



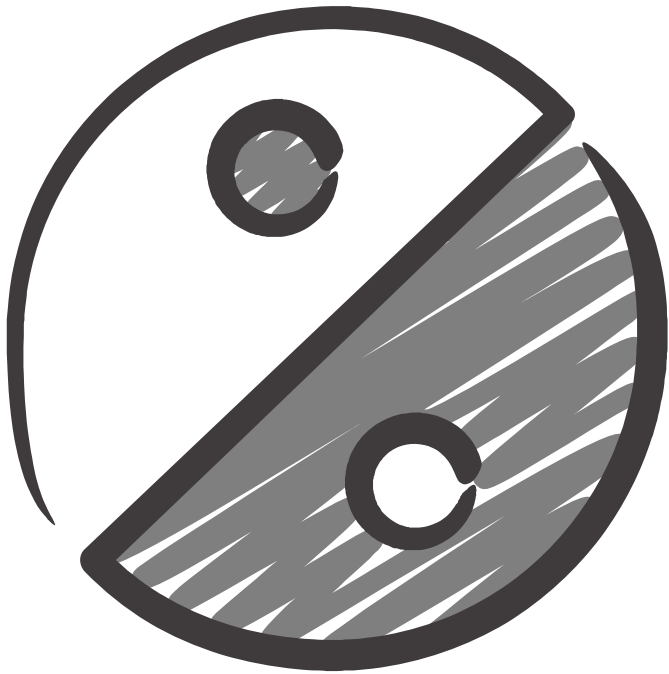
# Computational methods for lexical semantic change

Nina Tahmasebi, PhD

University of Gothenburg

Helsinki, Finland, Feb. 18, 2019

# Outline



Computational  
methods for LSC

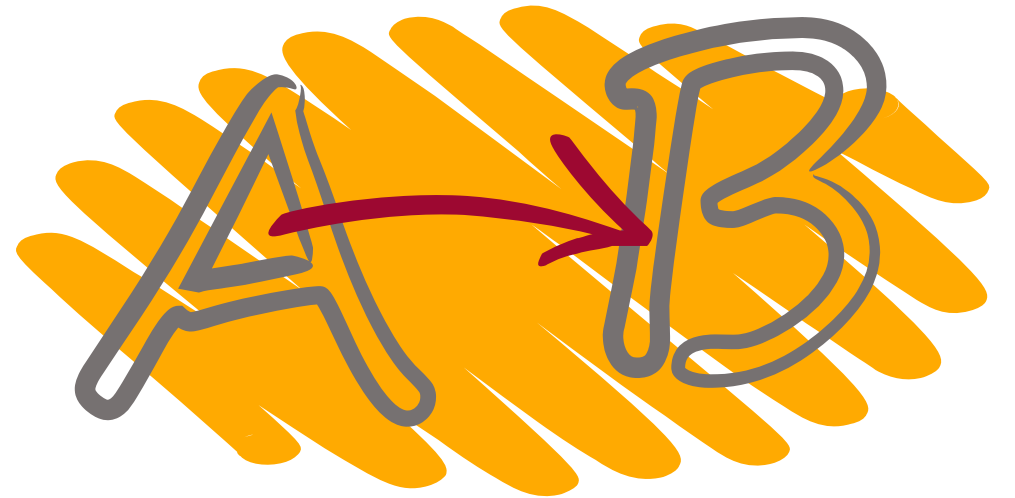


Historical Linguistic  
perspective




Evaluation

# Computational lexical semantic change




# LiWA – Living Web Archives

---

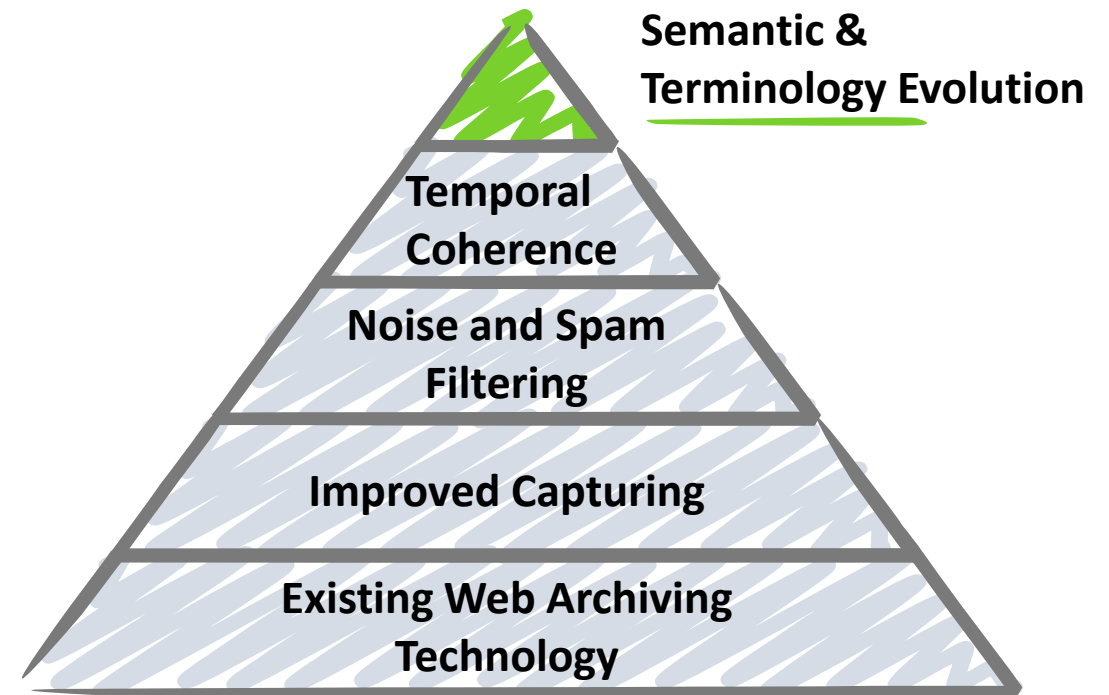


dealing with  
terminology  
evolution



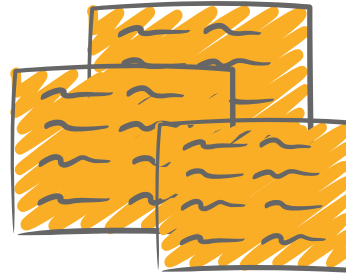
preparing for  
evolution aware  
access support

---



# Increasing amount of historical texts in digital format

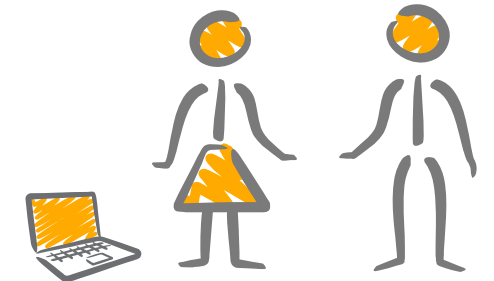
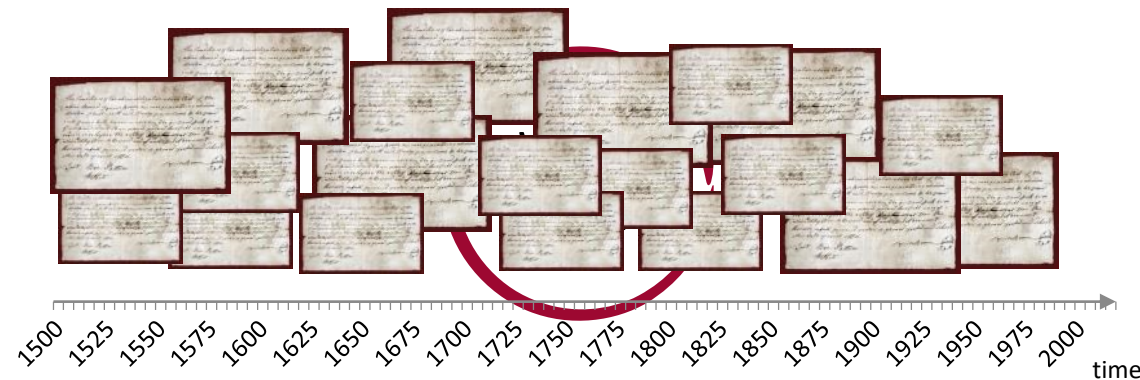
Easy digital access for anyone!  
**Not only scholars.**



Possibility to **digitally analyze**  
historical documents  
at **large scale.**

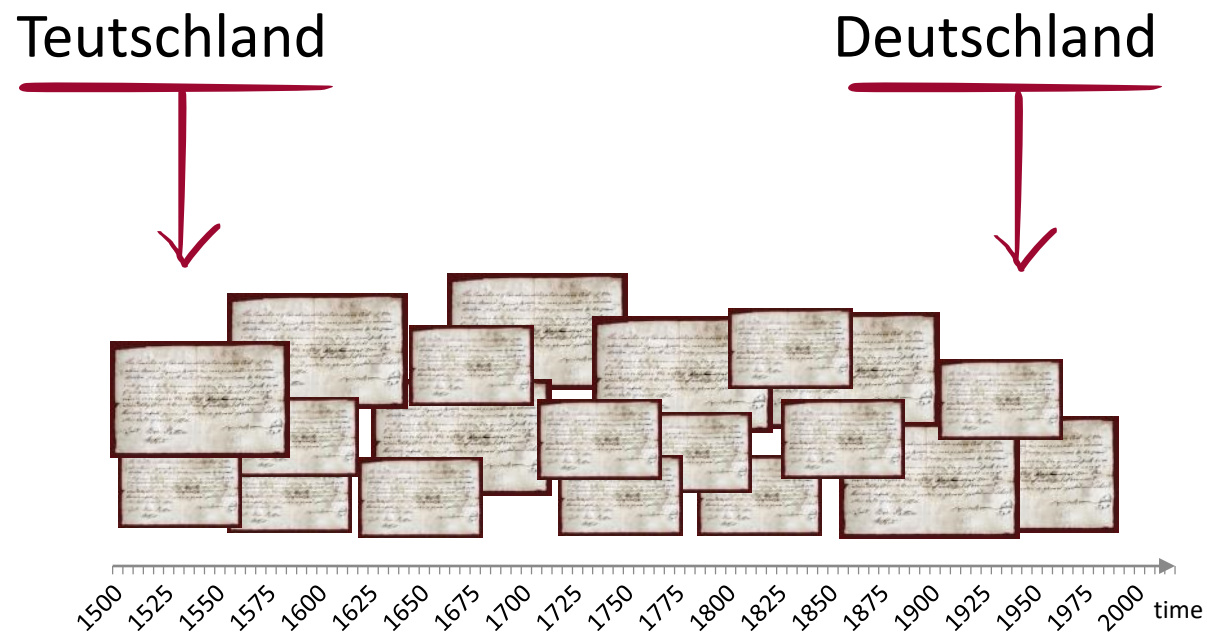
Information from primary sources  
**Not only modern interpretations.**

**Text-based  
Digital Humanities**

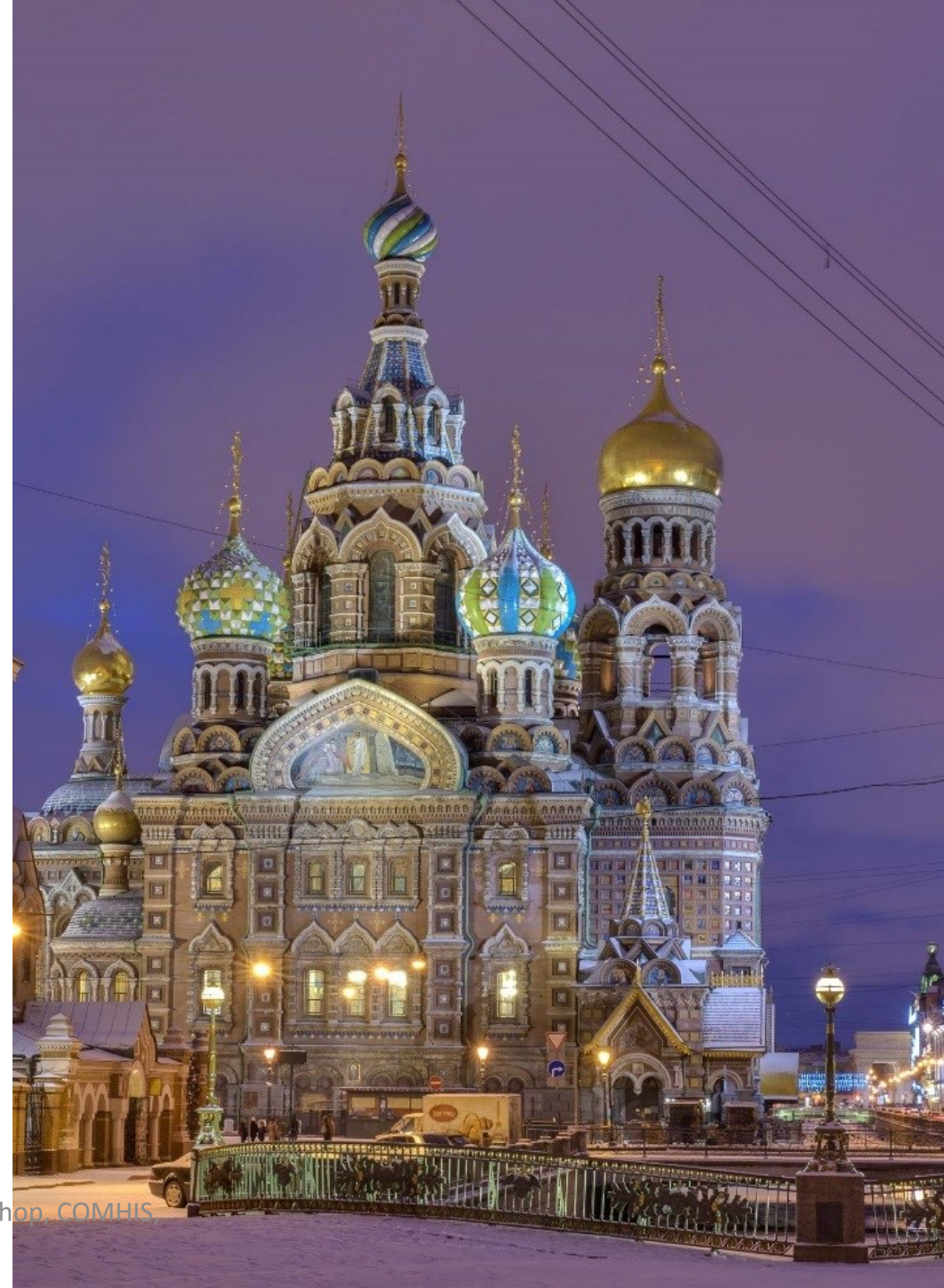
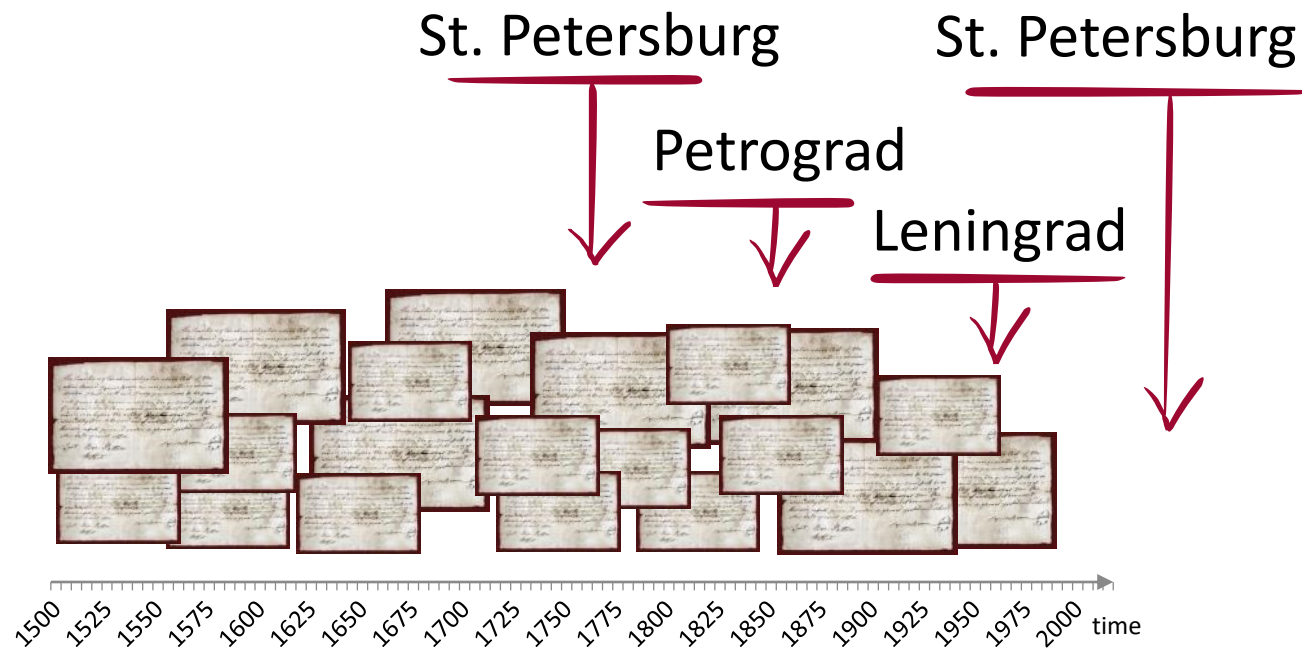




