

# Verbessertes Entity Linking mit neuronalen Word Embeddings

Bachelorarbeit

Verfasser: Robert Dziuba  
Betreuer: Prof. Dr. habil. Alexander Löser  
Gutachter: Prof. Dr.-Ing. Joachim Schimkat

# Problemstellung

# 3

## Das Wort - Mention



## Der Kontext und Word Embedding

President **Trump** is in the White House.

Hotel card games USA Tower  
administration conglomerat  
engineer President transport company

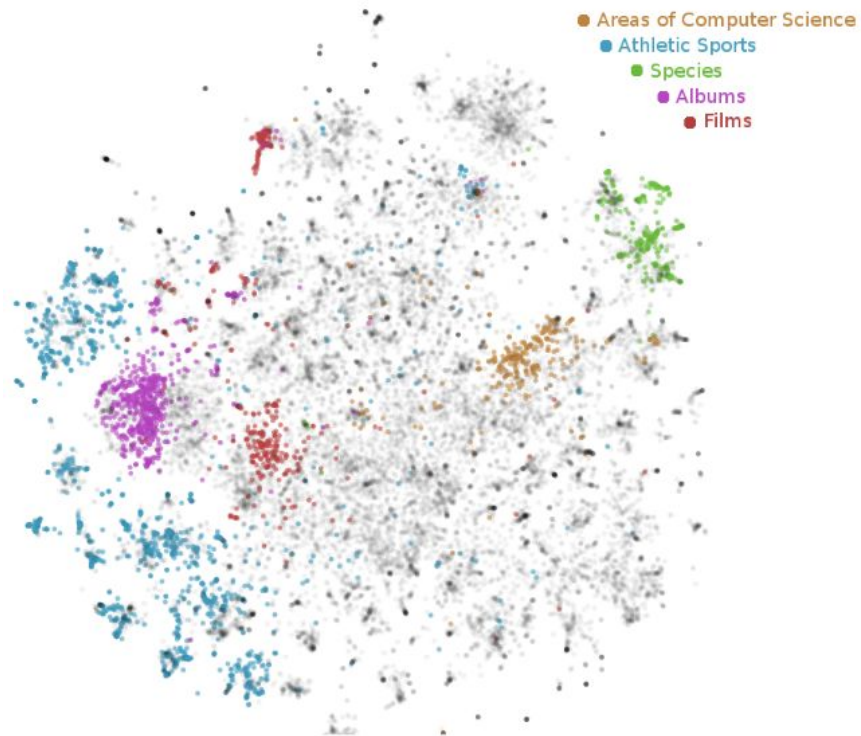
# Ziel der Arbeit

## Ziel der Arbeit

Das Ziel der Arbeit ist es, eine im Text genannten Mention, mit Hilfe ihres Kontextes, einer konkreten Entity zuzuordnen.

