

# Statistical Semantic Change Detection via Usage Similarities

Taichi Aida<sup>1</sup>, Daichi Mochihashi<sup>2</sup>, Hiroya Takamura<sup>3</sup>,  
Toshinobu Ogiso<sup>4</sup>, Mamoru Komachi<sup>1</sup>

<sup>1</sup>Hitotsubashi University

<sup>2</sup>The Institute of Statistical Mathematics

<sup>3</sup>National Institute of Advanced Industrial Science and Technology

<sup>4</sup>National Institute for Japanese Language and Linguistics





# Lexical Semantic Change

The meaning of words naturally evolve across different periods/domains.

- *plane* in 19th:

- “If a plane() be parallel to the horizontal...”
- “The sun is in the same plane() as the picture...”





- *plane* in 20th:

- “The President’s plane() landed at Goose Bay...”
- “The plane() kept climbing and climbing...”



# Lexical Semantic Change **Detection**

**Detecting semantically changed words** across different periods/domains.

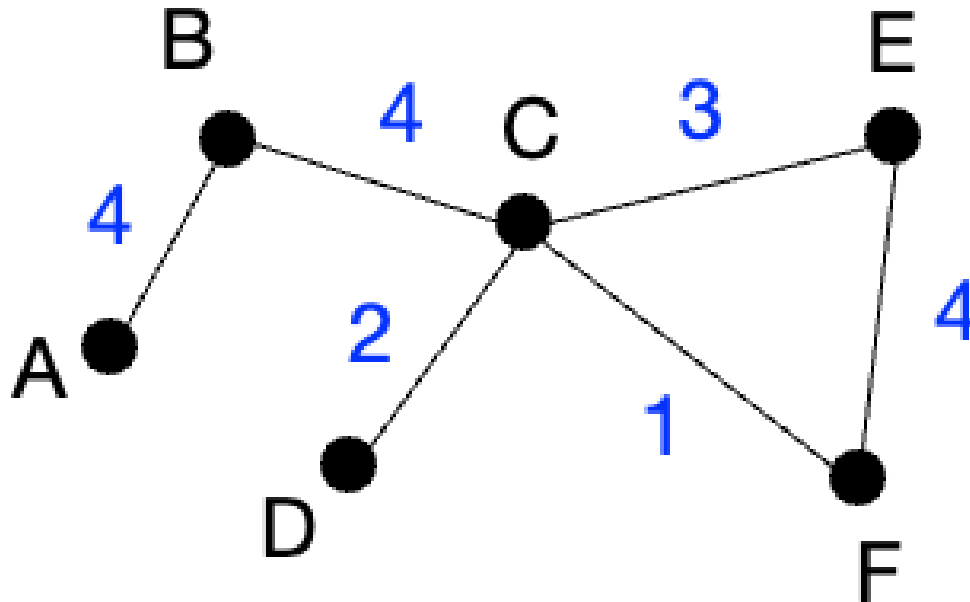
- *plane* in 19th:
  - “If a plane() be parallel to the horizontal...”
  - “The sun is in the same plane() as the picture...”
- *plane* in 20th:
  - “The President’s plane() landed at Goose Bay...”
  - “The plane() kept climbing and climbing...”