# Spelling change



# Lexical replacement: Named entity change

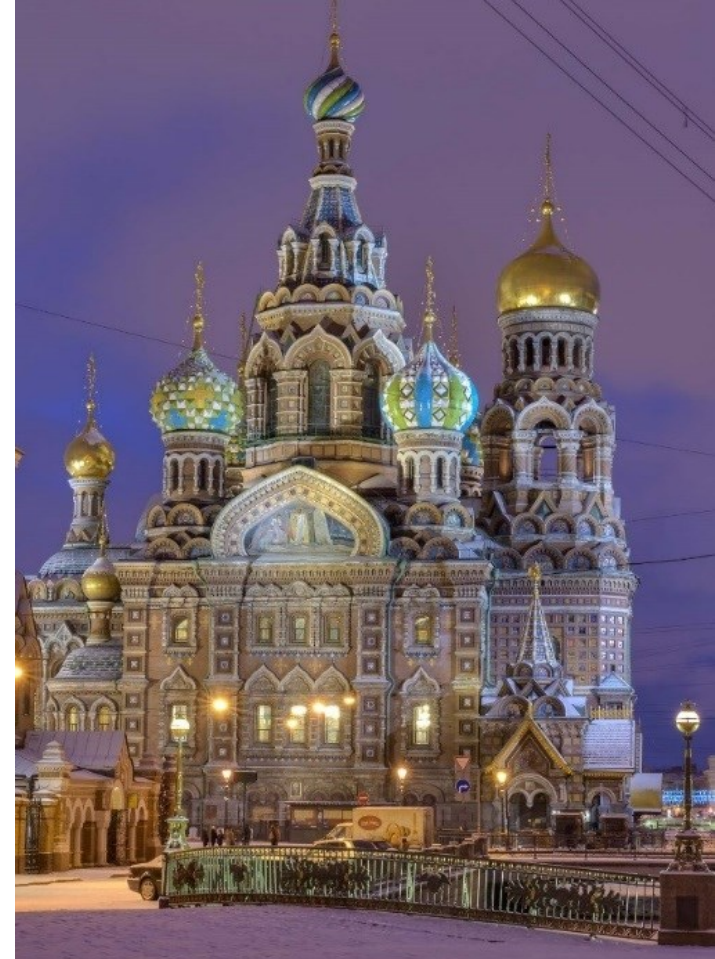
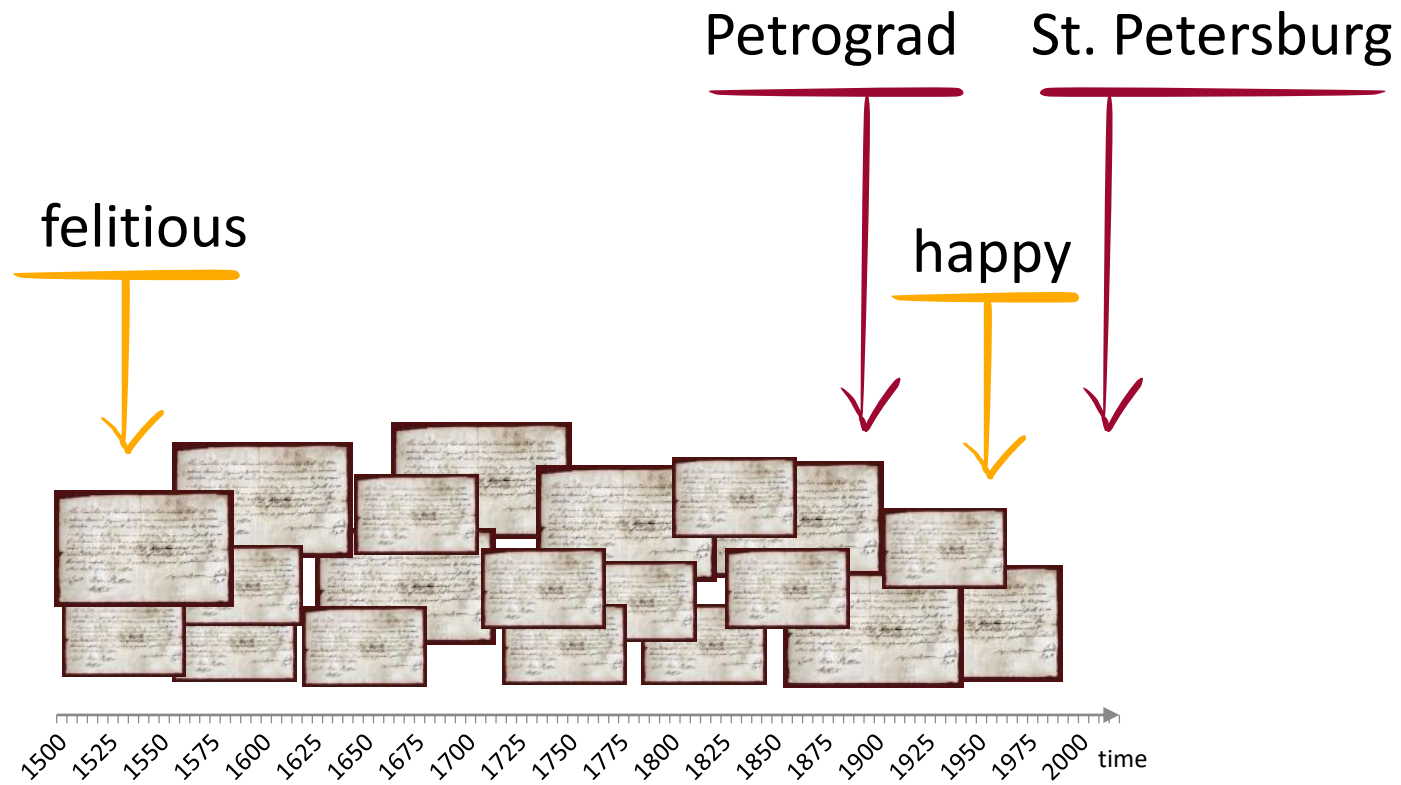








# Lexical replacement:





# awesome

He was an  
**awesome** leader!



He was an  
**awesome** leader!



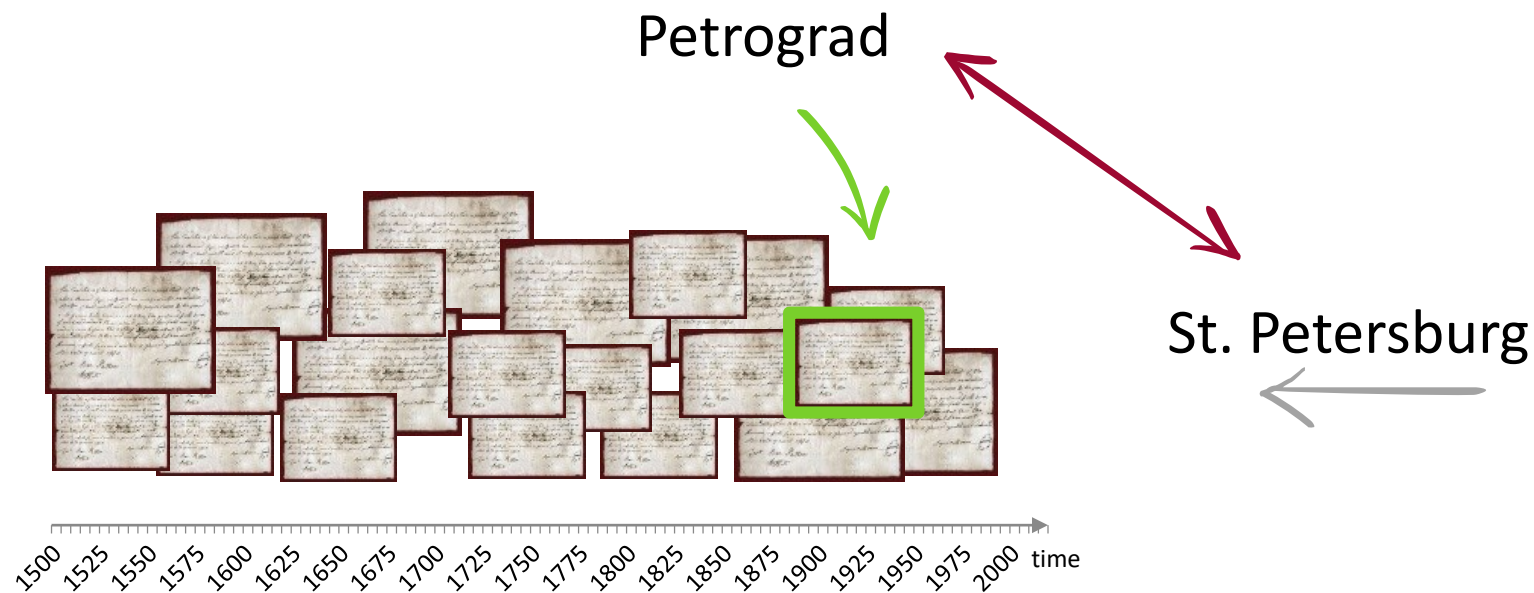
time





# What is the problem?

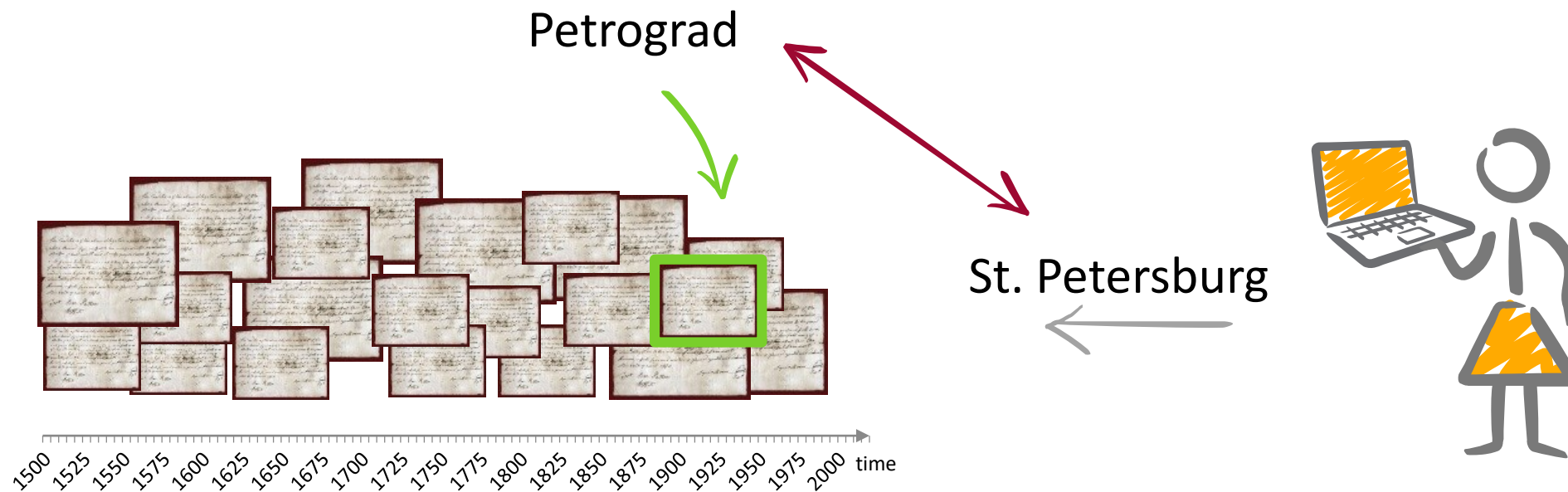
Finding



# What is the problem?

Finding

Interpreting



” Sebastini’s benefit last night at the  
Opera House was overflowing with  
the fashionable and **gay** ”









” Sebastini’s benefit last night at the Opera House was overflowing with the fashionable and **gay** ”