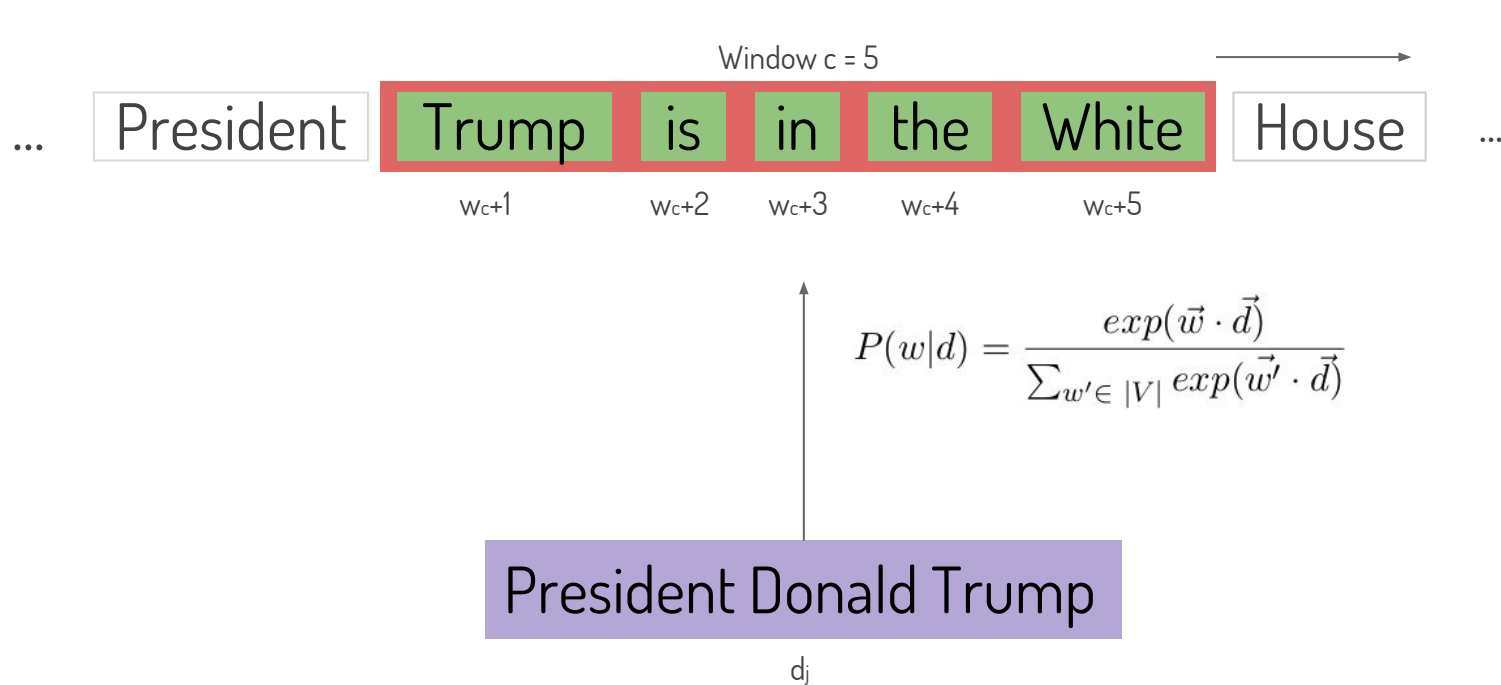
# Methodik

## Repräsentation von Kontext im Vektorraum - Entity Embeddings





## Paragraph Vector



# Disambiguierung im Vektorraum

**Mention** = "Trump"

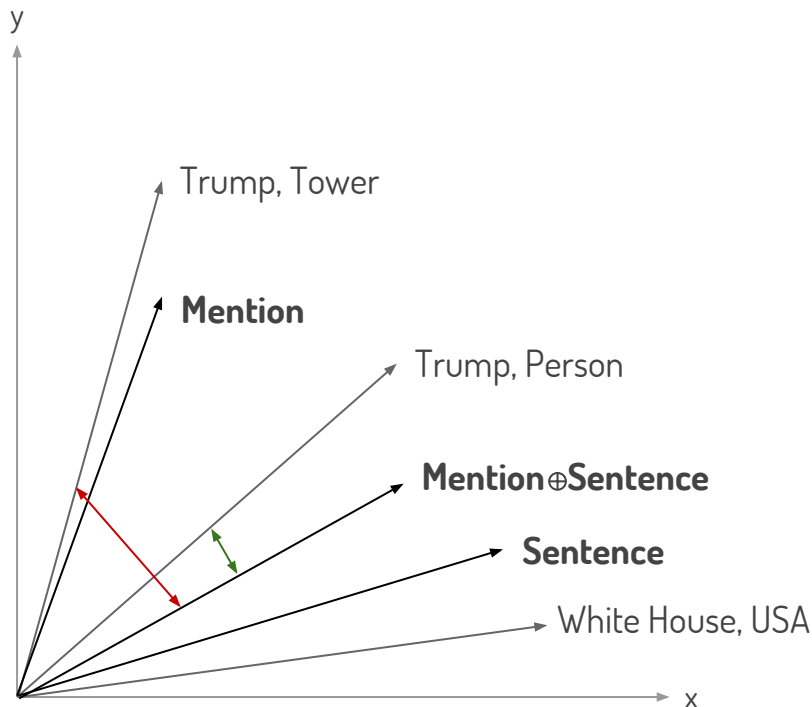
**Sentence** = "President Trump is in the White House."

$$\text{cosine}(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{|\vec{v}| |\vec{w}|} = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$

$$\hat{c} = \sum_{c \in |C_m|} \text{cosSim}(c)$$

mit  $\text{cosSim}(c) = \text{cosine}(\text{concat}(m, s), \text{concat}(c, c))$

und  $\text{concat}(a, b) = \text{vec}(a) \oplus \text{vec}(b)$



# Daten

## Wikipedia Korpus

Within American **political culture**, the **center-right** Republican Party is considered "conservative" and the **center-left** Democratic Party is considered "liberal".<sup>[375][376]</sup> The states of the **Northeast** and **West Coast** and some of the Great Lakes states, known as "**blue states**", are relatively liberal. The "**red states**" of the **South** and parts of the **Great Plains** and **Rocky Mountains** are relatively conservative.

Republican **Donald Trump**, the winner of the 2016 presidential election, is serving as the 45th President of the United States.<sup>[377]</sup> Leadership in the Senate includes Republican Vice President **Mike Pence**, Republican President Pro Tempore **Orrin Hatch**, **Majority Leader Mitch McConnell**, and Minority Leader **Chuck Schumer**.<sup>[378]</sup> Leadership in the House includes Speaker of the House **Paul Ryan**, **Majority Leader Kevin McCarthy**, and Minority Leader **Nancy Pelosi**.<sup>[379]</sup>

Wikipedia: Donald Trump - [https://en.wikipedia.org/wiki/United\\_States](https://en.wikipedia.org/wiki/United_States)

President Obama in 2017



**Donald Trump**  
45th **President**  
since January 20,  
2017

**Mike Pence**  
48th **Vice President**  
since January 20, 2017

## Trainingsdaten

### Beispielsatz

Q22686    Republican

Donald Trump, the winner of the 2016 presidential election, is serving as the 45th President of the United States.

### Datensatz

2.725.363 Entities

63.812.227 Sätzen

11 GB Gesamtumfang

# Evaluation

## Ergebnisse der Vektorkonkatenation-Strategie

Layer/Epochen	MRR	PREC (%)	REC (%)	F1 (%)
Lucene	0.534	48.24	45.10	46.60
	CommonPreprocessor			
100 Layer/10 Epochen	0.592	56.73	53.10	54.84
	MinimalLowercasePreprocessor			
100 Layer/10 Epochen	0.594	57.20	53.52	55.28
	CommonPreprocessor			
300 Layer/10 Epochen	<b>0.630</b>	<b>62.01</b>	<b>57.96</b>	<b>59.90</b>

# Fazit und Ausblick



## Fazit

- Mention Embeddings sind besser als Lucene Exact Match
- Konkatenierte Entity Embeddings sind besser als Lucene
- Mehr Layers führen zu besseren Ergebnissen
  
- Trainingszeit steigt überproportional mit Iteration/ Epochen
- Lucene limitiert das Suchergebnis
- Entities die nicht in der Wikipedia verlinkt sind werden nicht gelernt

## Ausblick

- Minimierung der Modellgröße
- Lucene beim Entity Linking weglassen
- Mehr Kontext
- Einfluss benachbarter Entities im Text nutzen

# Live Präsentation

Vielen Dank