantic Model',
Friedrich Schiller University Jena,
Jena, Germany

<http://www.modellromantik.uni-jena.de>



& Udo Hahn²

2: Jena University Language & Information
Engineering (JULIE) Lab
Friedrich Schiller University Jena,
Jena, Germany

<http://www.julielab.de>

You shall know a word by the company it keeps!

Firth, 1957

He reads a poem.

She reads a novel.

The novel has 312 pages.

The poem fits on two pages.

She listens to an opera.

He listens to jazz.

You shall know a word by the company it keeps!

Firth, 1957

He reads a poem.

She reads a novel.

The novel has 312 pages.

The poem fits on two pages.

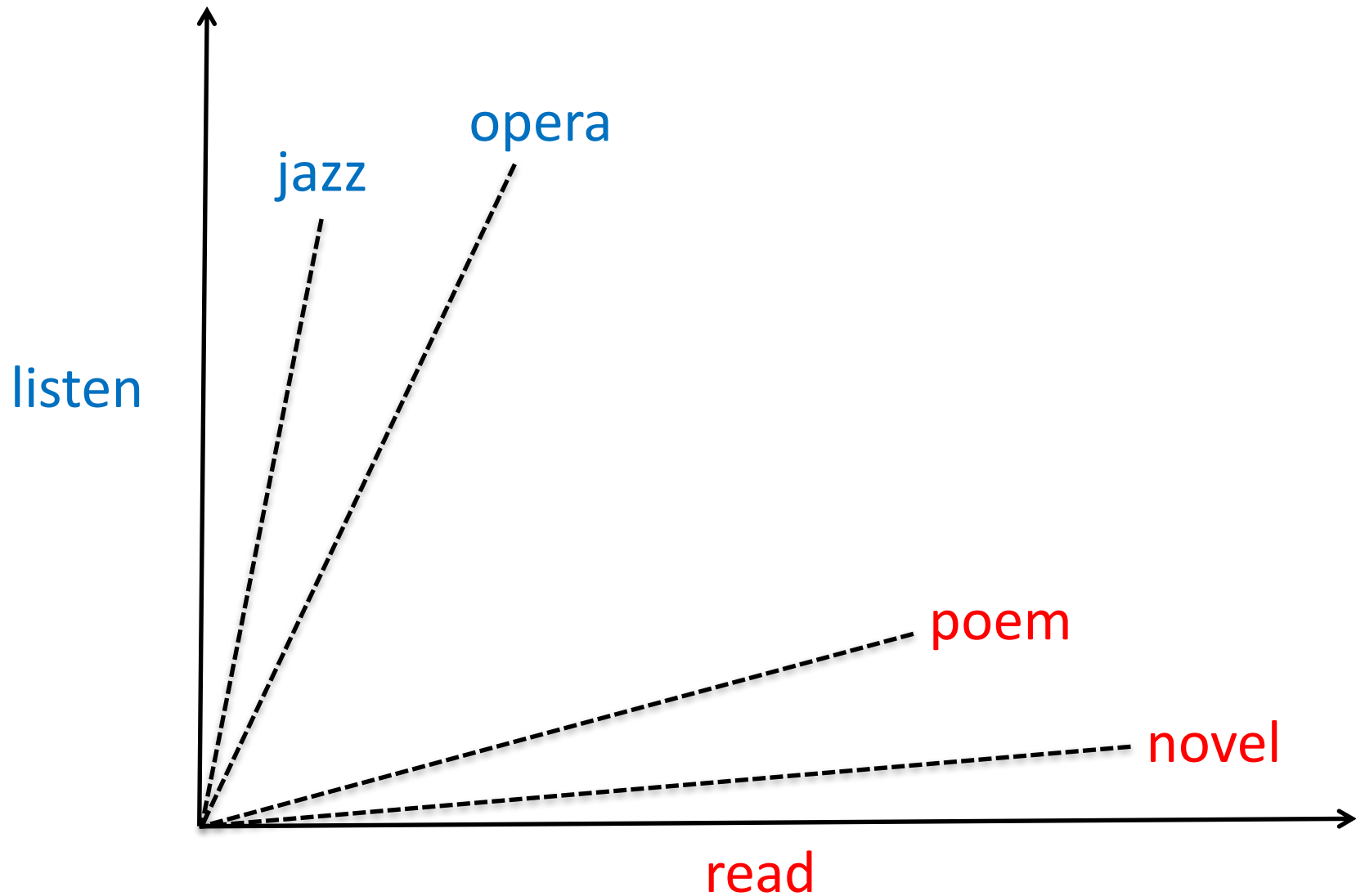
She listens to an opera.

He listens to jazz.