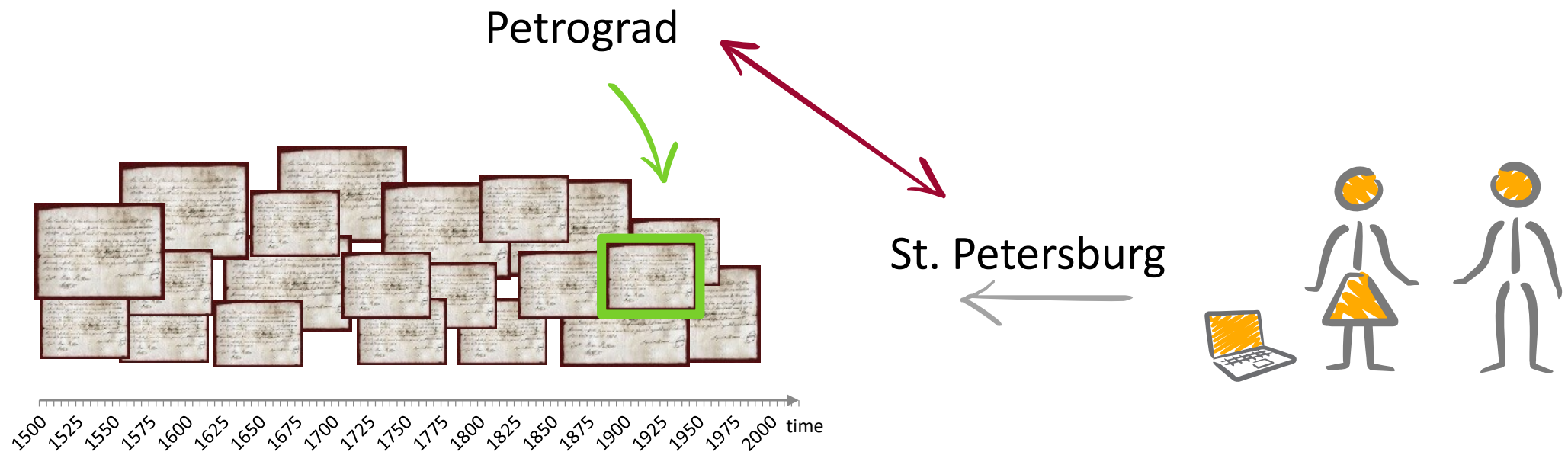
The Times, April 27th, 1787

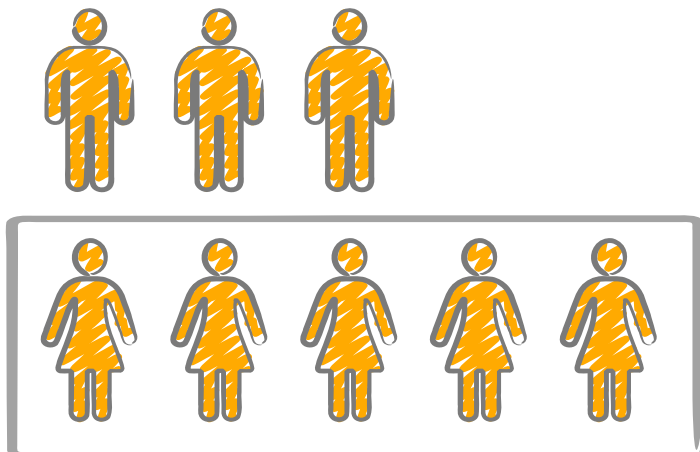


# What is the problem?

Finding

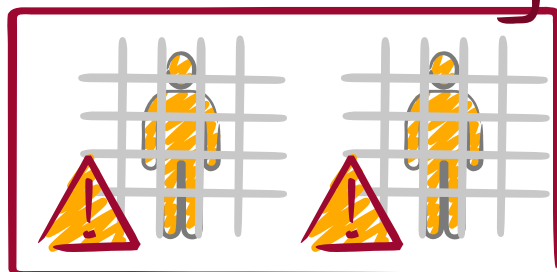
Interpreting



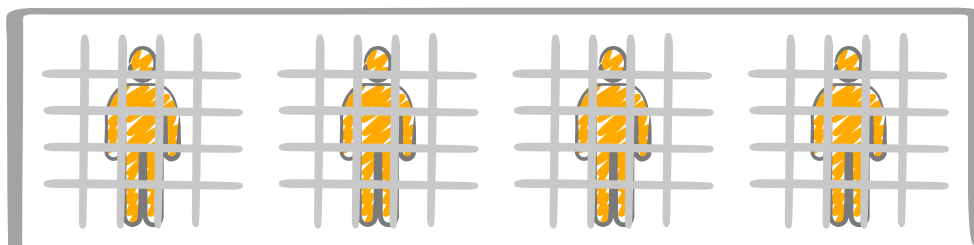


← girl

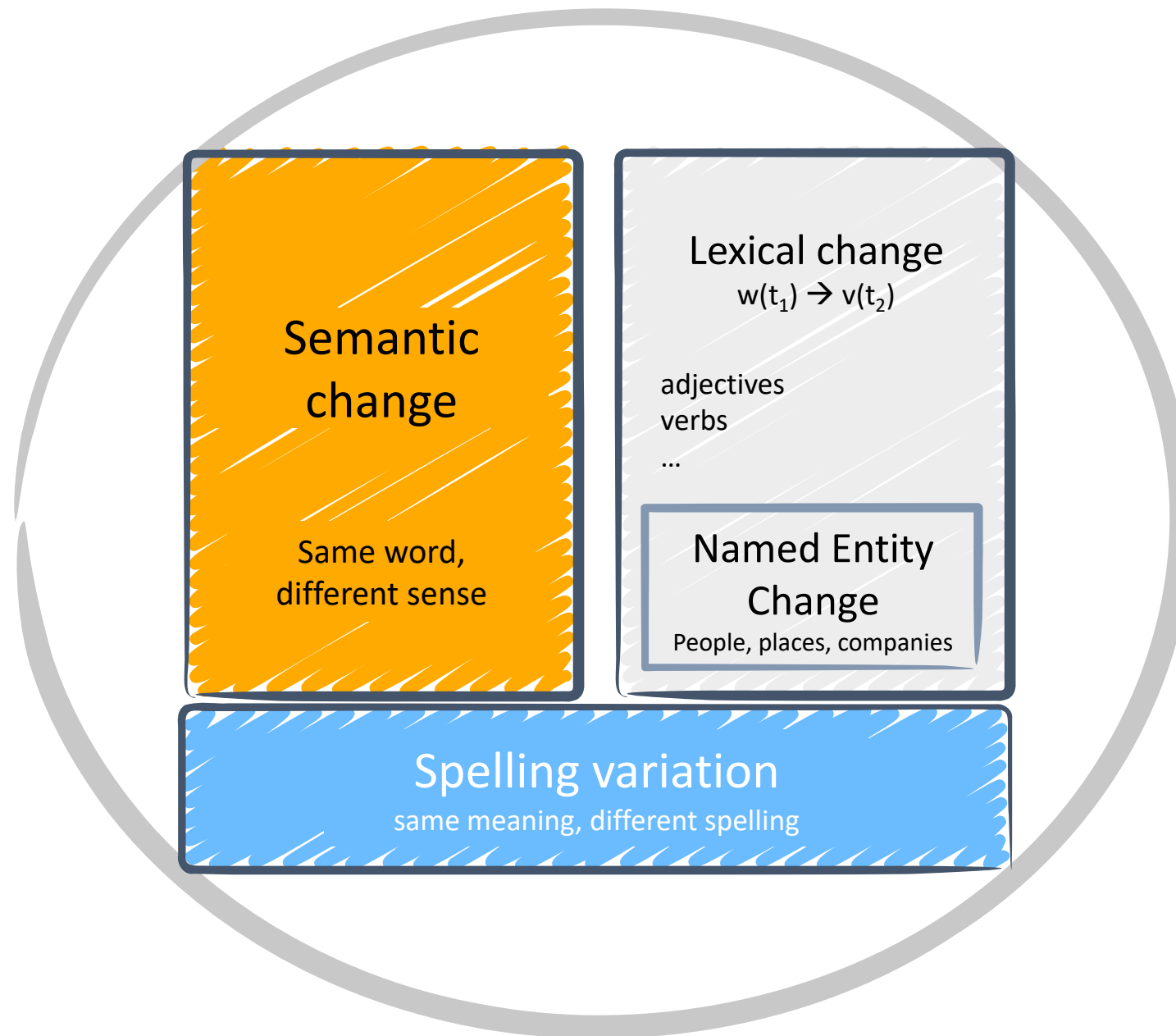
Wolf 'varg'



← criminal





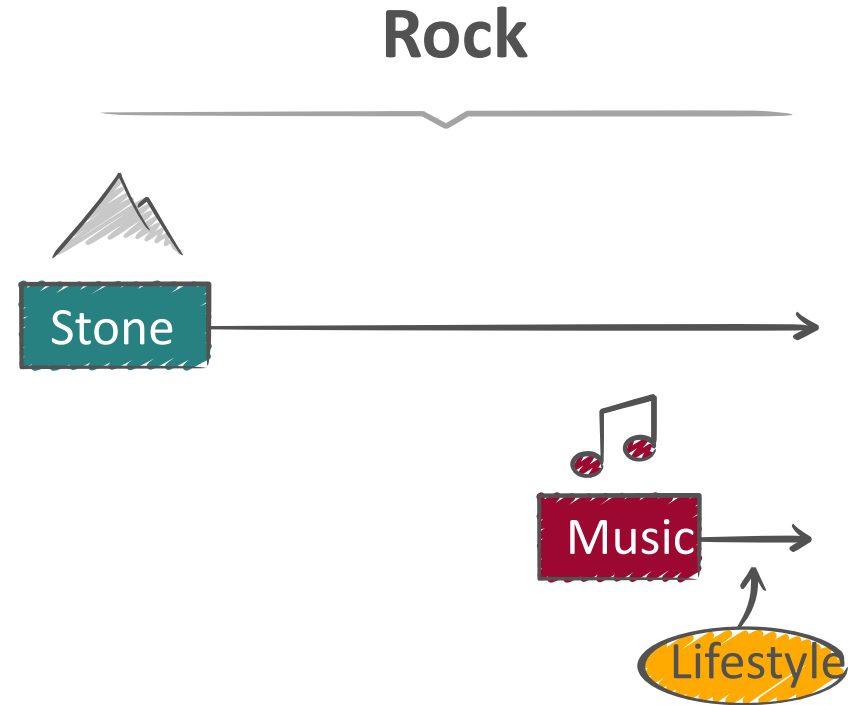


# Aims

To find word sense changes  
**automatically** by

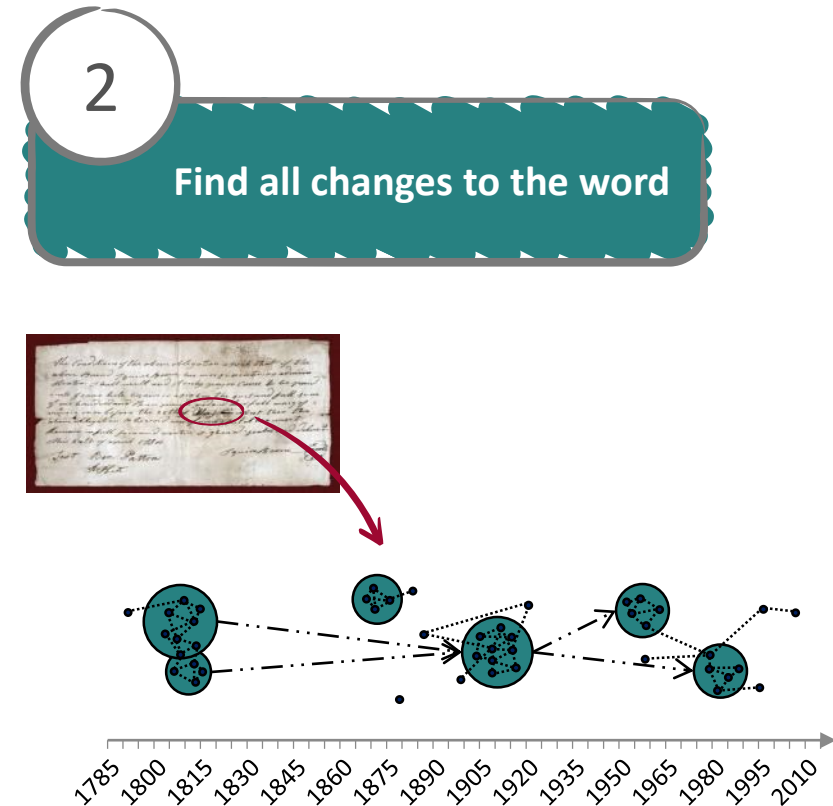
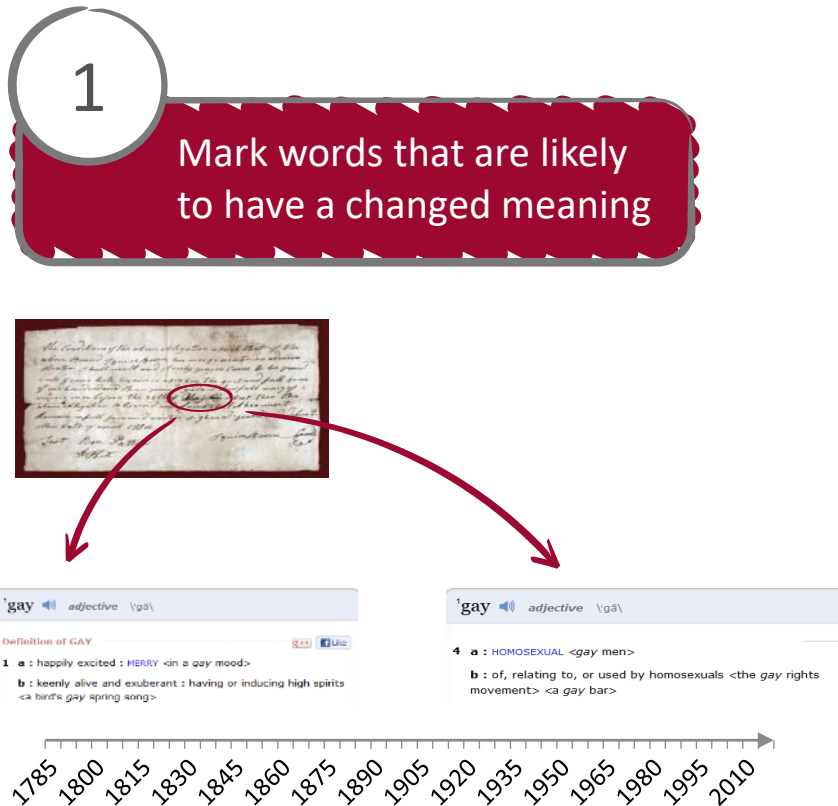
- 1 Modeling word senses
- 2 Comparing these over time