Key problem: **preparing annotated words**



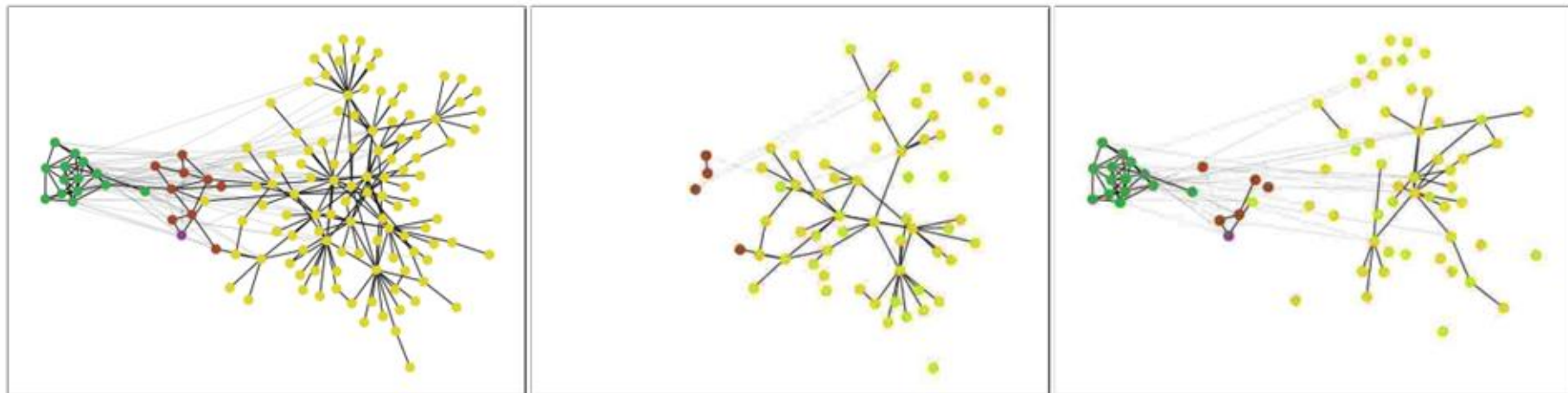
# DURel Annotation Pipeline [Schlechtweg+2018]

- **Annotating word usage** (A-F) via usage similarity (1: Unrelated - 4: Identical)
- It can be viewed as an annotation graph = **word usage graph**



# Word Usage Graph [Schlechtweg+2021]

- These graphs are clustered to identify senses
- **The degree of sense gain/loss** between different periods/domains  
= **The degree of semantic change**