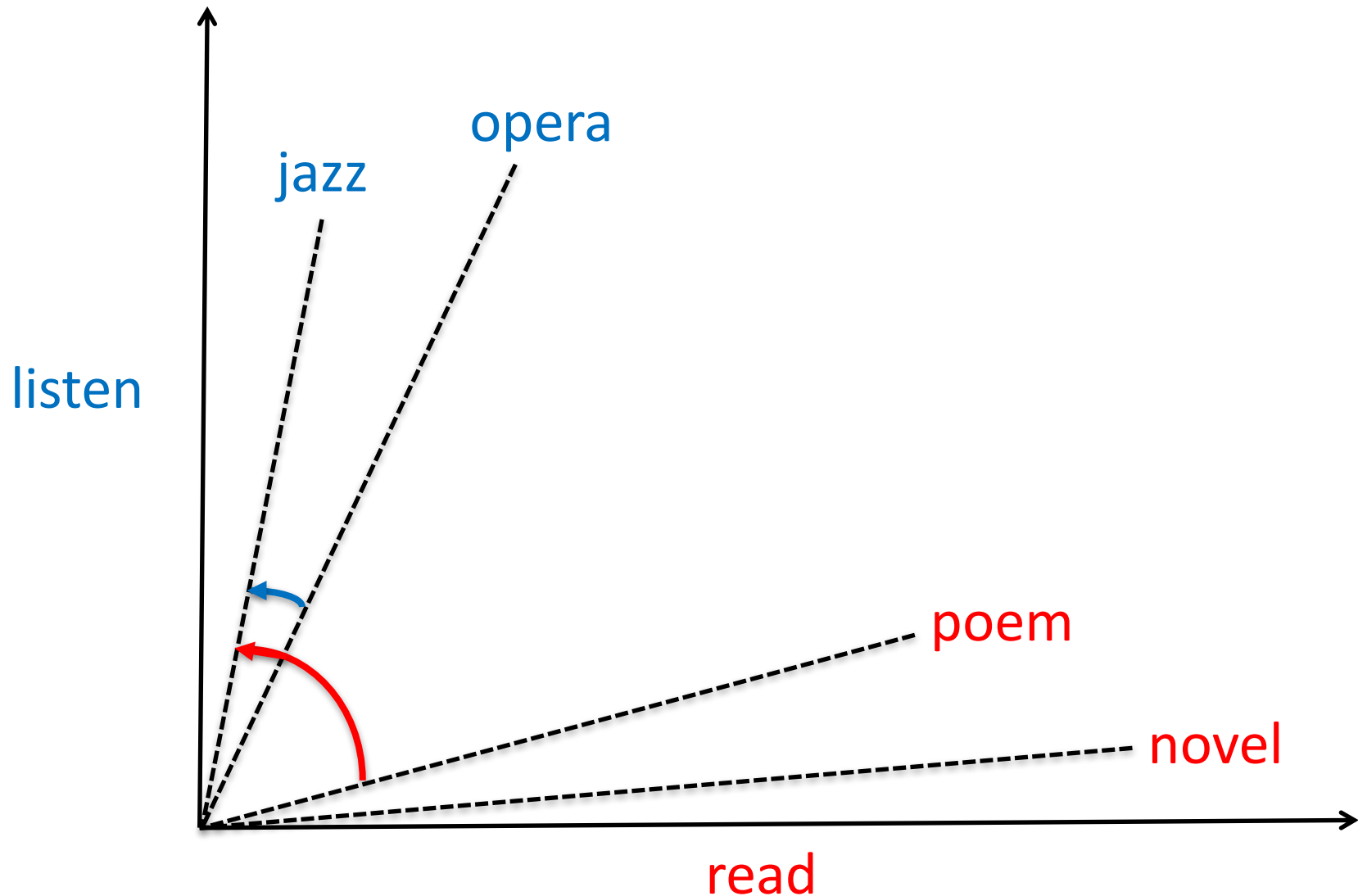
Counting Cooccurrences

	read	pages	hate	enjoy	listen	...
novel	98	60	3	56	2	
poem	67	10	1	47	8	
opera	4	8	0	42	38	
jazz	2	1	2	61	47	
...						

Vector Representation



Distance and Similarity



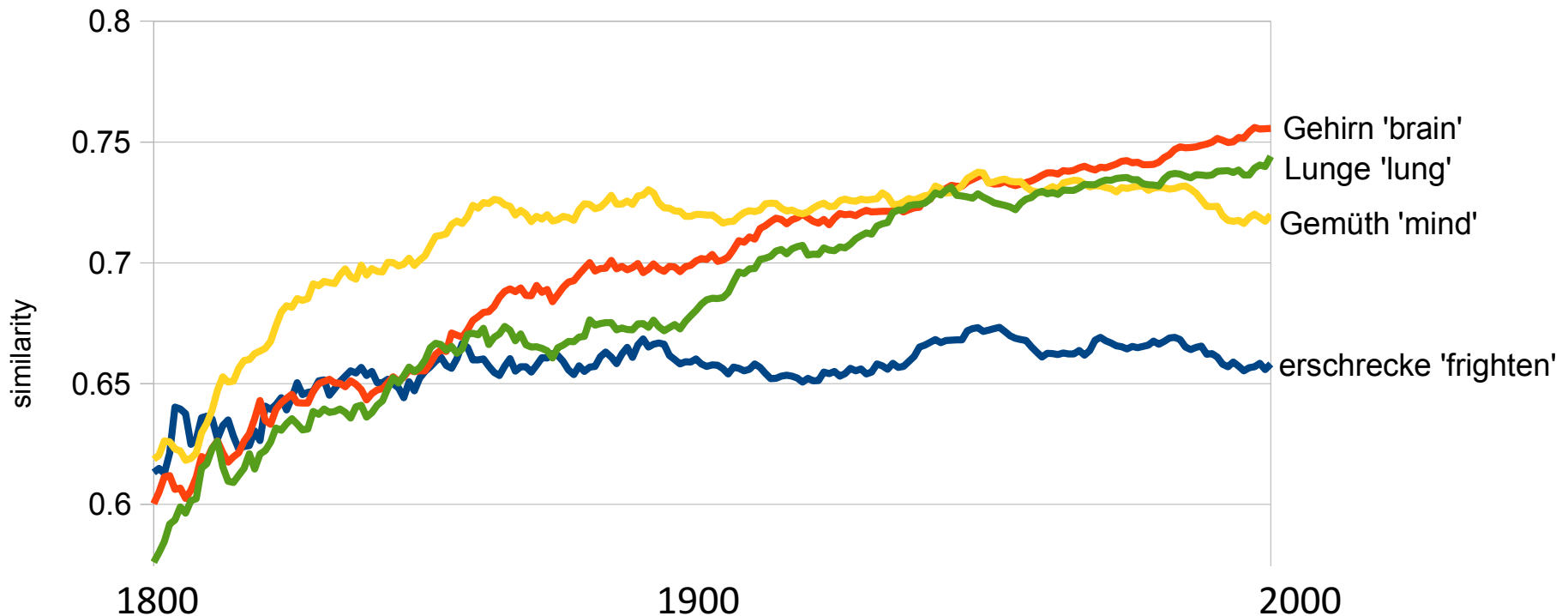
Dimensionality Problem

- One dimension per word
 - 50k to 100k dimensions
 - Large files and slow operations
- What about synonyms – it shouldn't matter if I **buy** or **purchase** a **novel**

Word Embeddings

- Represent words as dense vectors with 200–500 instead of 50k–100k dimensions
- Very popular in computational linguistics and digital humanities
- Better on judging word similarity

Application in DH: Semantic Development of *Herz* ,heart'



- Hellrich & Hahn, DH 2016
- First applied by Kim et al., *ACL 2014 Workshop on Language Technologies and Computational Social Science*