To find **what** changes, **how** it changed and **when** it changed



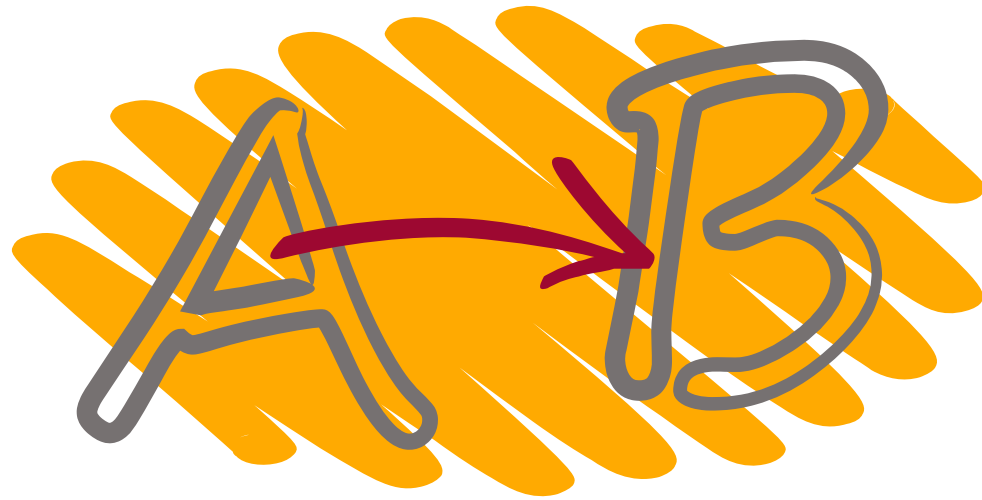
# Vision

Given a word in a document at time  $t$



# Lexical semantic change

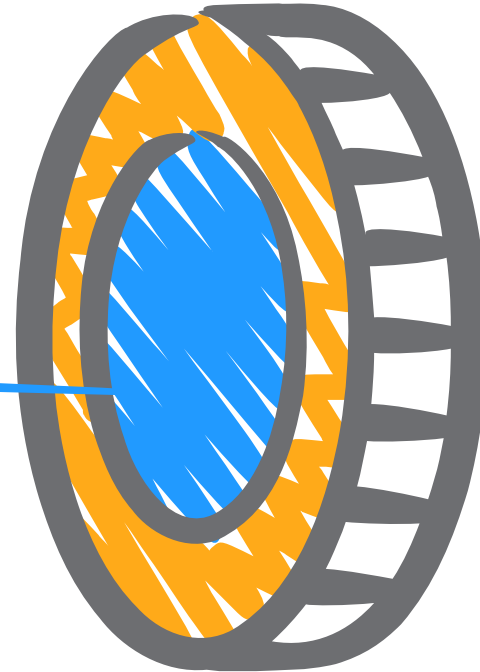
The (historical) linguistic perspective





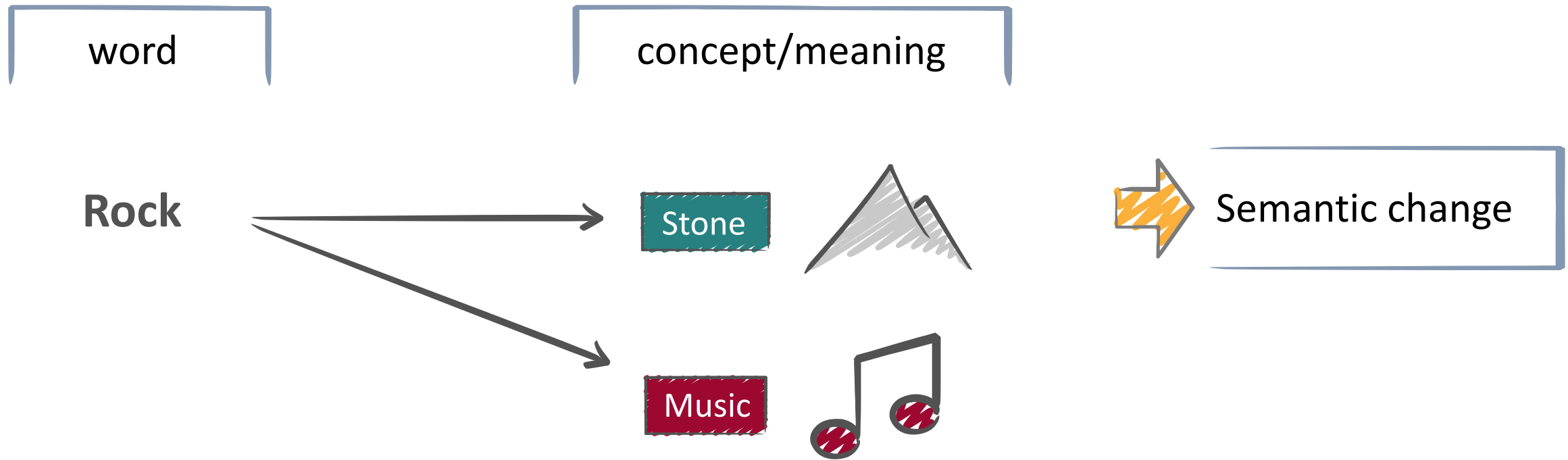


Semasiological

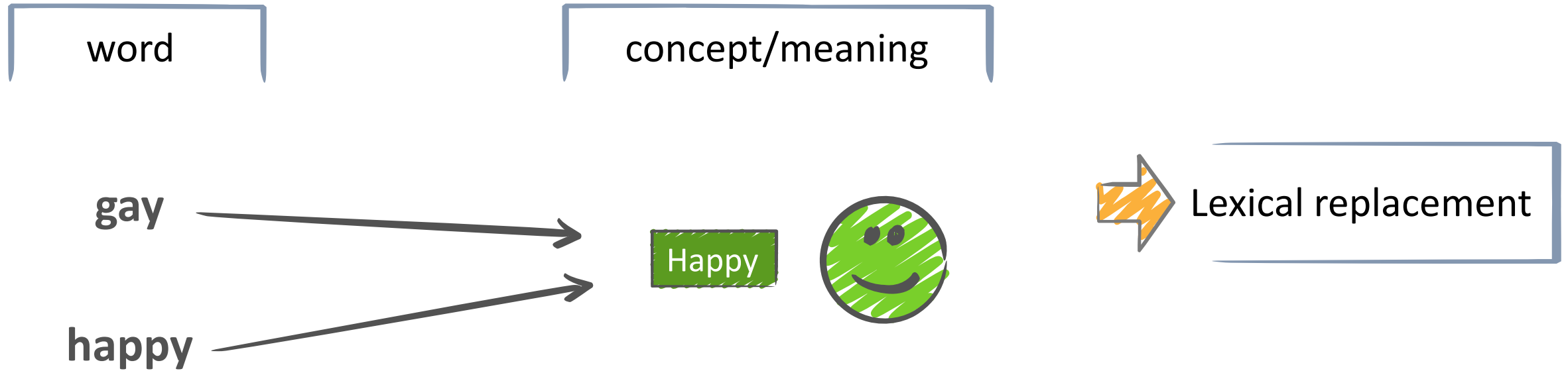


Onomasiological

# Semasiological perspective

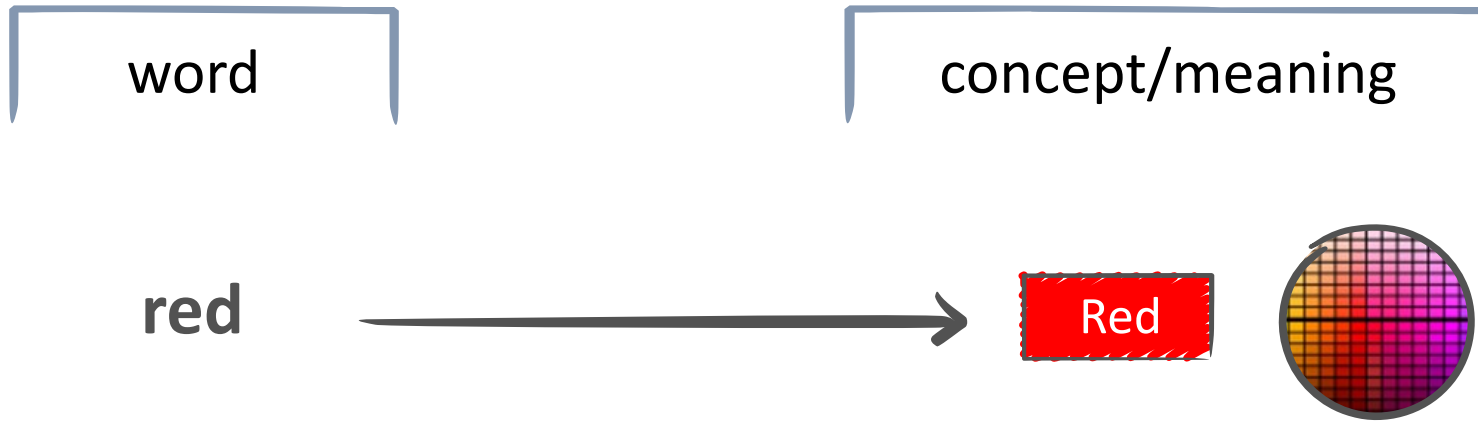


# Onomasiological perspective

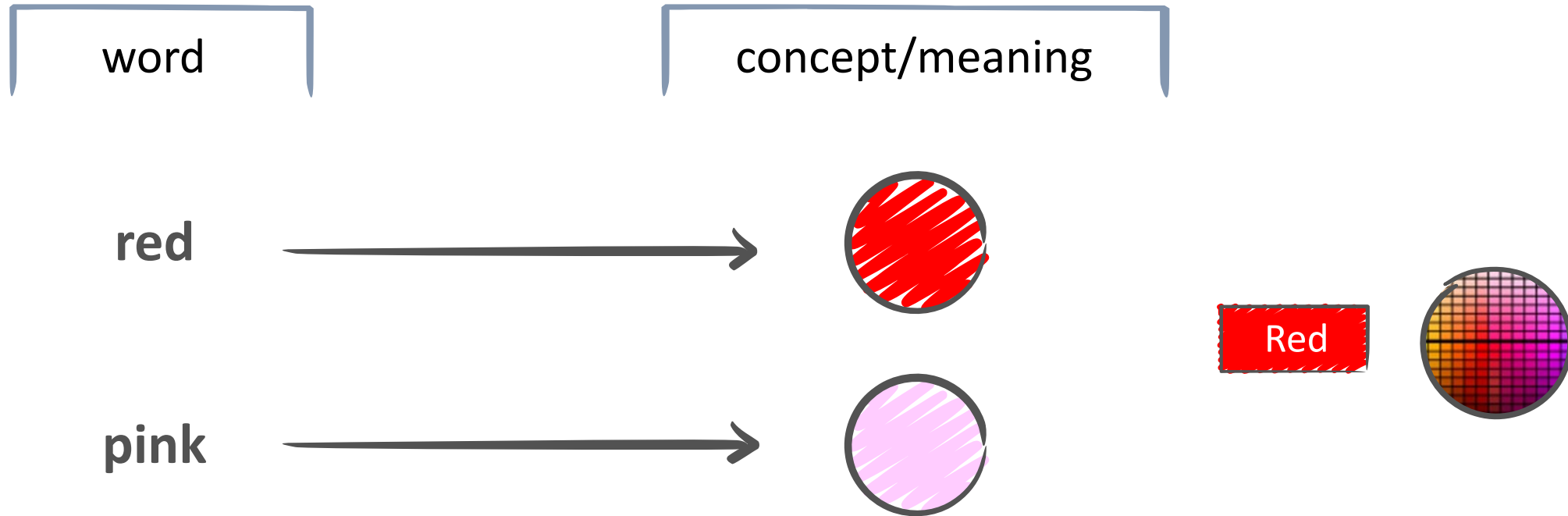




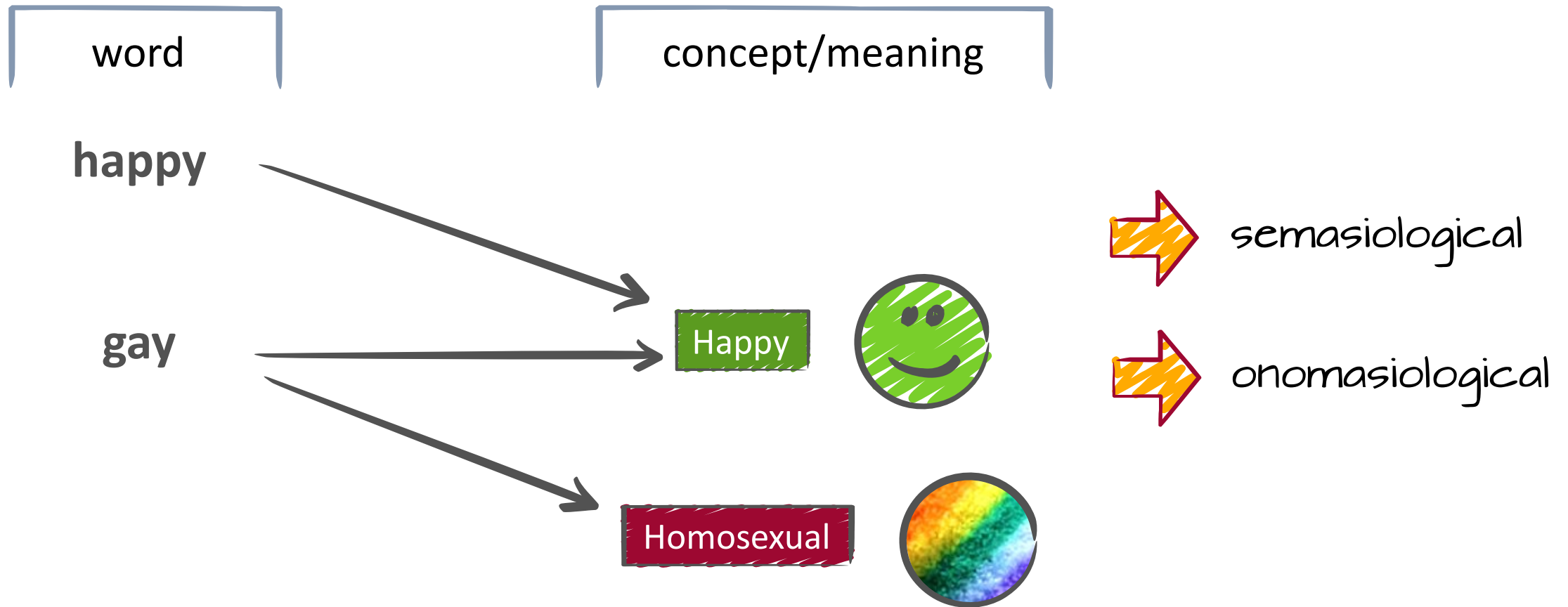
# Ono- and Semasiological are interlinked!



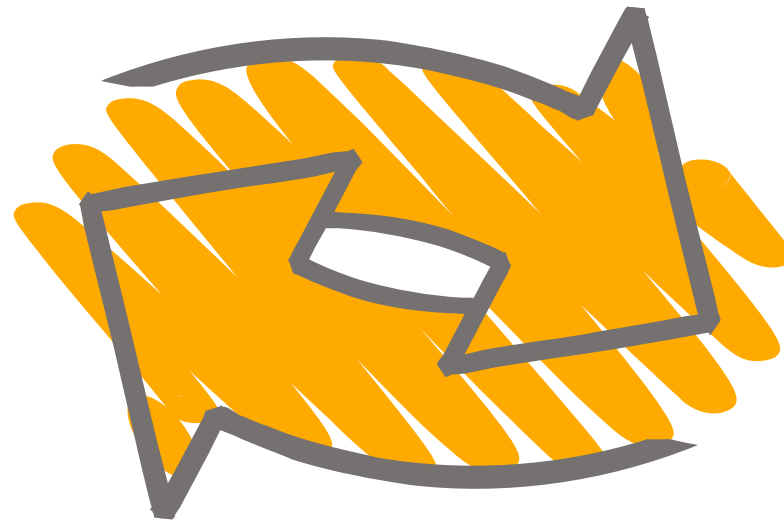
# Ono- and Semasiological are interlinked!



# One more example



Why?

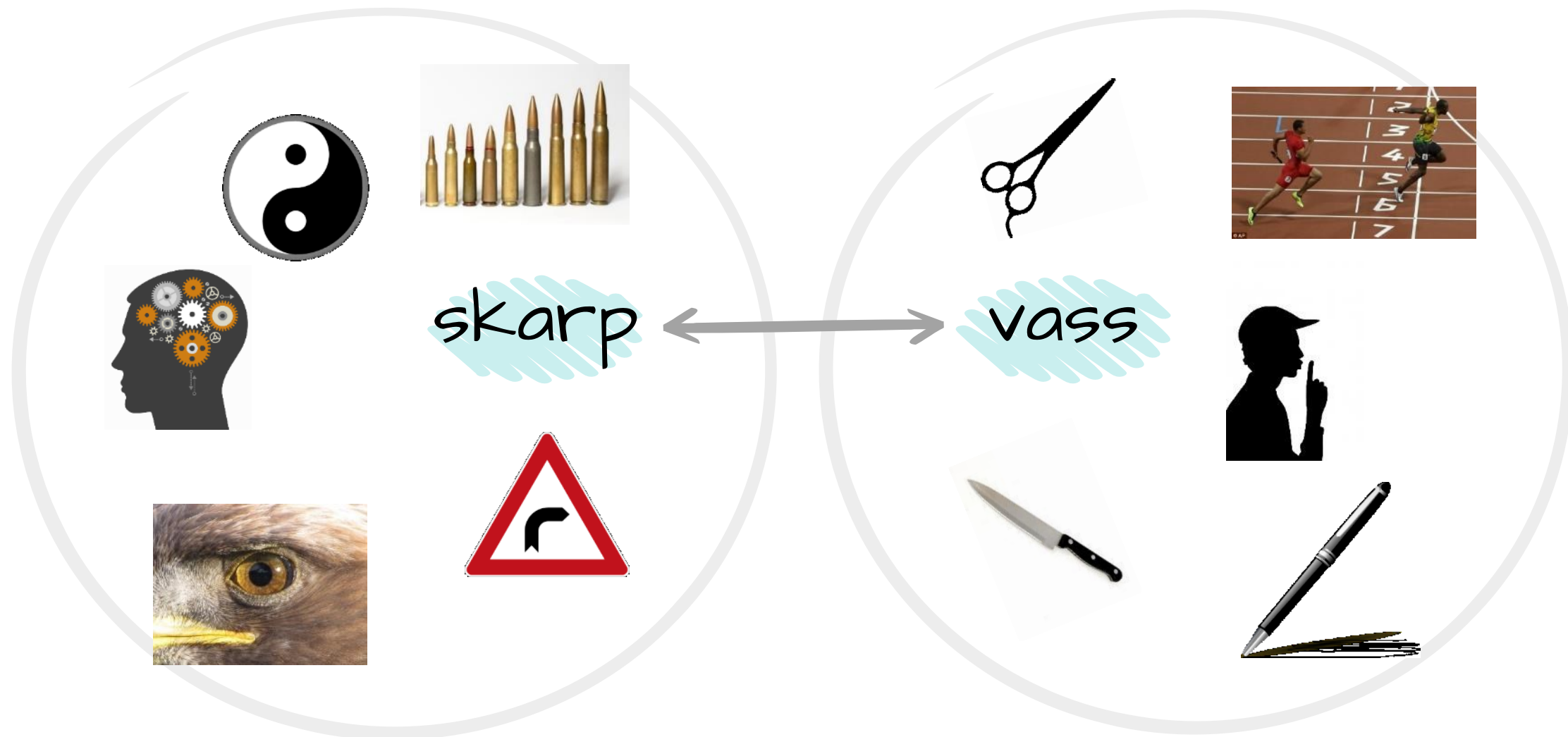




Nina Tahmasebi, Lexical Semantic Change Workshop, COMHIS,  
Helsinki



# A division of the semantic field 'sharp'



skarp



vass

skarp (adjektiv)