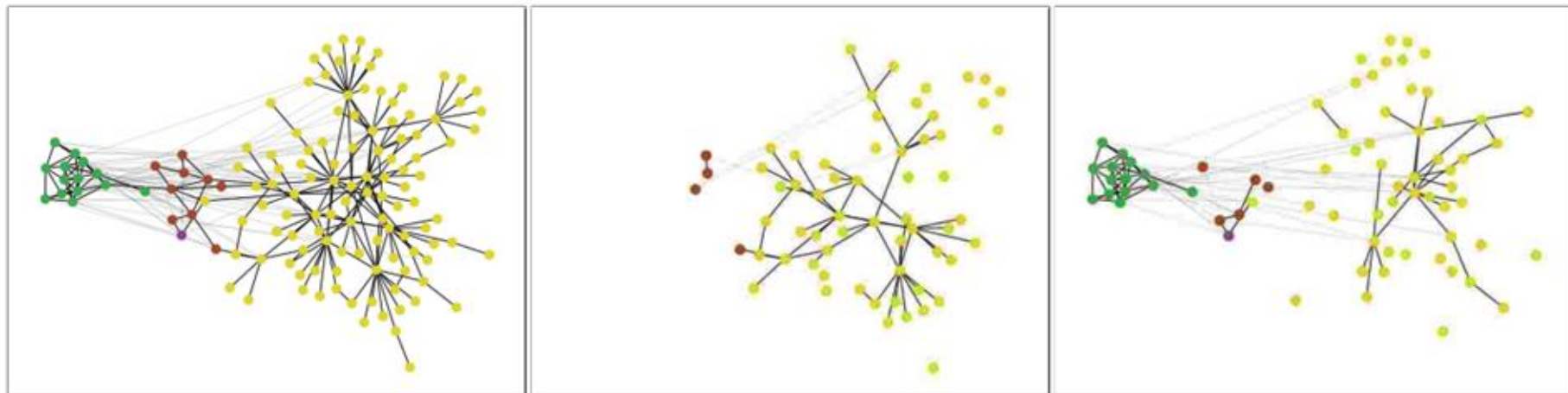


# Word Usage Graph [Schlechtweg+2021]

-

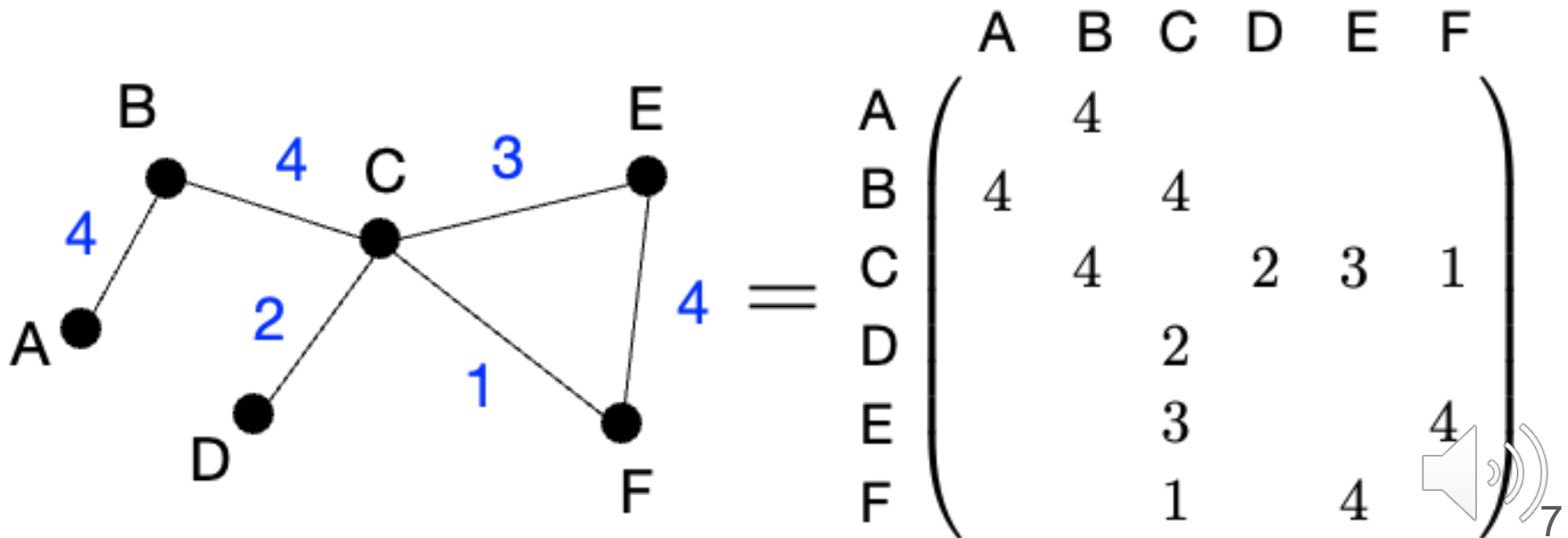
-

How can we model semantic change from word usage graphs directly?



# Proposal

- Represent the word usage graph as an adjacency matrix (A-C: Doc1, D-F: Doc2)
- **Model the usage similarity matrix directly**



# Proposal


- If the meaning of the word is **stable**
- All similarity scores **X** are generated from **the same distribution**

$$\begin{array}{c} \text{A} \\ \text{B} \\ \text{C} \\ \text{D} \\ \text{E} \\ \text{F} \end{array} \begin{pmatrix} \text{A} & \text{B} & \text{C} & \text{D} & \text{E} & \text{F} \\ 4 & 2 & 3 & 4 & 1 & 2 \\ 2 & 4 & 3 & 2 & 4 & 4 \\ 3 & 3 & 4 & 2 & 3 & 1 \\ 4 & 2 & 2 & 4 & 3 & 2 \\ 1 & 4 & 3 & 3 & 4 & 1 \\ 2 & 4 & 1 & 2 & 1 & 4 \end{pmatrix} = \mathbf{X}$$



# Proposal

- If the meaning of the word is **changed**
- All similarity scores **X1~X4** are generated from **the distinct distributions**

$$\begin{array}{c} \text{A} \\ \text{B} \\ \text{C} \\ \text{D} \\ \text{E} \\ \text{F} \end{array} \begin{array}{c} \text{A} \quad \text{B} \quad \text{C} \quad \text{D} \quad \text{E} \quad \text{F} \\ \left( \begin{array}{ccc|ccc} 4 & 4 & 4 & 1 & 1 & 1 \\ 4 & 4 & 4 & 1 & 1 & 1 \\ 4 & 4 & 4 & 1 & 1 & 1 \\ \hline 1 & 1 & 1 & 4 & 4 & 4 \\ 1 & 1 & 1 & 4 & 4 & 4 \\ 1 & 1 & 1 & 4 & 4 & 4 \end{array} \right)$$
$$= \begin{array}{c} \text{A} \\ \text{B} \\ \text{C} \\ \text{D} \\ \text{E} \\ \text{F} \end{array} \begin{array}{c} \text{A} \quad \text{B} \quad \text{C} \quad \text{D} \quad \text{E} \quad \text{F} \\ \left( \begin{array}{ccc|ccc} & & & & & \\ & \mathbf{X}_1 & & & \mathbf{X}_2 & \\ & & & & & \\ \hline & & & & & \\ & \mathbf{X}_3 & & & \mathbf{X}_4 & \\ & & & & & \end{array} \right)$$