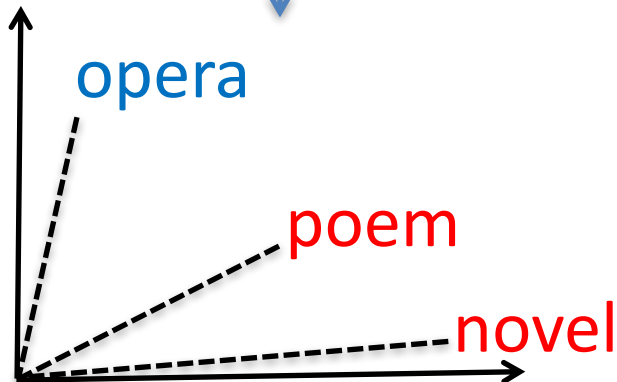
Types of Word Embeddings

Singular Value Decomposition

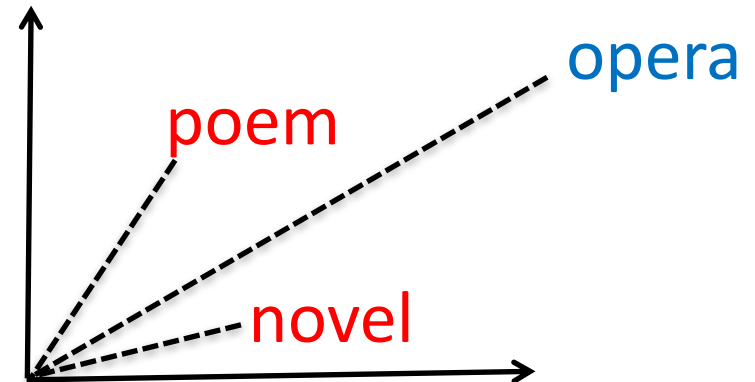
lots
of
text



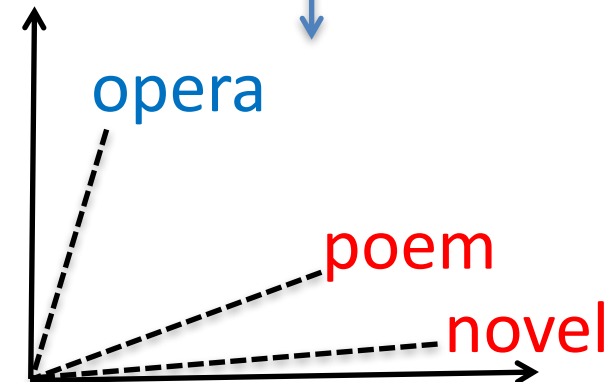
	read	pages	musician
poem	475	156	0
novel	823	492	3
opera	51	19	993



Neural Word Embeddings

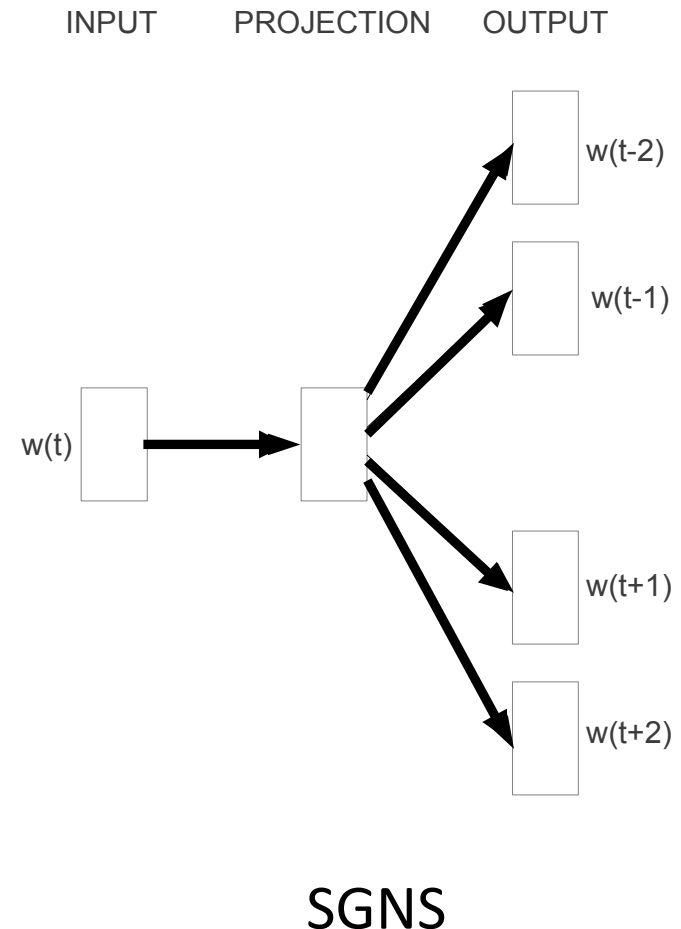


lots
of
text

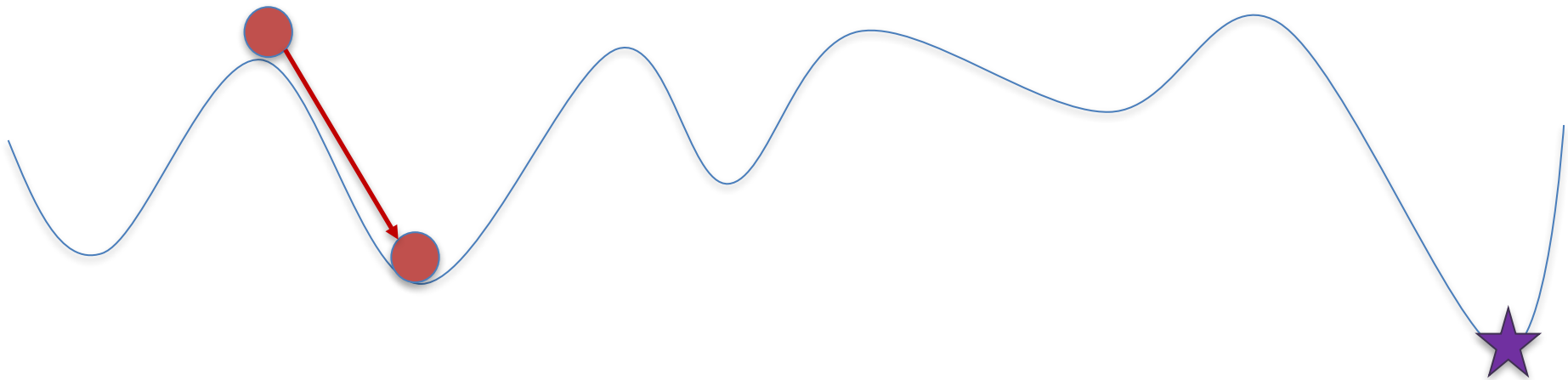


Neural Word Embeddings

- Extremely popular skip-gram negative sampling algorithm **SGNS/word2vec** (Mikolov et al., NIPS 2013)
- Alternative neural embeddings using an explicit cooccurrence matrix: **GloVe** (Pennington et al., EMNLP 2014)