Skarp

Attribut

Adverbial

skarp

1. kritik 1008
2. kontrast 822
3. ammunition 358
4. version 299
5. gräns 246
6. blick 213
7. kritiker 141
8. varning 151
9. kurva 125
10. kant 77
11. analys 89
12. sväng 73
13. bild 169
14. protest 82
15. tillsägelse 34

1. lika 120
2. mycket 227
3. ganska 82
4. så 270
5. alltför 24
6. liten 34
7. föga 32
8. riktig 75
9. osedvanlig 9
10. tillräcklig 18
11. i går 7
12. oväntad 10
13. samtidigt 8
14. oerhörd 17
15. uppknappt 4

vass (adjektiv)

Vass

Attribut

Adverbial

vass

1. kniv 803
2. spurt 200
3. penna 114
4. avslutning 118
5. tunga<sup>2</sup> 79
6. tunga 79
7. avslutare 54
8. kant 72
9. egg 45
10. speed 36
11. slutspeed 26
12. sax 30
13. satir 34
14. målskytt 44
15. sten<sup>2</sup> 36

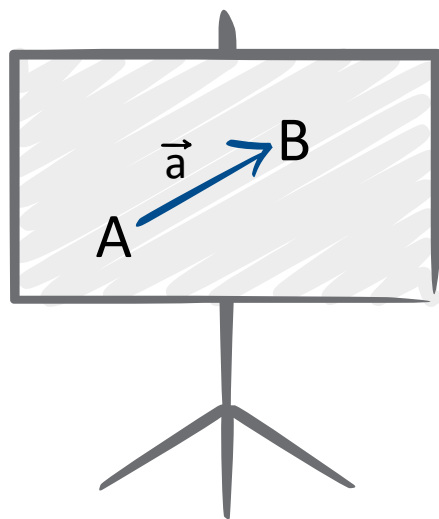
1. riktig 529
2. lika 230
3. tillräcklig 31
4. jävligt 29
5. jävlig 29
6. ruskig 14
7. jäklig 14
8. grön 5
9. oerhörd 16
10. ovanlig 11
11. ganska 41
12. ruggig 5
13. speciell 8
14. onekligen 5
15. invändig 2



# Methods for computational semantic change



# Some terminology

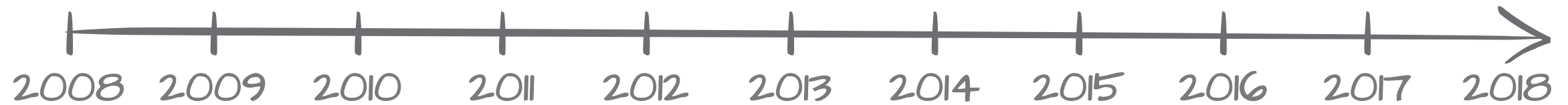


Vector (1, 4, 3) (=3 dimensions)

Topic modeling

- embeddings
- neural embeddings
- dynamic embeddings

Single-sense

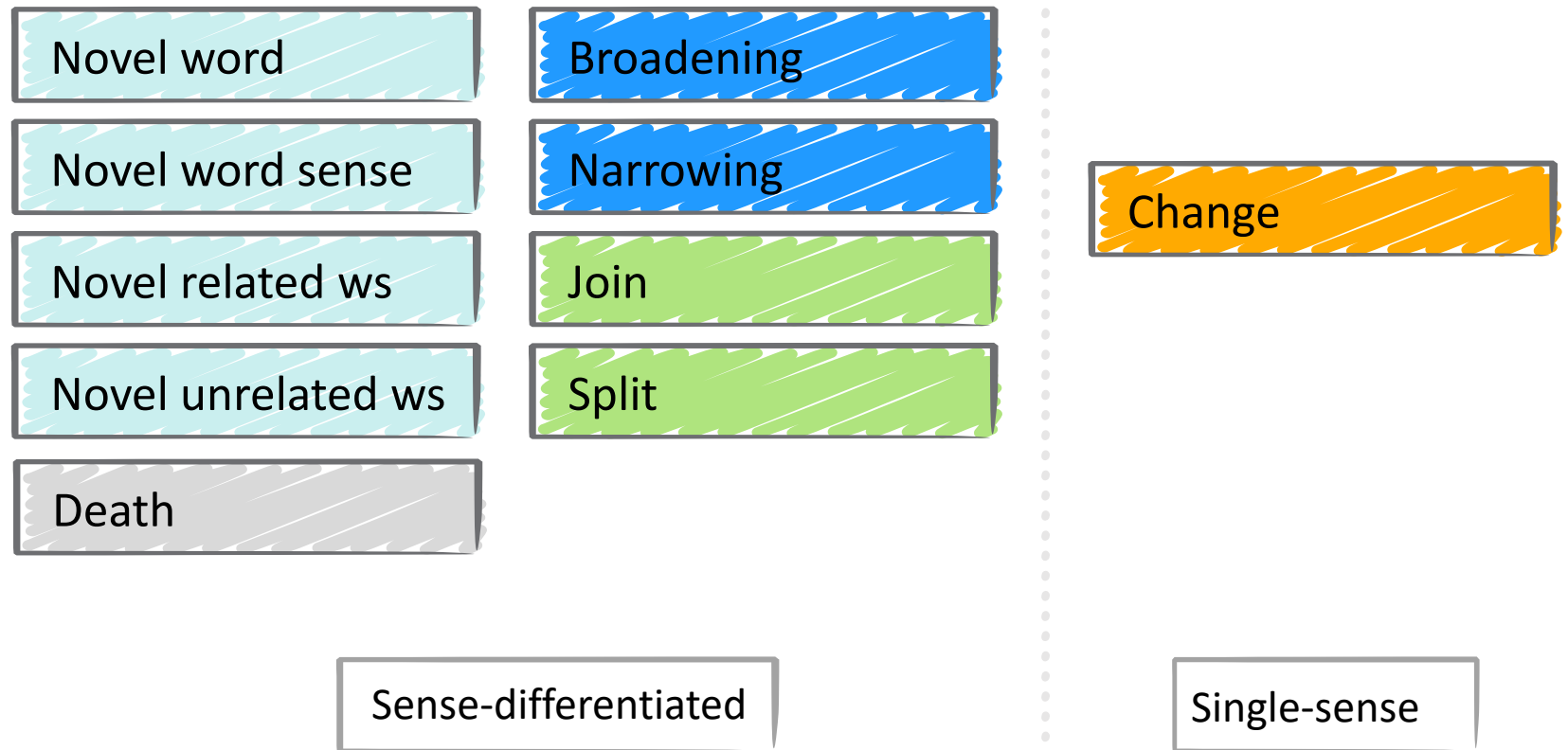
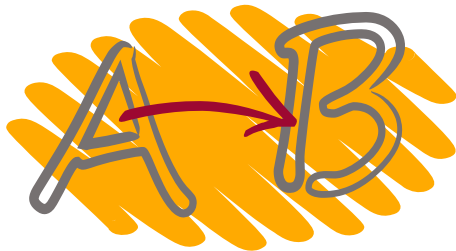


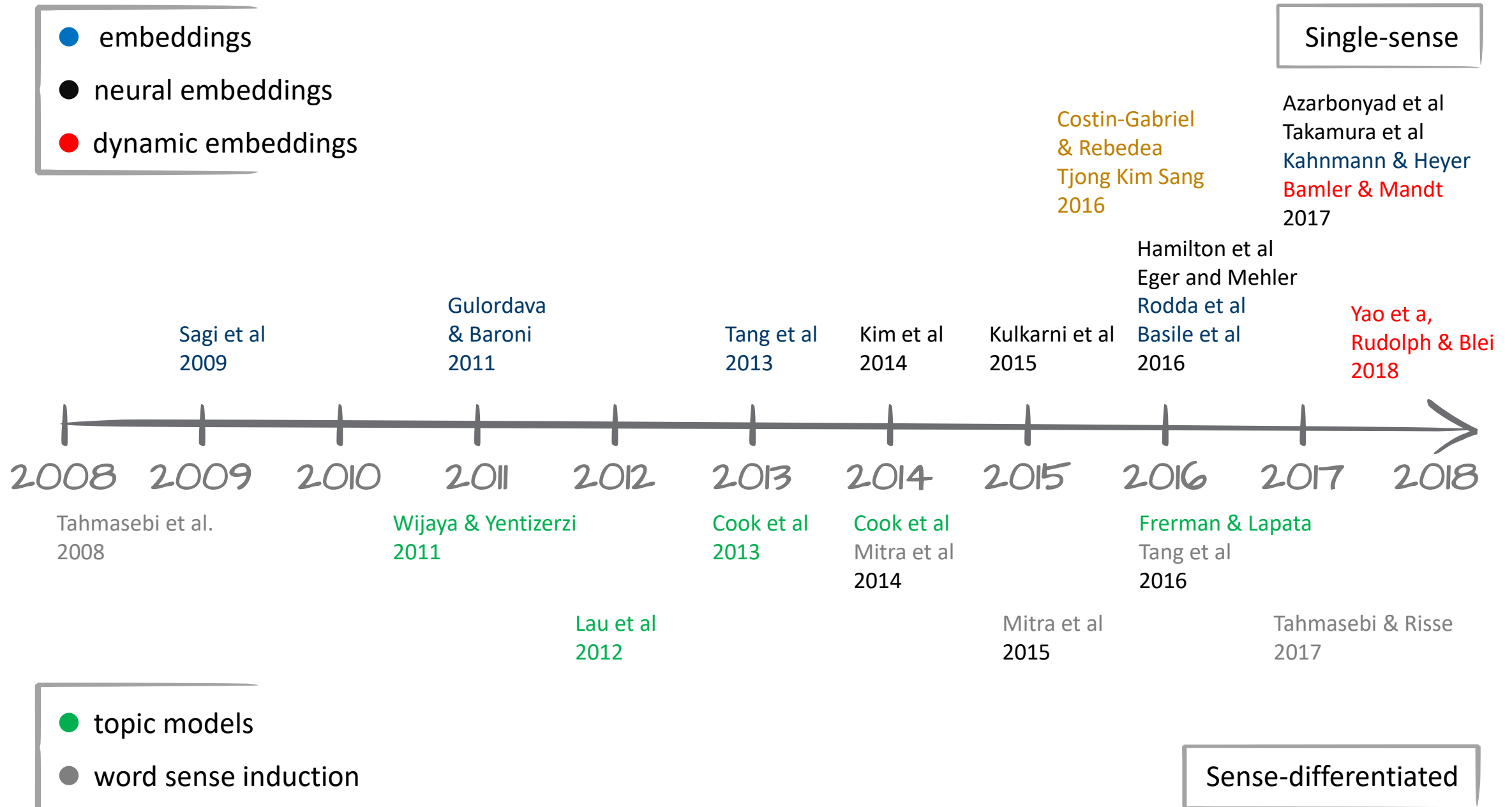
- topic models
- word sense induction

Sense-differentiated



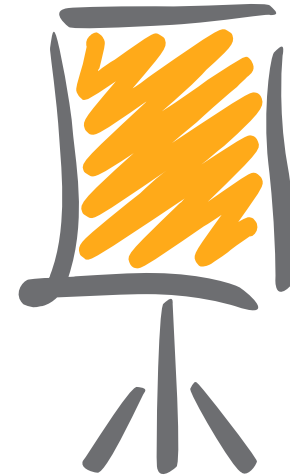
# Change type





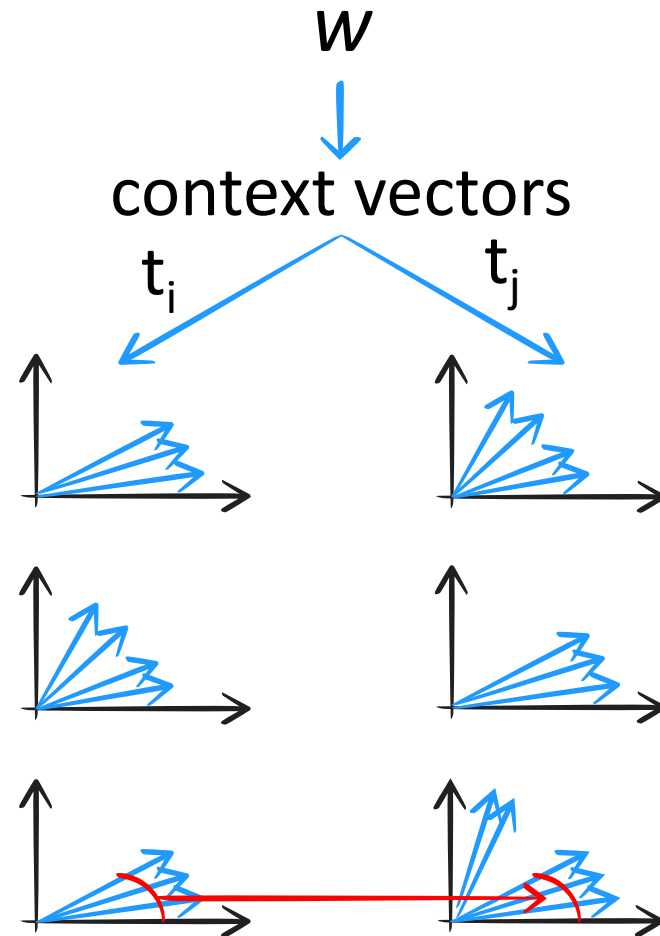
# Outline

- embeddings / context-based methods
  - neural embeddings
  - dynamic embeddings
- topic models
- word sense induction



# Context-based method

Sagi et al.  
GEMS 2009



Data set split in approp. sets

Broadening of sense

Narrowing of sense

With grouping:  
Added/removed sense

BUT: 1.

No discrimination between senses

2.

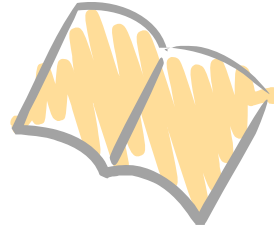
No alignment of senses over time!