# Proposal

- Our goal: estimate  $p(\theta|\mathbf{X})$  from similarity matrix  $\mathbf{X}$  ( $\theta \in \{0, 1\}$ : semantic change indicator)

$$\frac{p(\theta|\mathbf{X}) \propto p(\mathbf{X}|\theta)p(\theta) \propto \begin{cases} p(\mathbf{X}|\theta=0) \\ p(\mathbf{X}|\theta=1) \end{cases}}{\text{Bayes' theorem}}$$

- To predict semantic change, we can compare

$$\left\{ \begin{array}{l} p(\mathbf{X}) = \int p(\mathbf{X}|\mathbf{p})p(\mathbf{p})d\mathbf{p} = \frac{\Gamma(\sum_k \alpha_k)}{\Gamma(L + \sum_k \alpha_k)} \prod_{k=1}^K \frac{\Gamma(\alpha_k + n_k)}{\Gamma(\alpha_k)} \\ p(\mathbf{X}|\theta=0) = p(\mathbf{X}), \quad p(\mathbf{X}|\theta=1) = \prod_{n=1}^4 p(\mathbf{X}_n) \end{array} \right. \quad \text{Pólya distribution}$$



# Results

- Our method **outperforms** previous SotA approaches using word embeddings

Method	EN	DE	LA	SV
SGNS [11]	73.0	54.2	45.0	61.3
SGNS [12]	62.2	75.0	<b>70.0</b>	67.7
BERT [13]	70.3	75.0	55.0	74.2
Pólya (ours)	<b>76.1</b>	<b>80.0</b>	N/A	<b>88.6</b>



# Results

Data	Language	Grouping 1	Grouping 2	Accuracy (%)	
				MostFreq	Pólya
DWUG EN	EN	1810–1860	1960–2010	54.3	<b>76.1</b>
DWUG EN Resampled	EN	1810–1860	1960–2010	60.0	<b>80.0</b>
DWUG DE	DE	1800–1899	1946–1990	60.0	<b>80.0</b>
DWUG DE Resampled	DE	1800–1899	1946–1990	60.0	<b>73.3</b>
DiscoWUG	DE	1800–1899	1946–1990	51.0	<b>72.0</b>
RefWUG	DE	1750–1800	1850–1900	<b>54.5</b>	45.5
DURel	DE	1750–1800	1850–1900	<b>63.6</b>	<b>63.6</b>
SURel	DE	general	domain specific	63.6	<b>68.2</b>
RuSemShift 1	RU	1682–1916	1918–1990	<b>77.5</b>	<b>77.5</b>
RuSemShift 2	RU	1918–1990	1991–2016	<b>62.3</b>	<b>62.3</b>
RuShiftEval 1	RU	1700–1916	1918–1990	<b>74.8</b>	<b>74.8</b>
RuShiftEval 2	RU	1918–1990	1992–2016	<b>70.3</b>	<b>70.3</b>
RuShiftEval 3	RU	1700–1916	1992–2016	<b>68.5</b>	<b>68.5</b>
DWUG ES	ES	1810–1906	1994–2020	55.5	<b>78.0</b>
DiaWUG	ES	Spanish variant 1	Spanish variant 2	65.6	<b>81.3</b>
DWUG SV	SV	1790–1830	1895–1903	68.1	<b>88.6</b>
DWUG SV Resampled	SV	1790–1830	1895–1903	60.0	<b>73.3</b>
ChiWUG	ZH	1954–1978	1979–2003	<b>57.5</b>	52.5
DWUG IT	IT	1948–1970	1990–2014	<b>69.2</b>	N/A
DWUG LA	LA	–200–0	0–2000	<b>55.5</b>	N/A
NorDiaChange 1	NO	1929–1965	1970–2013	67.5	<b>75.0</b>
NorDiaChange 2	NO	1980–1990	2012–2019	<b>77.5</b>	70.0



# Discussion

- Only 5→2% of the total entries work!