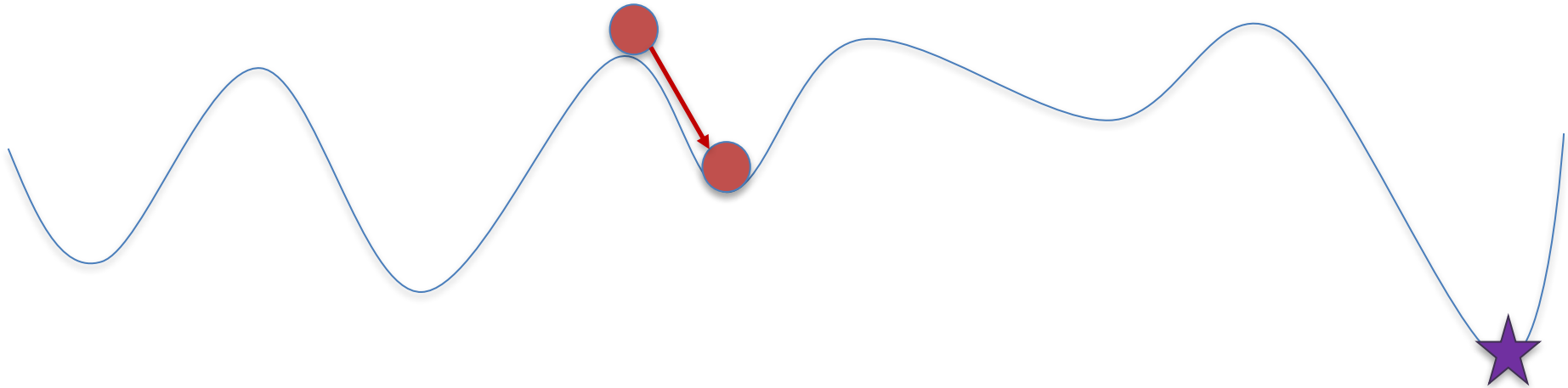


Training Neural Word Embeddings



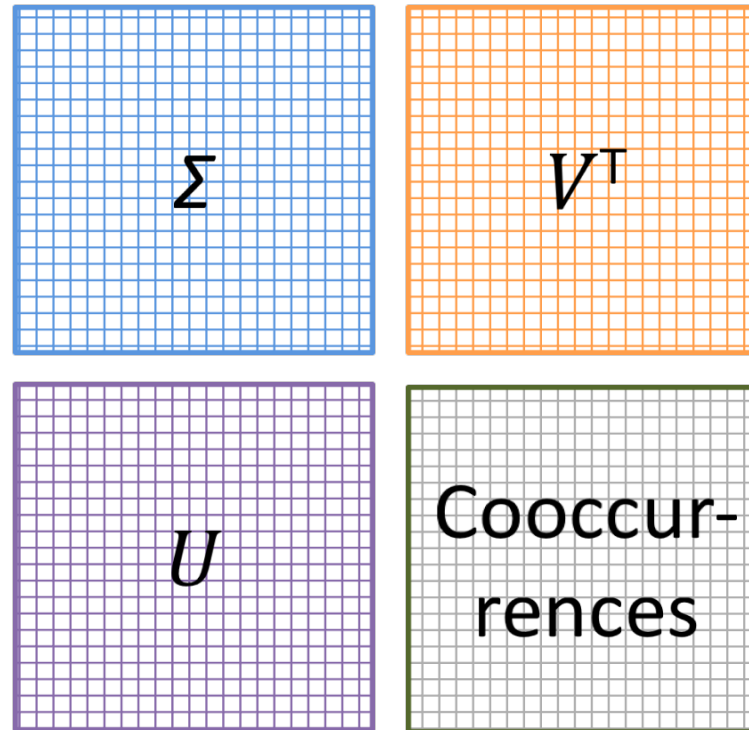
- Word Embeddings are updated after looking at the text
- Tries to minimize false predictions (cost function)
- Will lead us to a **local**, yet rarely to the **global** minimum

Training Neural Word Embeddings



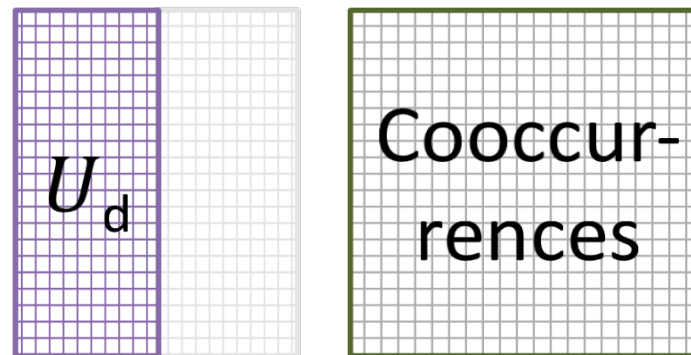
- Word Embeddings are updated after looking at the text
- Tries to minimize false predictions (cost function)
- Will lead us to a **local**, yet rarely to the **global** minimum

Singular Value Decomposition



- Express Cooccurrences as $U\Sigma V^T$
 - U represents words, V^T context words
 - Σ measures importance of dimensions

Singular Value Decomposition



- Classical SVD embeddings: U_d , selecting only d dimensions from U based on Σ

SVD_{PPMI}

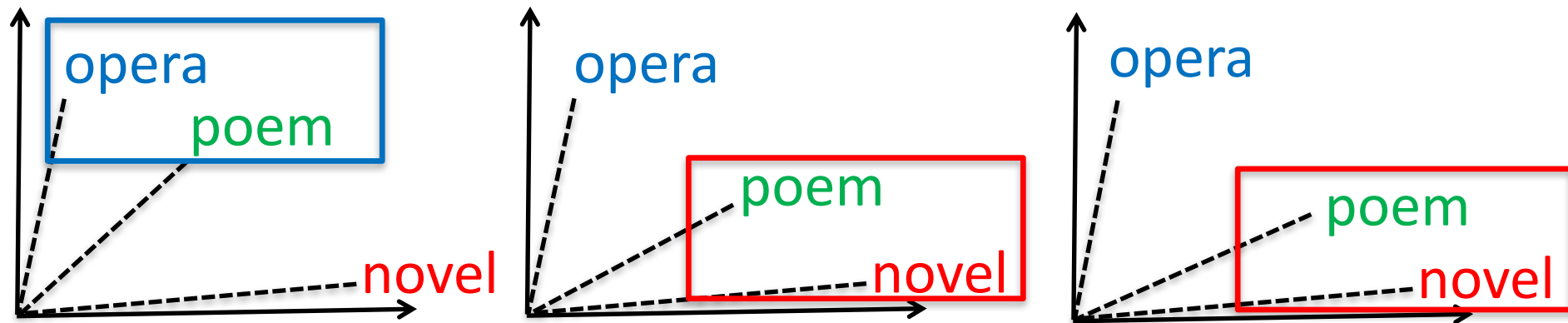
- Levy et al., TACL 2015
- Positive pointwise mutual information instead of frequency
- Post-/preprocessing inspired by SGNS and GloVe

Measuring Reliability

- Train multiple models with identical parameters on one corpus
- Measure percentage of identical neighborhoods for each word between models
- Hellrich&Hahn, COLING 2016