# Word embedding-based models

---

Kulkarni et al. WWW'15

---



Project a word onto a vector/point  
(POS, frequency and embeddings)



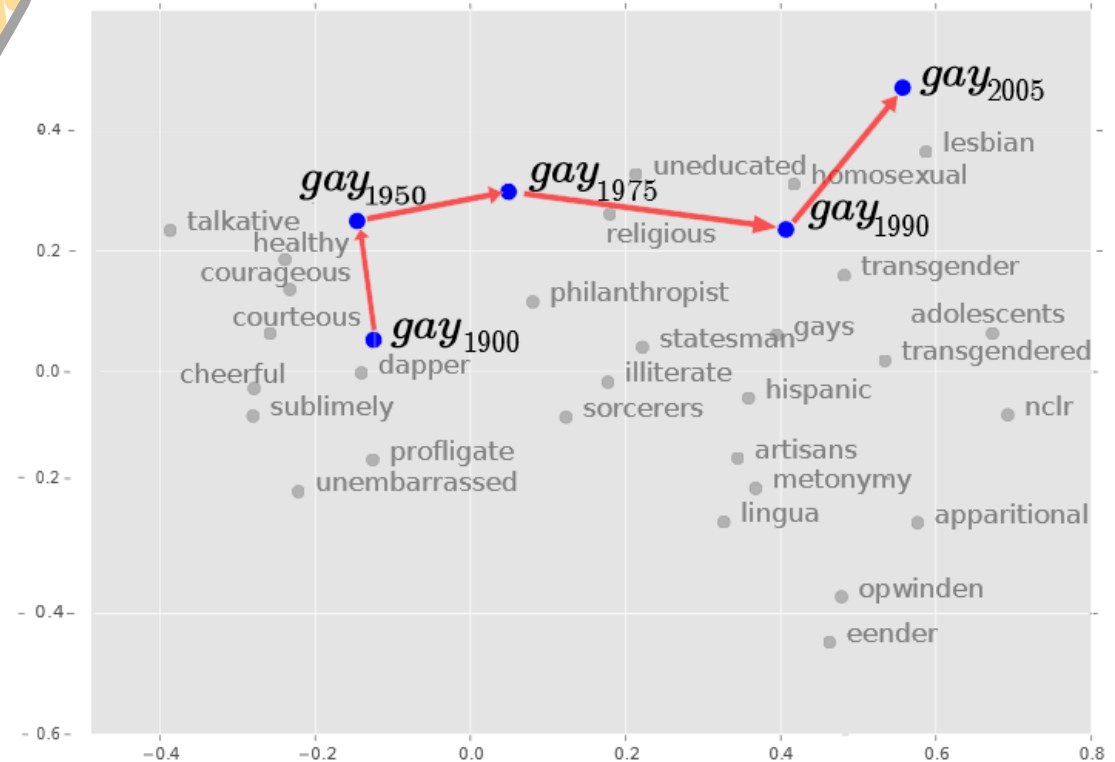
Track vectors over time

Kim et al. LACSS 2014

Basile et al. CLiC-it 2016

Hamilton et al. ACL 2016

---





# Dynamic Embeddings

---

Share data across all time points  
Avoids aligning

---

Bamler & Mandt:

- Bayesian Skip-gram

Yao et al:

- PPMI embeddings

Rudolph & Blei:

- Exponential family embeddings  
(Bernoulli embeddings)



Sharing data is **highly beneficial!**

# Topic-based methods

- 1 Topic model (HDP)
- 2 Assign topics to all instances of a word.
- 3 If a word sense  $WS_i$  is assigned to collection 2 but not 1 then  $WS_i$  is a **novel** word sense.

**BUT:**

- A Only two time points (typically there is much noise!)
- B **No alignment** of senses over time!

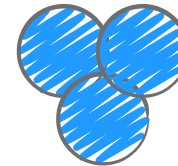
Lau et al.  
EACL 2014

Wijaya & Yeniterzi  
DETECT '11

Cook et al.  
Coling 2014

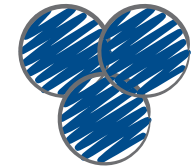
Frermann & Lapata  
TACL 2016

BNC



Finally, we conduct a preliminary evaluation in which we apply our methods to the task of

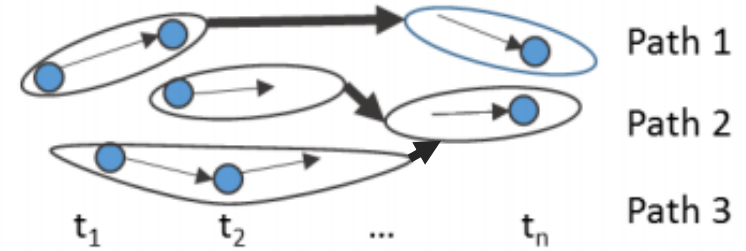
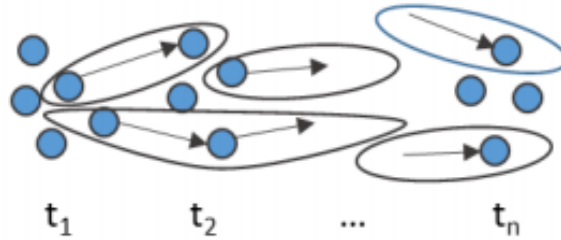
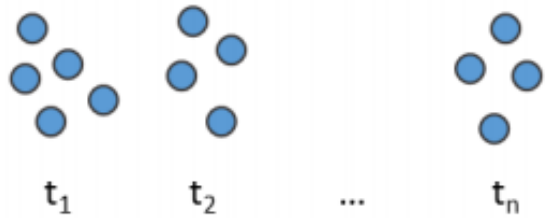
ukWaC



meanings of words are not fixed but in fact they do change



# Word sense induction



## Step 1:

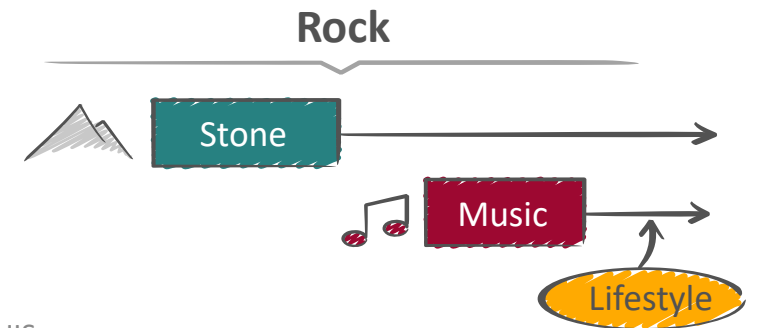
Word sense discr.  
(curvature clustering)  
individual time slices

## Step 2:

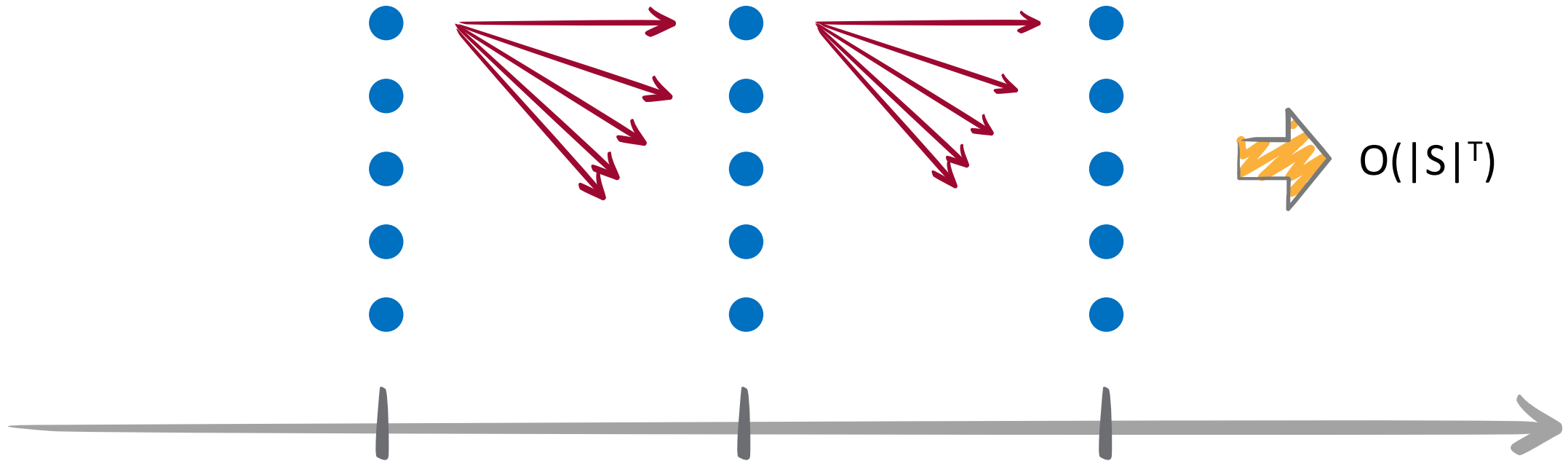
Detecting stable  
senses  
→ units

## Step 3:

Relating units  
→ Paths



# Complexity



# How?



Nina Tahmasebi, Lexical Semantic Change Workshop, COMHIS,  
Helsinki







dec. ut. r. q. m. q.  
 quoniam. et  
**C**onsuetudo est.  
 autem dicitur iustitia  
 tot hinc uelut  
 tio. phetia. eli  
 si. ue possit ma  
 liguata. sic tunc  
 phetia. ducere  
 pupilla. ante  
 reddita. r. q. m. q.  
 ue possit ma  
 gnari. C. te. m. r.  
 dicto. m. r. q. m. q.  
 tot. sic uoc. a. d.  
 ton. licet. s. pig  
 un. adicere. C.  
 qui. bonis. co. d.  
 pos. l. m. Jte. m. r.  
 tator. u. d. s. ho  
 uore. am. m. e. b.  
 y. q. r. m. r. q. m. q.  
 tot. Jte. q. regu  
 la. illa. p. m. r. q. m. q.

# NLP pipeline: From text to result

**Text-mining method** 

**Dimensions**

Filtering: Function words

Filtering: Stopwords

Part-of-speech tagging

Lemmatization

Tokenization





like

(only verbs)

room

(frequency filtering)

room

sheet. (only nouns)

I like room

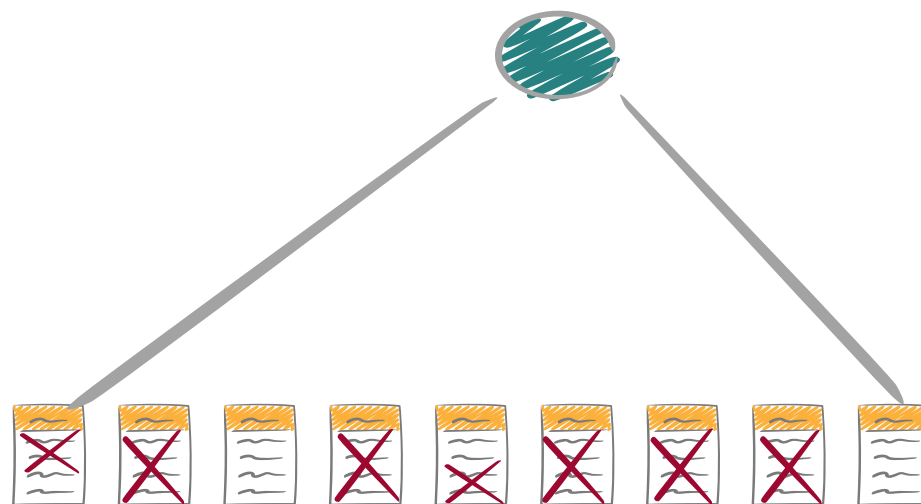
sheet. (after lemmatization)

I like room

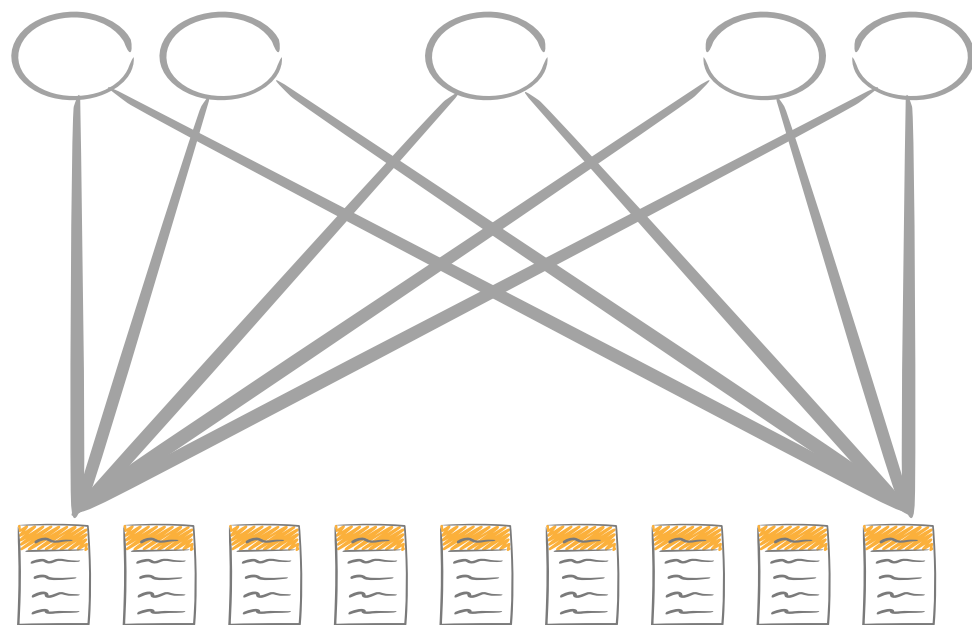
sheets. (after stop word filtering)

I like the room but not the sheets.

# Viewpoint on the data

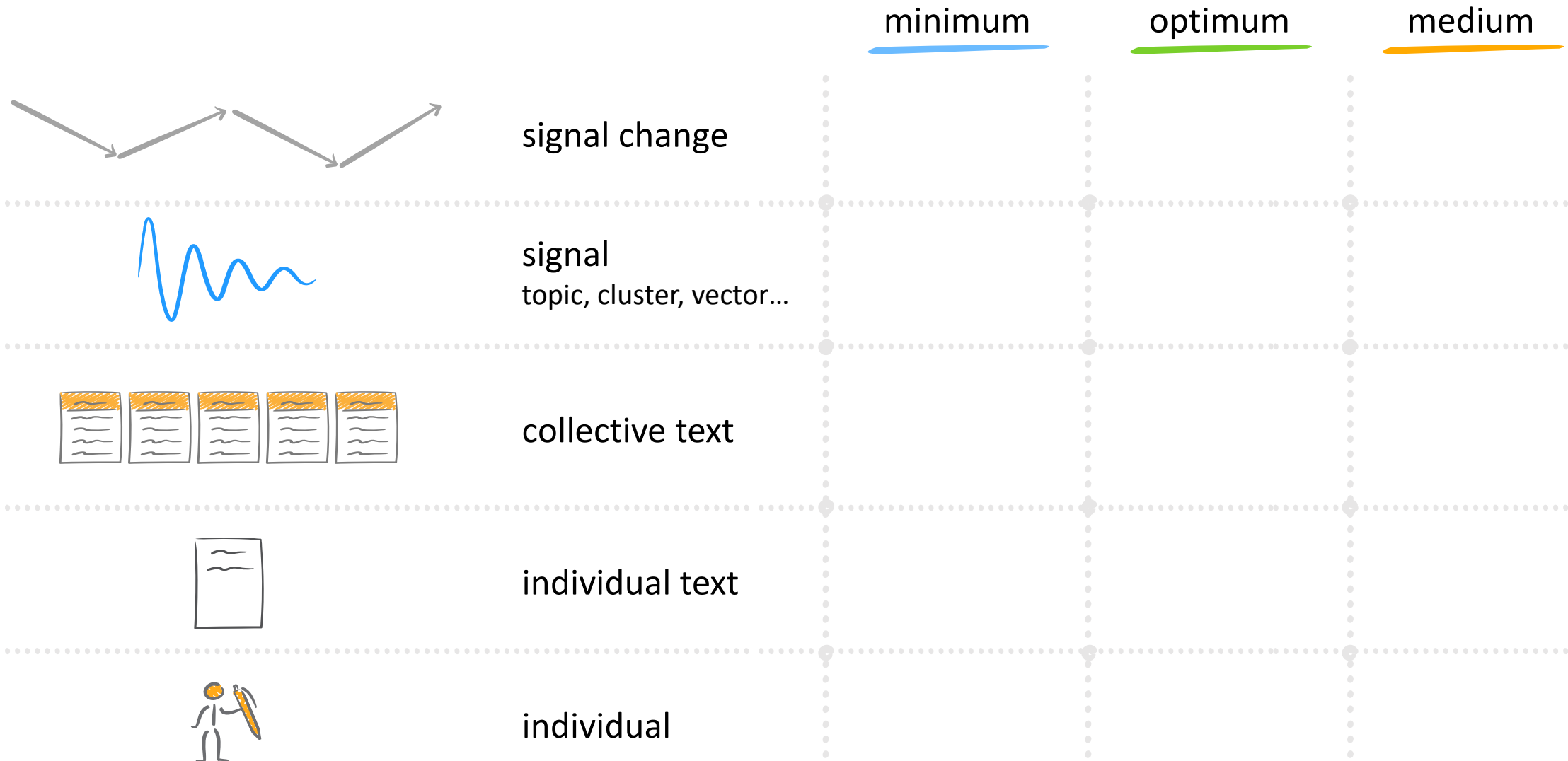


# Viewpoint on the data (cont'd)

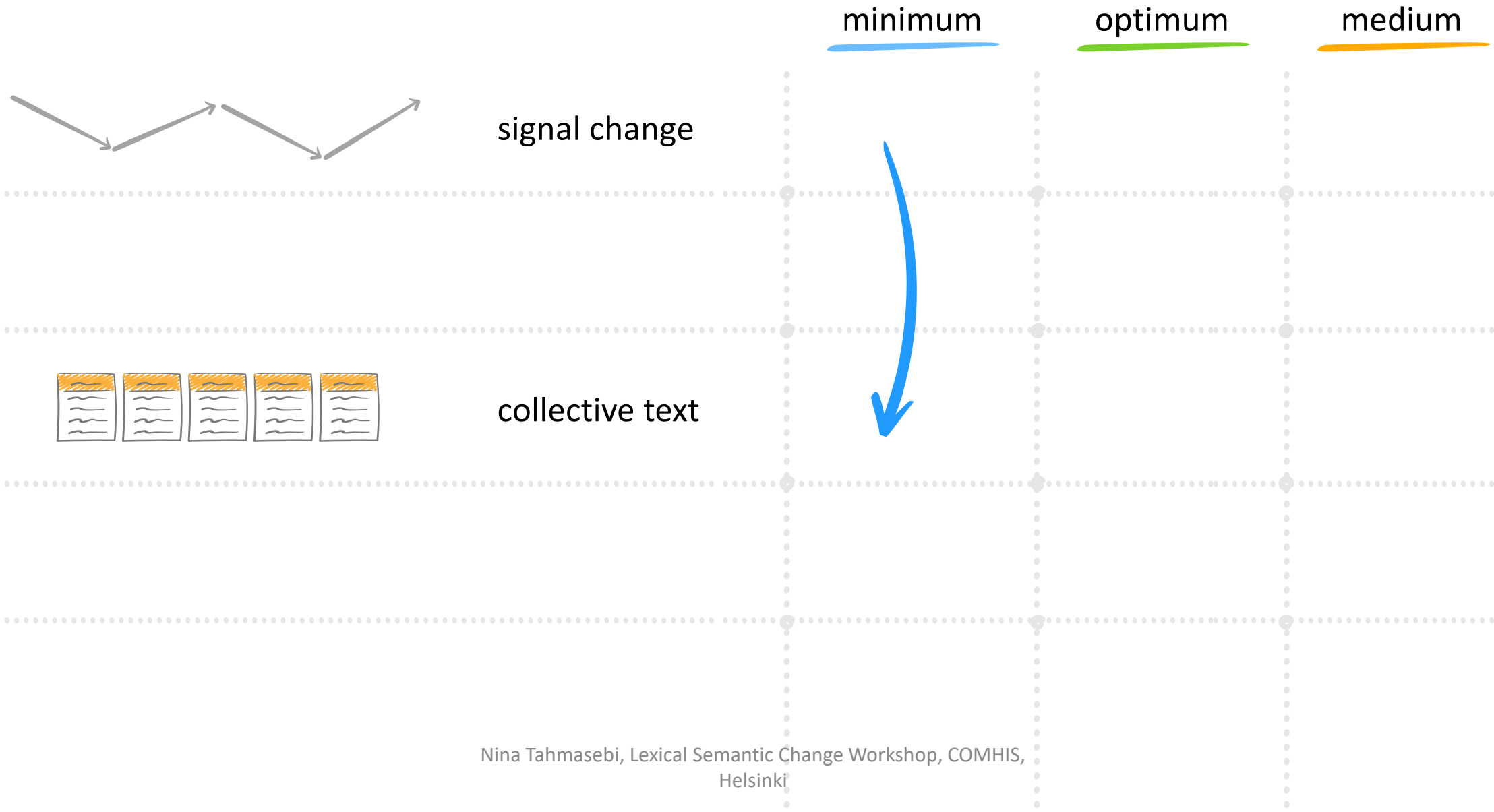




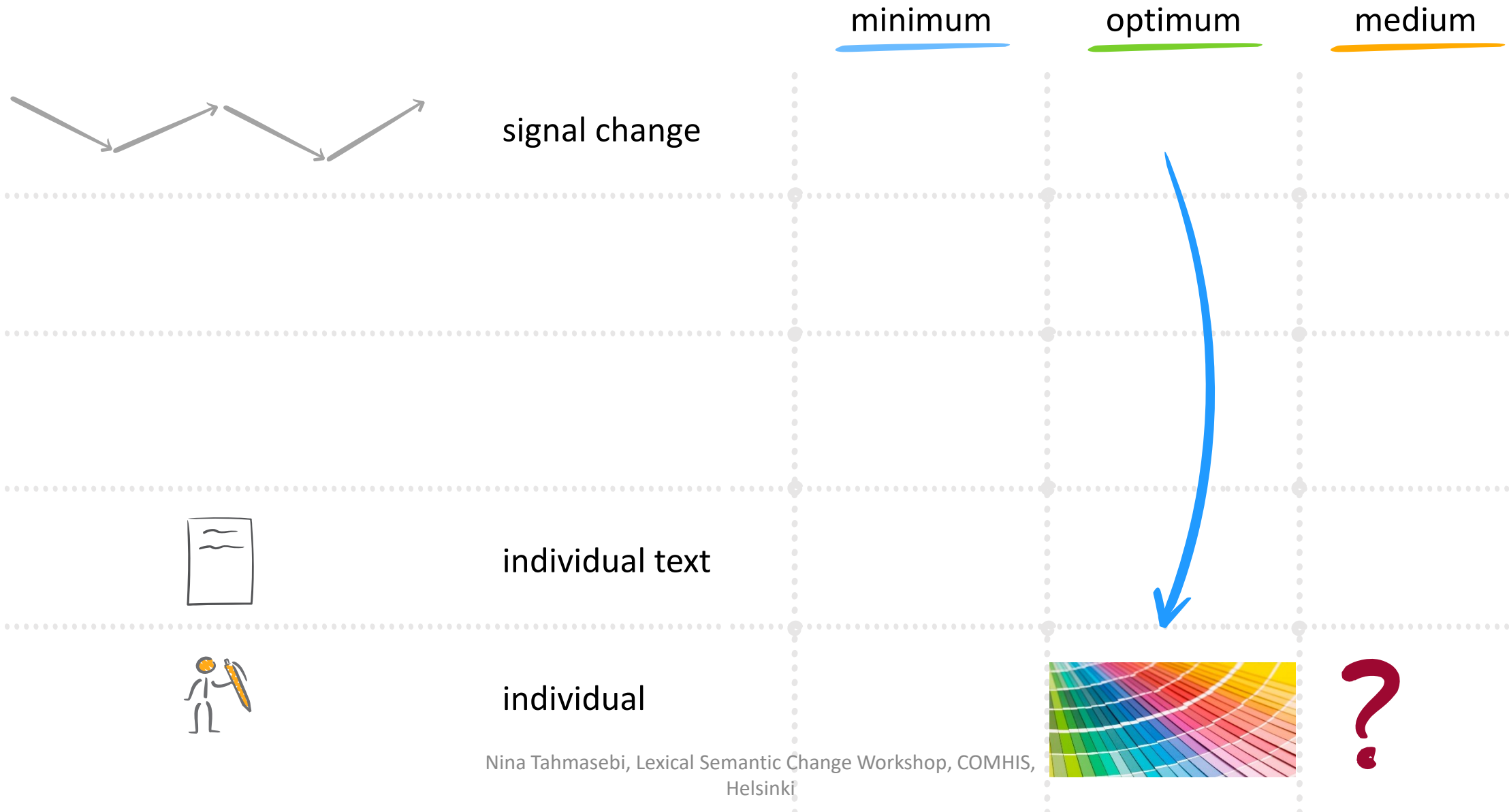
# Evaluation



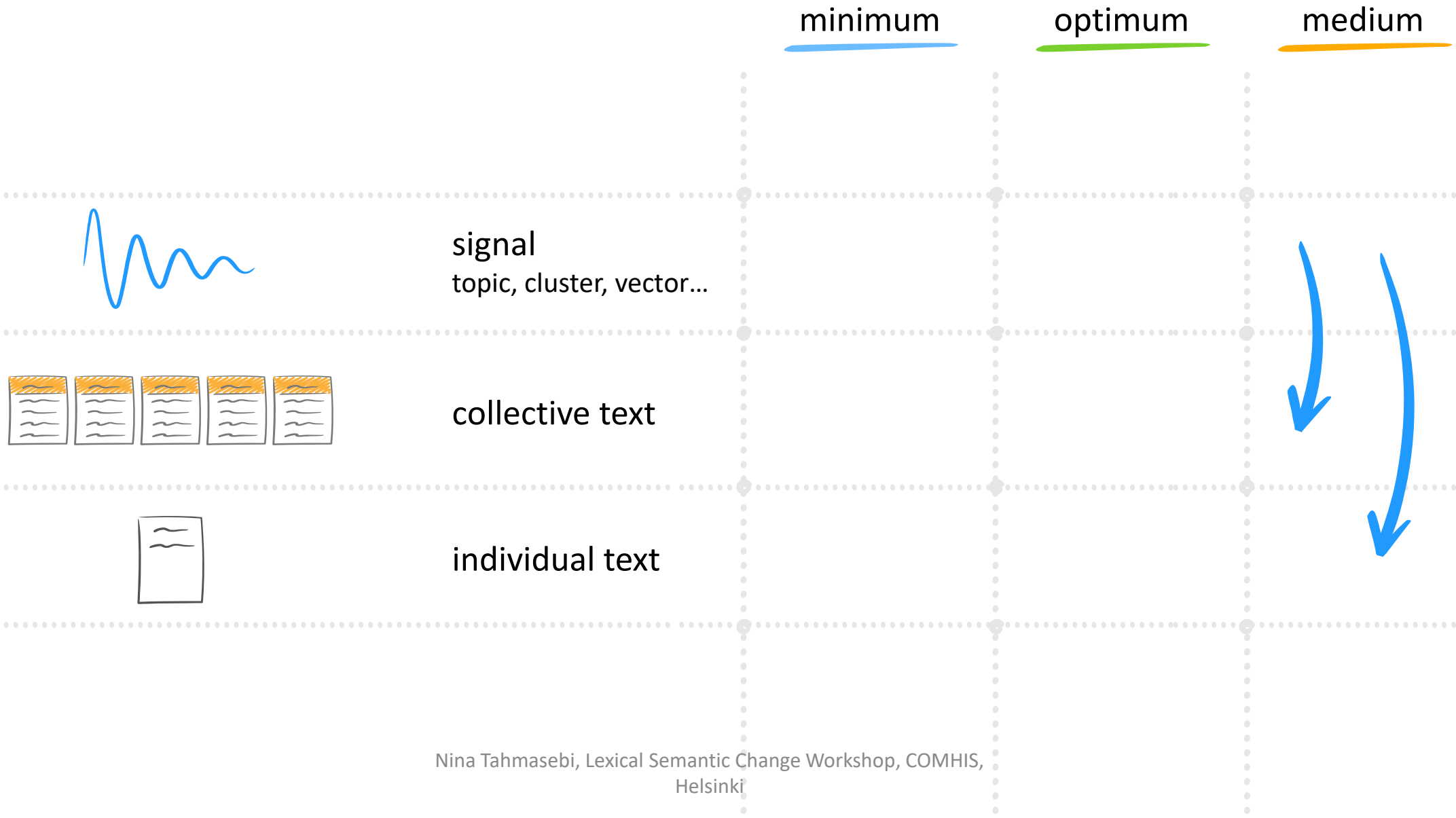
# Evaluation



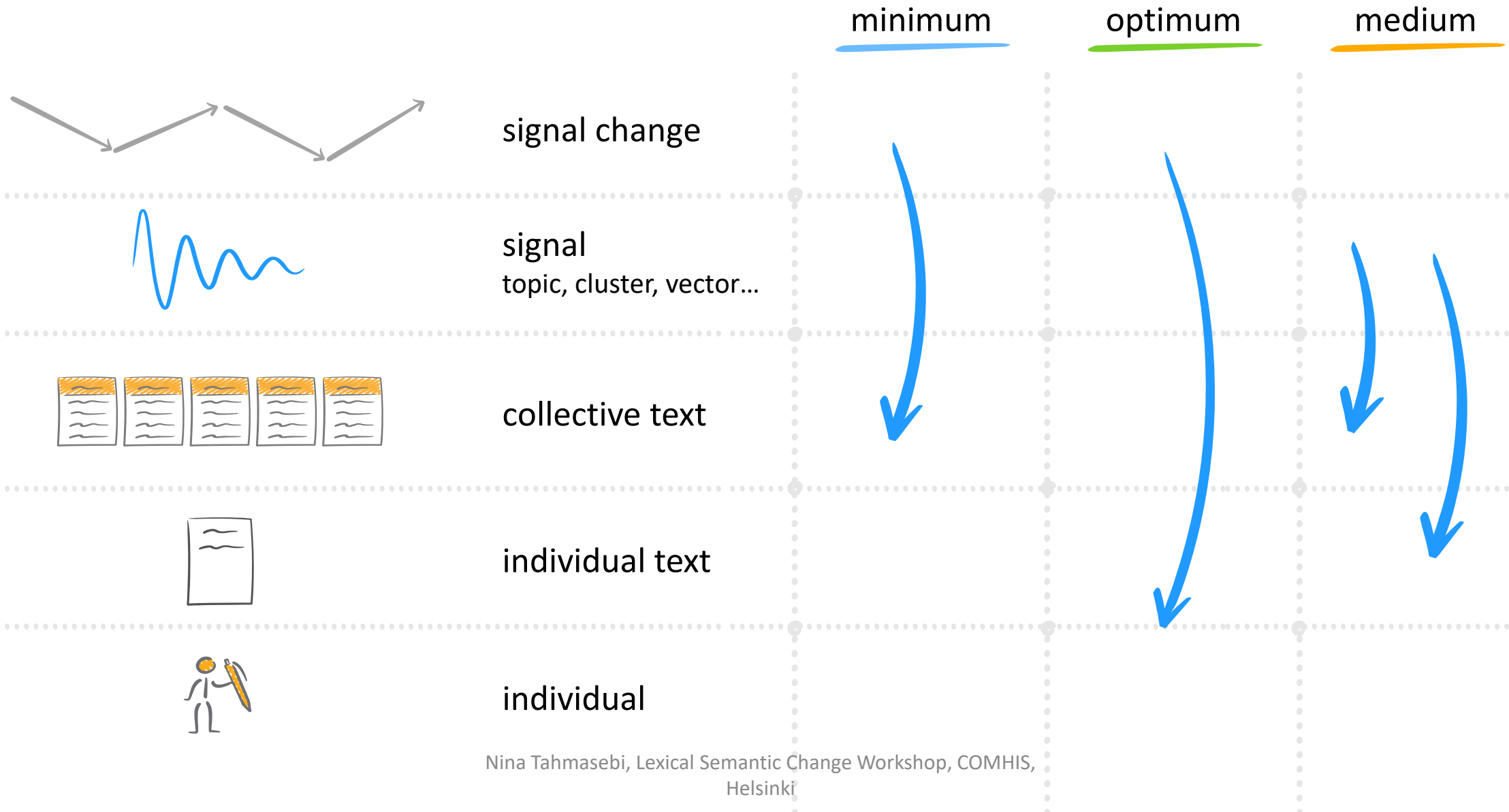
# Evaluation



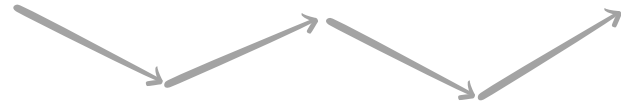
# Evaluation



# Evaluation



# Evaluation



signal change

minimum

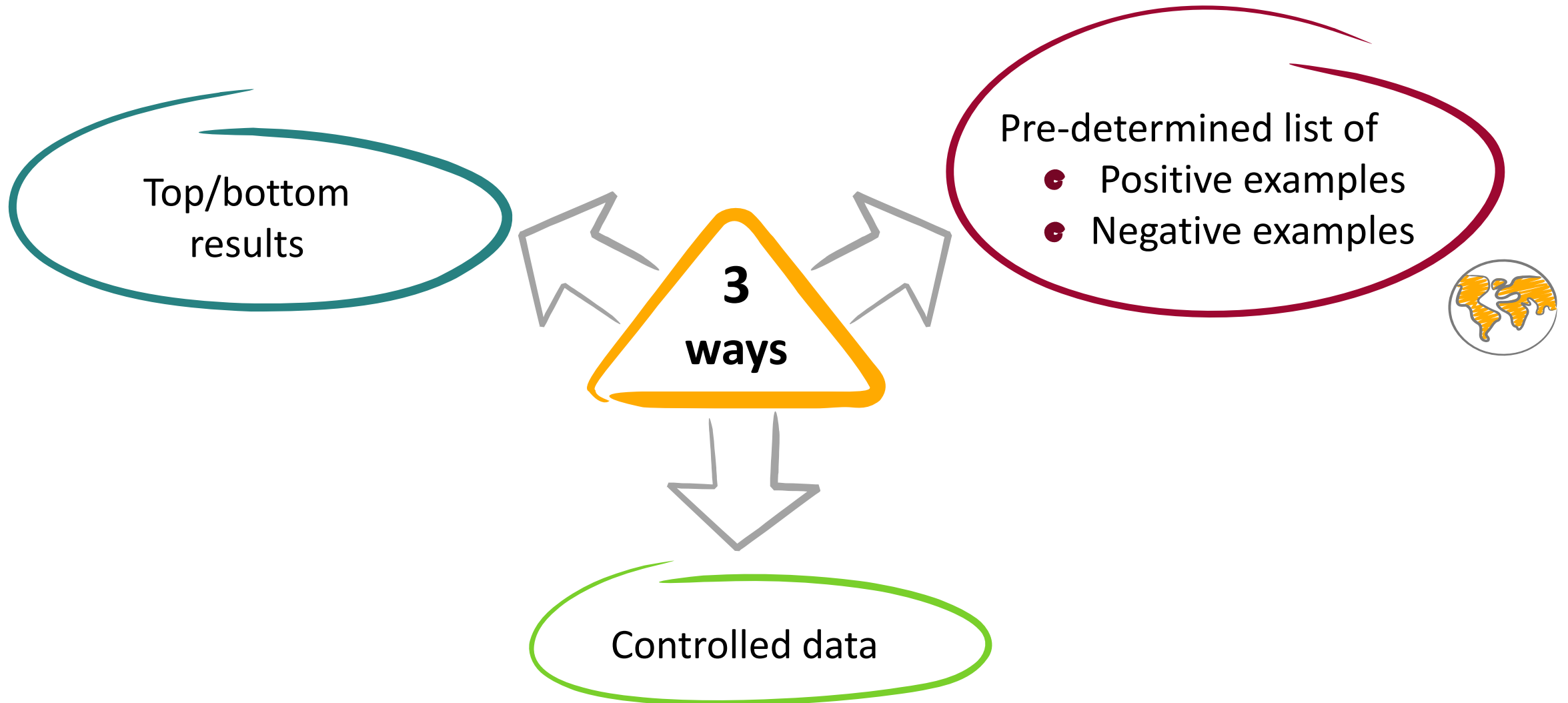


optimum

medium



# Evaluation



	prechosen		top	entity	eval. method	span	time	# points	# classes	classes	modes	
	# pos	# neg		(S)ingle/ (P)airs	(M)anual/ (A)utomatic						time / sense	aware / diff
Sagi, Kaufmann, and Clark (2009a)	4	0		S	M	569y		4	2	broad./narrow.	no	no
Gulordava and Baroni (2011)	0	0	100 <sup>54</sup>	S	M	40y		2	1	change	no	no
Tang, Qu, and Chen (2013)	33	12		S	M	59		59	3	B/N/novel/change <sup>55</sup>	no	no
Kim et al. (2014)	0	0	10/10 <sup>56</sup>	S/P <sup>57</sup>	M	110		110	1	change	yes <sup>58</sup>	no
Kulkarni et al. (2015)	20	0	20 <sup>59</sup>	S	M/A	105y/12y/2y		21/13/24	1	change	yes	no
Hamilton, Leskovec, and Jurafsky (2016b)	28	0	10 <sup>60</sup>	S/P	M	200/190		20	1	change	no	no
Rodda, Senaldi, and Lenci (2016)	0	0	50	S	M	1200y		2	1	change	no	no
Eger and Mehler (2016)	0	0	21 <sup>61</sup>	S/P	M	200/190		20/19	1	change	no	no
Basile et al. (2016)	40	0		S	M	170		17	1	change	yes	no
Azarbonyad et al. (2017)	24	0	5/5 <sup>62</sup>	S	M	20/11		2/2	1	change	no	no
Takamura, Nagata, and Kawasaki (2017)	10	0	100/20 <sup>63</sup>	S/P	M	- <sup>64</sup>		2	1	change	no	no
Kahmann, Niekler, and Heyer (2017)	4	0		S	M	≤ 1 <sup>65</sup>		48	1 <sup>66</sup>	change	no	no
Bamler and Mandt (2017)	6	0	10	S/P	M <sup>67</sup>	209/230/7		209/230/21	1	change	no	no
Yao et al. (2018)	4/1888 <sup>68</sup>	0		S	M/A	27		27	1	change	no	no
Wijaya and Yeniterzi (2011)	4	2		S	M	500 <sup>69</sup>		500	2 <sup>70</sup>	change novel	yes	yes <sup>71</sup>
Lau et al. (2012)	5	5		S	M	43 y		2	1	novel	no	yes
Cook et al. (2013)	0	0	30	S	M	14		2	1	novel	no	yes
Cook et al. (2014)	7/13	50/164		S	M	43y/17y		2/2	1	novel	no	yes
Mitra et al. (2015) <sup>72</sup>	0	0	69/50	S	M/A	488/2		8/2	3	split/join/novel <sup>73</sup>	no	yes
Frermann and Lapata (2016)	4	0	200	S	M/A	311		16	2	change/novel	no	yes
Tang, Qu, and Chen (2016) <sup>74</sup>	197	0		S	M	59		59	6	B/N/novel/change <sup>75</sup>	no	yes