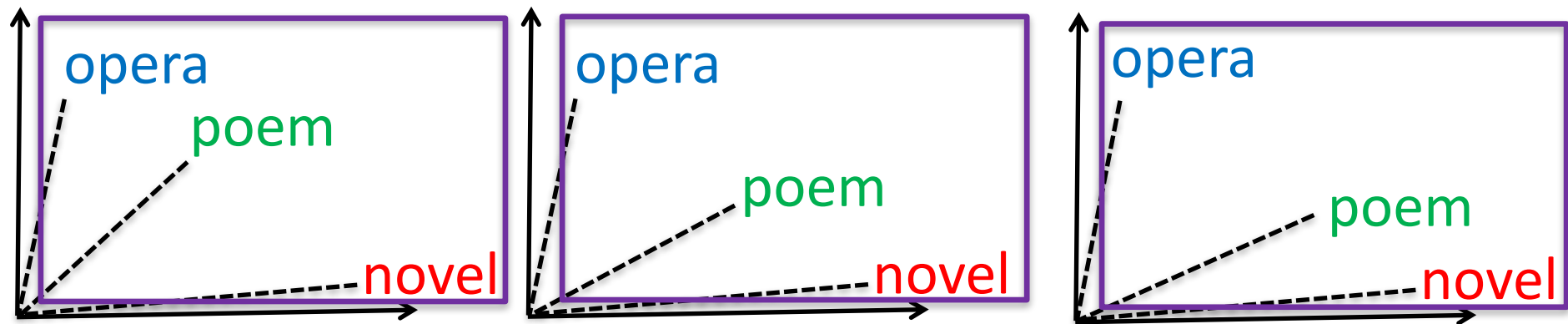
Measuring Reliability

- Train multiple models with identical parameters on one corpus
- Measure percentage of identical neighborhoods for each word between models
- Example: No agreement at neighborhood size 1 for **poem**



Measuring Reliability

- Train multiple models with identical parameters on one corpus
- Measure percentage of identical neighborhoods for each word between models
- Example: Agreement at neighborhood size 2 for poem



Experiment

- 3 models each for SGNS, GloVe and SVD_{PPMI}
- Trained on corpus of 645 German texts from 19th century, subset of Deutsches Textarchiv ‘German Text Archive’
- Technical Details:
 - Window size 5,
 - 300 dimensions
 - hyperwords toolkit

Reliability for *Herz* ‘heart’

Embedding Model	First Neighbor	Second Neighbor	Third Neighbor	Fourth Neighbor	Fifth Neighbor
SGNS 1	schmerzen ‘pain’	bekommen ‘anxious’	busen ‘bosom’	bluten ‘to bleed’	herzen ‘to caress’
SGNS 2	bluten ‘to bleed’	klopfend ‘beating’	busen ‘bosom’	bekommen ‘anxious’	herzen ‘to caress’
SGNS 3	herzen ‘to caress’	busen ‘bosom’	klopfend ‘beating’	bekommen ‘anxious’	bluten ‘to bleed’
GloVe 1	gemüt ‘mind’	mein ‘my’	seele ‘soul’	liebe ‘love’	brust ‘chest’
GloVe 2	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
GloVe 3	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
SVD_{PPMI} all	busen ‘bosom’	fühlen ‘to feel’	liebe ‘love’	schmerzen ‘pain’	menschenherz ‘human heart’

Reliability for *Herz* ‘heart’

Embedding Model	First Neighbor	Second Neighbor	Third Neighbor	Fourth Neighbor	Fifth Neighbor
SGNS 1	schmerzen ‘pain’	bekommen ‘anxious’	busen ‘bosom’	bluten ‘to bleed’	herzen ‘to caress’
SGNS 2	bluten ‘to bleed’	klopfend ‘beating’	busen ‘bosom’	bekommen ‘anxious’	herzen ‘to caress’
SGNS 3	herzen ‘to caress’	busen ‘bosom’	klopfend ‘beating’	bekommen ‘anxious’	bluten ‘to bleed’
GloVe 1	gemüt ‘mind’	mein ‘my’	seele ‘soul’	liebe ‘love’	brust ‘chest’
GloVe 2	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
GloVe 3	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
SVD_{PPMI} all	busen ‘bosom’	fühlen ‘to feel’	liebe ‘love’	schmerzen ‘pain’	menschenherz ‘human heart’

Reliability for *Herz* ‘heart’

Embedding Model	First Neighbor	Second Neighbor	Third Neighbor	Fourth Neighbor	Fifth Neighbor
SGNS 1	schmerzen ‘pain’	bekommen ‘anxious’	busen ‘bosom’	bluten ‘to bleed’	herzen ‘to caress’
SGNS 2	bluten ‘to bleed’	klopfend ‘beating’	busen ‘bosom’	bekommen ‘anxious’	herzen ‘to caress’
SGNS 3	herzen ‘to caress’	busen ‘bosom’	klopfend ‘beating’	bekommen ‘anxious’	bluten ‘to bleed’
GloVe 1	gemüt ‘mind’	mein ‘my’	seele ‘soul’	liebe ‘love’	brust ‘chest’
GloVe 2	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
GloVe 3	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
SVD_{PPMI} all	busen ‘bosom’	fühlen ‘to feel’	liebe ‘love’	schmerzen ‘pain’	menschenherz ‘human heart’

Reliability for *Herz* ‘heart’

Embedding Model	First Neighbor	Second Neighbor	Third Neighbor	Fourth Neighbor	Fifth Neighbor
SGNS 1	schmerzen ‘pain’	bekommen ‘anxious’	busen ‘bosom’	bluten ‘to bleed’	herzen ‘to caress’
SGNS 2	bluten ‘to bleed’	klopfend ‘beating’	busen ‘bosom’	bekommen ‘anxious’	herzen ‘to caress’
SGNS 3	herzen ‘to caress’	busen ‘bosom’	klopfend ‘beating’	bekommen ‘anxious’	bluten ‘to bleed’
GloVe 1	gemüt ‘mind’	mein ‘my’	seele ‘soul’	liebe ‘love’	brust ‘chest’
GloVe 2	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
GloVe 3	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
SVD_{PPMI} all	busen ‘bosom’	fühlen ‘to feel’	liebe ‘love’	schmerzen ‘pain’	menschenherz ‘human heart’

Reliability for *Herz* ‘heart’

Embedding Model	First Neighbor	Second Neighbor	Third Neighbor	Fourth Neighbor	Fifth Neighbor
SGNS 1	schmerzen ‘pain’	bekommen ‘anxious’	busen ‘bosom’	bluten ‘to bleed’	herzen ‘to caress’
SGNS 2	bluten ‘to bleed’	klopfend ‘beating’	busen ‘bosom’	bekommen ‘anxious’	herzen ‘to caress’
SGNS 3	herzen ‘to caress’	busen ‘bosom’	klopfend ‘beating’	bekommen ‘anxious’	bluten ‘to bleed’
GloVe 1	gemüt ‘mind’	mein ‘my’	seele ‘soul’	liebe ‘love’	brust ‘chest’
GloVe 2	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
GloVe 3	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
SVD_{PPMI} all	busen ‘bosom’	fühlen ‘to feel’	liebe ‘love’	schmerzen ‘pain’	menschenherz ‘human heart’

Reliability for *Herz* ‘heart’

Embedding Model	First Neighbor	Second Neighbor	Third Neighbor	Fourth Neighbor	Fifth Neighbor
SGNS 1	schmerzen ‘pain’	bekommen ‘anxious’	busen ‘bosom’	bluten ‘to bleed’	herzen ‘to caress’
SGNS 2	bluten ‘to bleed’	klopfend ‘beating’	busen ‘bosom’	bekommen ‘anxious’	herzen ‘to caress’
SGNS 3	herzen ‘to caress’	busen ‘bosom’	klopfend ‘beating’	bekommen ‘anxious’	bluten ‘to bleed’
GloVe 1	gemüt ‘mind’	mein ‘my’	seele ‘soul’	liebe ‘love’	brust ‘chest’
GloVe 2	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
GloVe 3	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
SVD_{PPMI} all	busen ‘bosom’	fühlen ‘to feel’	liebe ‘love’	schmerzen ‘pain’	menschenherz ‘human heart’

Reliability for *Herz* ‘heart’

Embedding Model	First Neighbor	Second Neighbor	Third Neighbor	Fourth Neighbor	Fifth Neighbor
SGNS 1	schmerzen ‘pain’	bekommen ‘anxious’	busen ‘bosom’	bluten ‘to bleed’	herzen ‘to caress’
SGNS 2	bluten ‘to bleed’	klopfend ‘beating’	busen ‘bosom’	bekommen ‘anxious’	herzen ‘to caress’
SGNS 3	herzen ‘to caress’	busen ‘bosom’	klopfend ‘beating’	bekommen ‘anxious’	bluten ‘to bleed’
GloVe 1	gemüt ‘mind’	mein ‘my’	seele ‘soul’	liebe ‘love’	brust ‘chest’
GloVe 2	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
GloVe 3	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
SVD_{PPMI} all	busen ‘bosom’	fühlen ‘to feel’	liebe ‘love’	schmerzen ‘pain’	menschenherz ‘human heart’

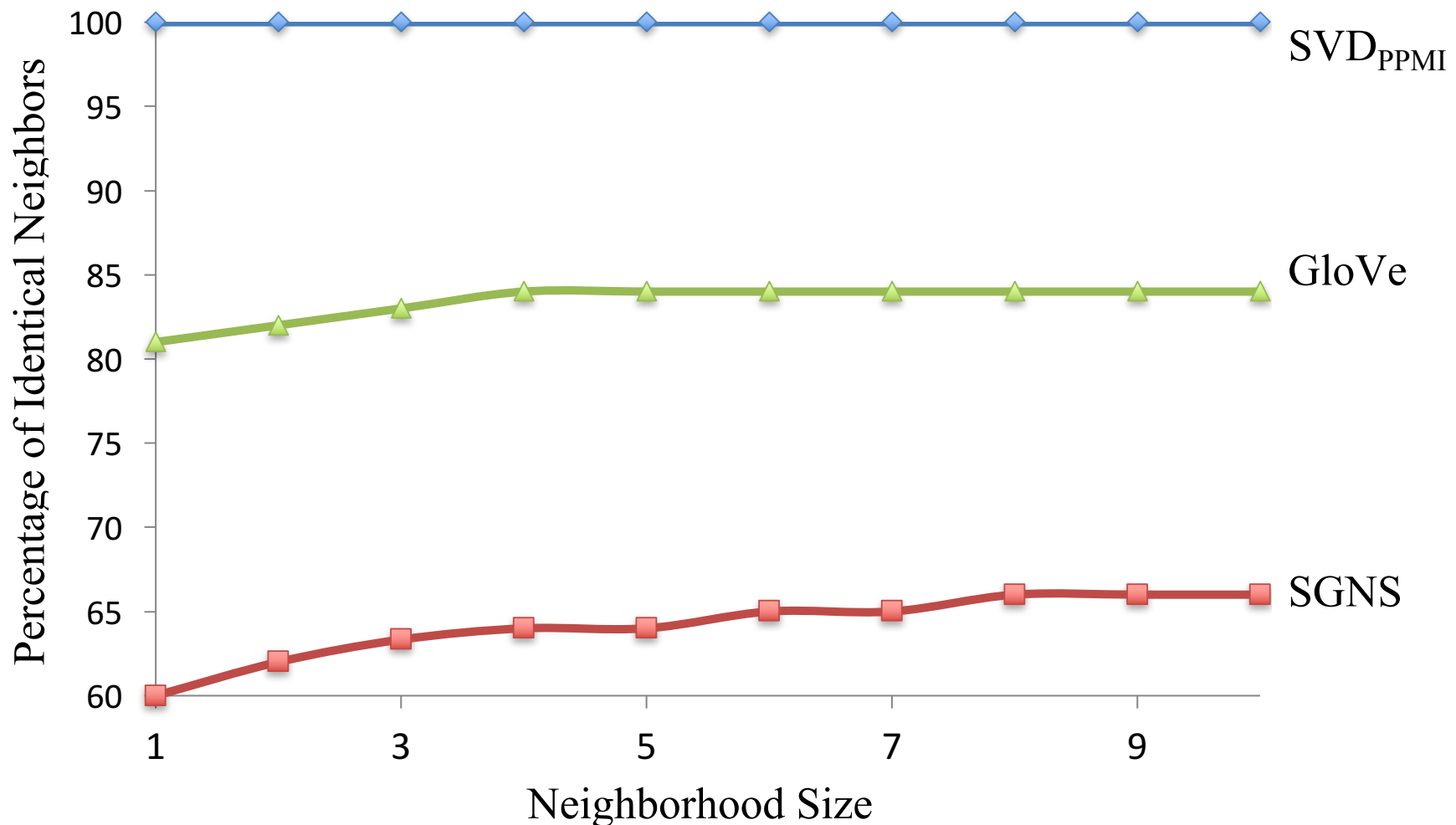
Reliability for *Herz* ‘heart’

Embedding Model	First Neighbor	Second Neighbor	Third Neighbor	Fourth Neighbor	Fifth Neighbor
SGNS 1	schmerzen ‘pain’	bekommen ‘anxious’	busen ‘bosom’	bluten ‘to bleed’	herzen ‘to caress’
SGNS 2	bluten ‘to bleed’	klopfend ‘beating’	busen ‘bosom’	bekommen ‘anxious’	herzen ‘to caress’
SGNS 3	herzen ‘to caress’	busen ‘bosom’	klopfend ‘beating’	bekommen ‘anxious’	bluten ‘to bleed’
GloVe 1	gemüt ‘mind’	mein ‘my’	seele ‘soul’	liebe ‘love’	brust ‘chest’
GloVe 2	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
GloVe 3	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
SVD_{PPMI} all	busen ‘bosom’	fühlen ‘to feel’	liebe ‘love’	schmerzen ‘pain’	menschenherz ‘human heart’

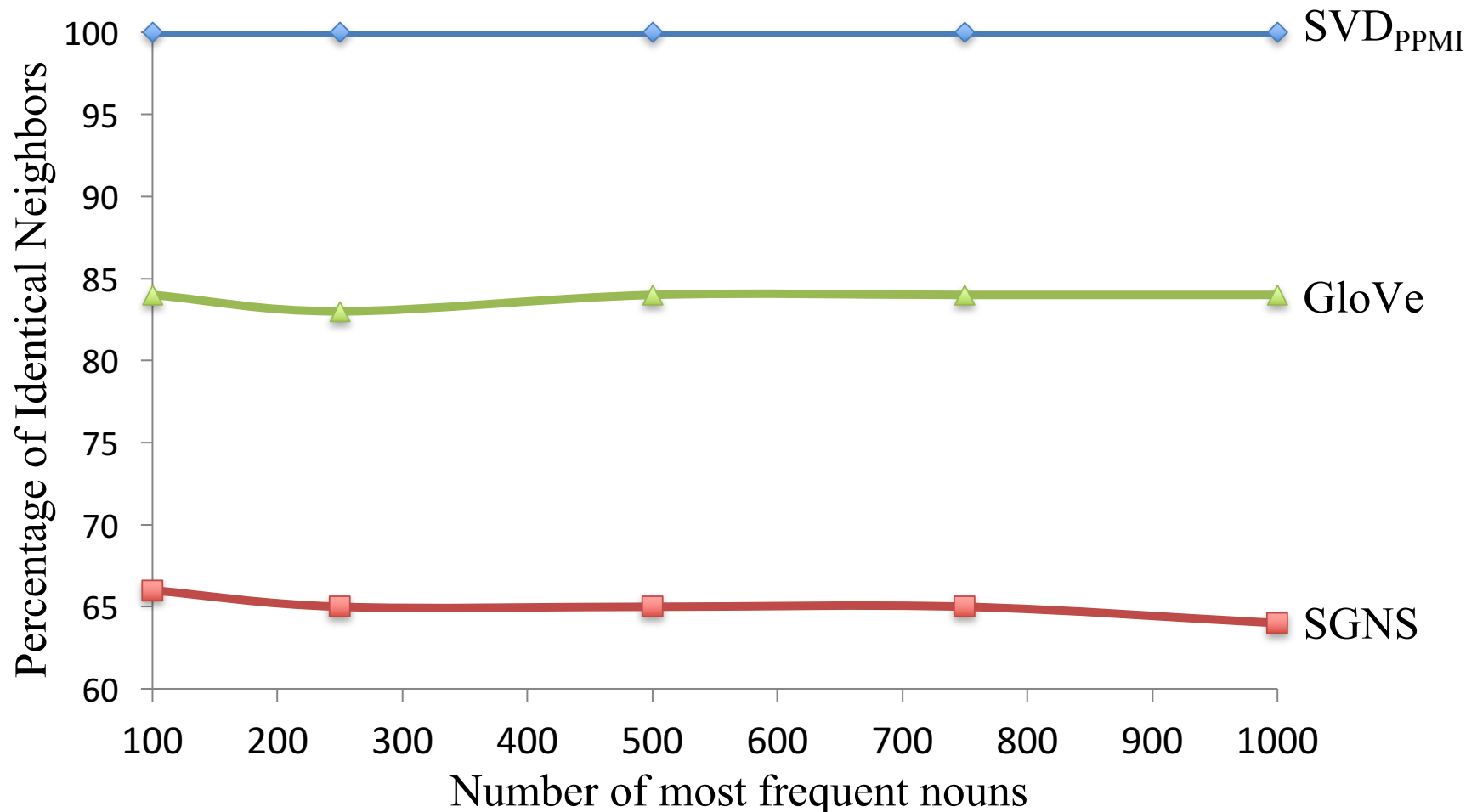
Reliability for *Herz* ‘heart’

Embedding Model	First Neighbor	Second Neighbor	Third Neighbor	Fourth Neighbor	Fifth Neighbor
SGNS 1	schmerzen ‘pain’	bekommen ‘anxious’	busen ‘bosom’	bluten ‘to bleed’	herzen ‘to caress’
SGNS 2	bluten ‘to bleed’	klopfend ‘beating’	busen ‘bosom’	bekommen ‘anxious’	herzen ‘to caress’
SGNS 3	herzen ‘to caress’	busen ‘bosom’	klopfend ‘beating’	bekommen ‘anxious’	bluten ‘to bleed’
GloVe 1	gemüt ‘mind’	mein ‘my’	seele ‘soul’	liebe ‘love’	brust ‘chest’
GloVe 2	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
GloVe 3	gemüt ‘mind’	mein ‘my’	seele ‘soul’	brust ‘chest’	liebe ‘love’
SVD_{PPMI} all	busen ‘bosom’	fühlen ‘to feel’	liebe ‘love’	schmerzen ‘pain’	menschenherz ‘human heart’

Reliability for 1000 most frequent nouns depending on neighborhood size



Reliability for 100–1000 most frequent nouns depending on word frequency



Conclusion

- Neural word embeddings are unreliable
- SVD_{PPMI} is reliable and performs very similar on evaluation tasks
- Also think about: Preprocessing often includes random sampling

Accessible SVD_{PPMI} embeddings for diachronic linguistics

Welcome to JeSemE

The Jena Semantic Explorer

☒ COHA ☐ DTA ☐ GB Fiction ☐ GB German ☐ RSC

JeSemE allows you to explore the semantic development of words over time. An interesting example is searching "heart" in the COHA corpus.

<http://jeseme.org>

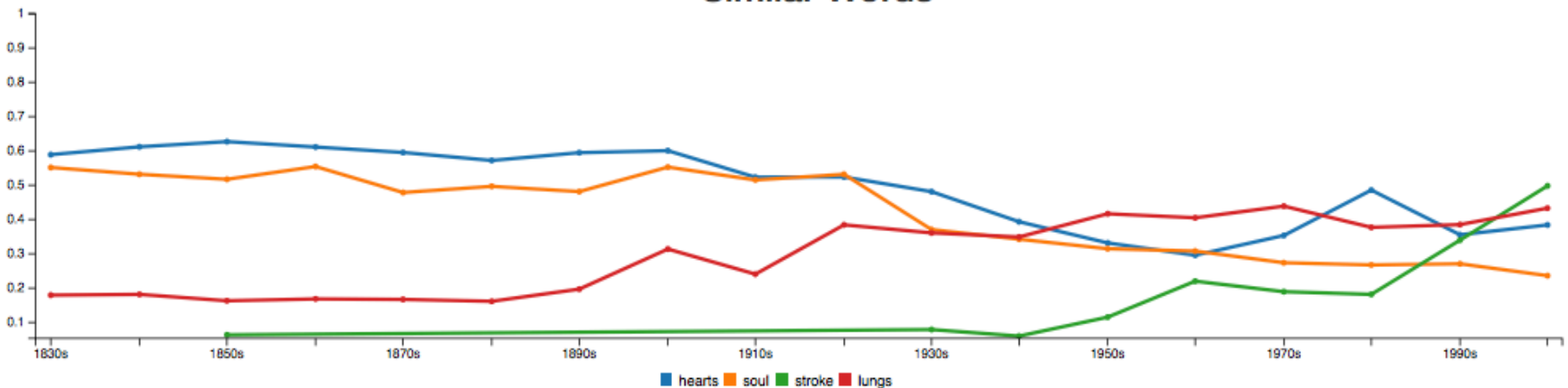
Hellrich & Hahn, ACL 2017

Accessible SVD_{PPMI} embeddings for diachronic linguistics

JeSemE - The Jena Semantic Explorer

Results for "heart" in Corpus of Historical American English
Note: lowercased
Search in [Corpus of Historical American English](#)

Similar Words



<http://jeseme.org>

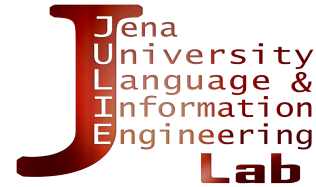
Hellrich & Hahn, ACL 2017



FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA



MODELL
ROMANTIK



Don't Get Fooled by Word Embeddings— Better Watch Their Neighborhood



Johannes Hellrich^{1,2}

1: Graduate School 'The Romantic Model',
Friedrich Schiller University Jena,
Jena, Germany

<http://www.modellromantik.uni-jena.de>



& Udo Hahn²

2: Jena University Language & Information
Engineering (JULIE) Lab
Friedrich Schiller University Jena,
Jena, Germany

<http://www.julielab.de>

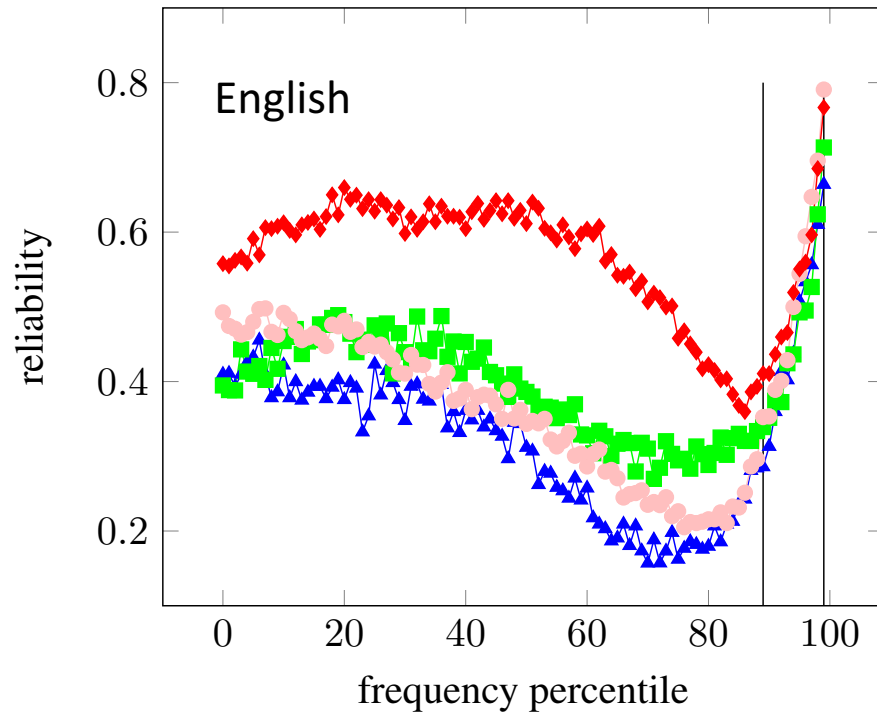
Word Embedding Performance

Method	WordSim Similarity	WordSim Relatedness	Bruni et al. MEN	Radinsky et al. M. Turk	Luong et al. Rare Words	Hill et al. SimLex	Google Add / Mul	MSR Add / Mul
PPMI	.755	.697	.745	.686	.462	.393	.553 / .679	.306 / .535
SVD	.793	.691	.778	.666	.514	.432	.554 / .591	.408 / .468
SGNS	.793	.685	.774	.693	.470	.438	.676 / .688	.618 / .645
GloVe	.725	.604	.729	.632	.403	.398	.569 / .596	.533 / .580

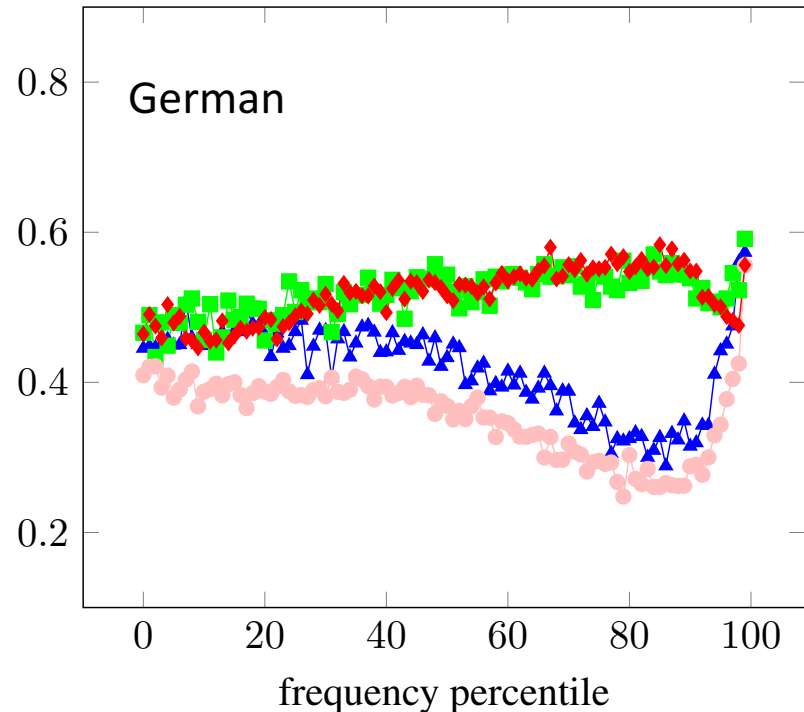
Table 4: Performance of each method across different tasks using the best configuration for that method and task combination, assuming $\text{win} = 2$.

From Levy et al. (2015)

Reliability of word2vec at different frequencies



—▲— SGHS 1900–1904 —■— SGNS 1900–1904
—●— SGHS 2005–2009 —◆— SGNS 2005–2009



—▲— SGHS 1900–1904 —■— SGNS 1900–1904
—●— SGHS 2005–2009 —◆— SGNS 2005–2009

- Hellrich&Hahn, COLING 2016
- word2vec models trained on Google Books corpora

Warning: Automatic word change research is focused on high frequency words