Tahmasebi and Risse (2017a)	35	25		S	M	222y		221	4	novel,B/N,stable	yes	yes

<https://languagechange.org/publication/2018-surveypaper/>

Table 3

Datasets used for diachronic conceptual change detection. Non-English ·

Sagi, Kaufmann, and Clark (2009a)	Helsinki corpus
Gulordava and Baroni (2011)	Google Ngram
Wijaya and Yeniterzi (2011)	Google Ngram
Lau et al. (2012)	British National Corpus (BNC), ukWaC
Cook et al. (2013)	Gigawords corpus
Cook et al. (2014)	BNC, ukWaC, Sibol/Port
Mihalcea and Nastase (2012)	Google books
· Basile et al. (2016)	Google Ngram (Italian)
· Tang, Qu, and Chen (2013, 2016)	Chinese People's Daily
Kim et al. (2014)	Google Ngram
Kulkarni et al. (2015)	Google Ngram, Twitter, Amazon movie reviews
Mitra et al. (2015)	Google Ngram, Twitter
Hamilton, Leskovec, and Jurafsky (2016b)	COHA, Google Ngram
· Eger and Mehler (2016)	COHA, Süddeutsche Zeitung, PL <sup>76</sup>
Azarbonyad et al. (2017)	New York Times Annotated Corpus, Hansard
· Rodda, Senaldi, and Lenci (2016)	Thesaurus Linguae Graecae
Frermann and Lapata (2016)	DATE corpus
Takamura, Nagata, and Kawasaki (2017)	Wikipedia (English and Japanese)
Kahmann, Niekler, and Heyer (2017)	Guardian (non-public)
Tahmasebi and Risse (2017a)	Times Archive, New York Times Annotated Corpus
Bamler and Mandt (2017)	Google Ngram, State of the Union addresses, Twitter
Yao et al. (2018)	New York Times (non-public)
Rudolph and Blei (2018)	ACM abstracts, ML papers ArXiv, U.S. Senate speech



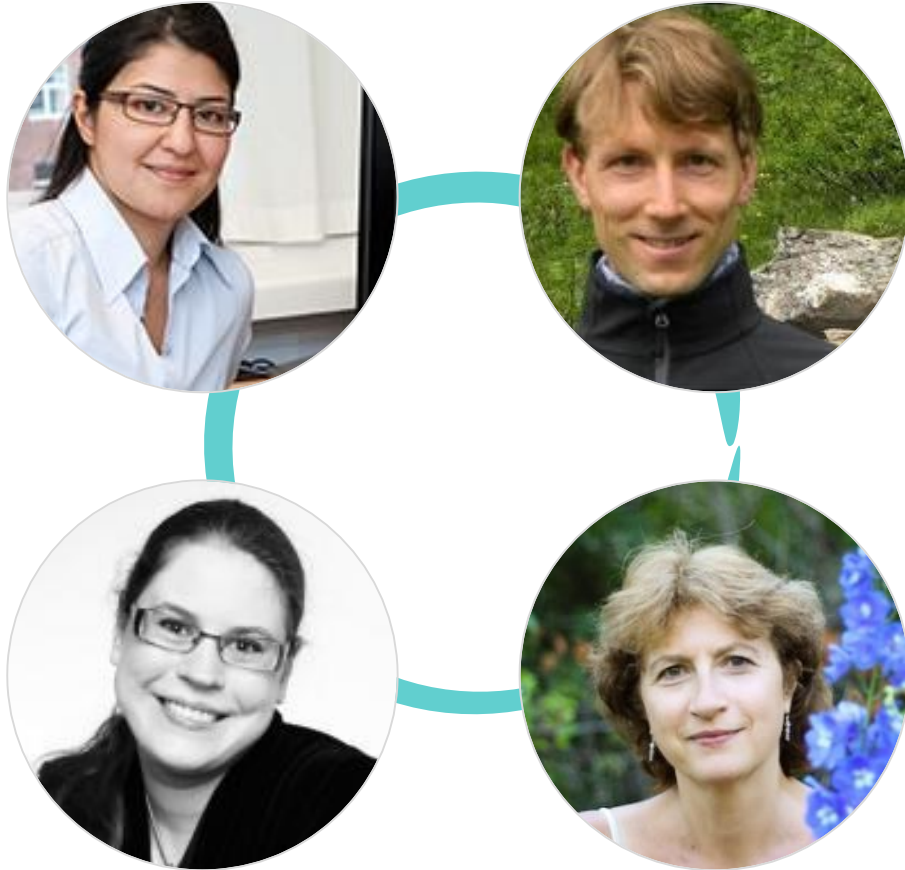
# Towards automatic language change detection

VR funded

6 million sek (+ cofunding Språkbanken ~700k sek)

2019 – 2022

# 4 year project: <https://languagechange.org/>



Overall goal is to bridge the gap between the four of us and all that can benefit from the results.

# Main goals

## Wp1: Swedish word sense induction

- Using sense-differentiated dynamic embeddings

## Wp3: Lexical replacements

- On the basis of Wp1
- Or using other textual clues



## Wp2: Semantic change

- On the basis of Wp1

## Wp4: Applications

- Applied sociology, historical linguistics, history of concepts, ...

## WP\*: Evaluation

- Integrated in all work packages



# Planned activities

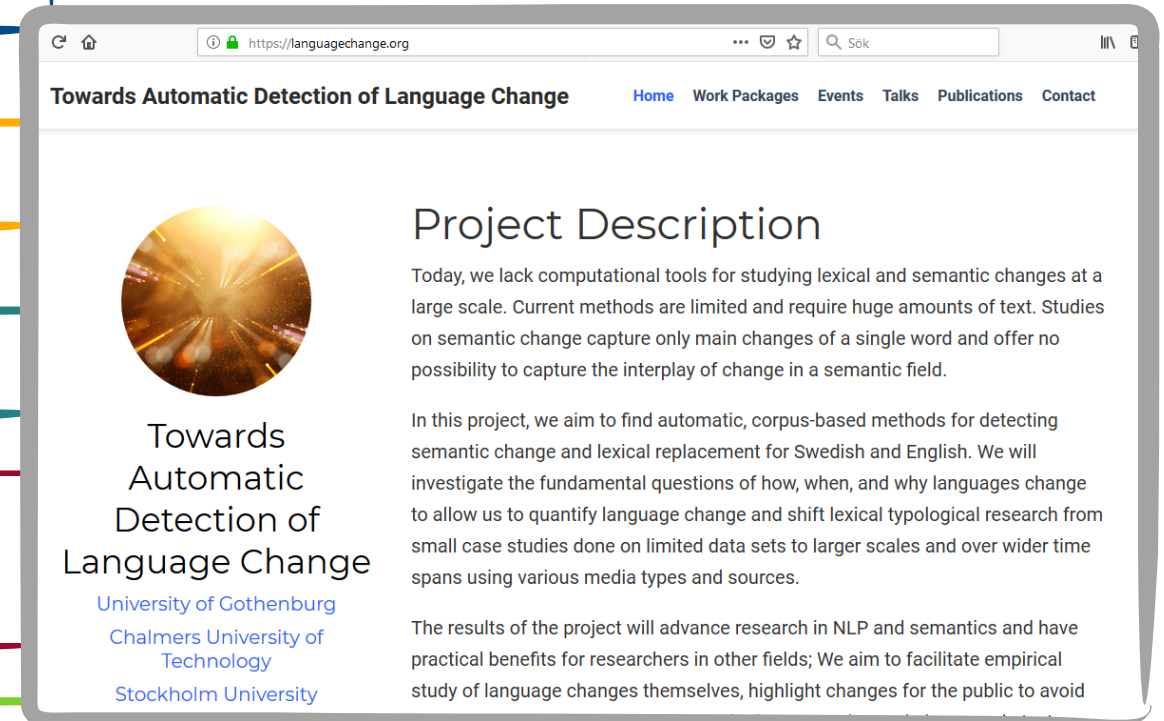
News-list ([news@languagechange.org](mailto:news@languagechange.org))

Introductory videos to LS change

Workshops (next at ACL2019)

Work on evaluation  
(possibly in a SemEval task)

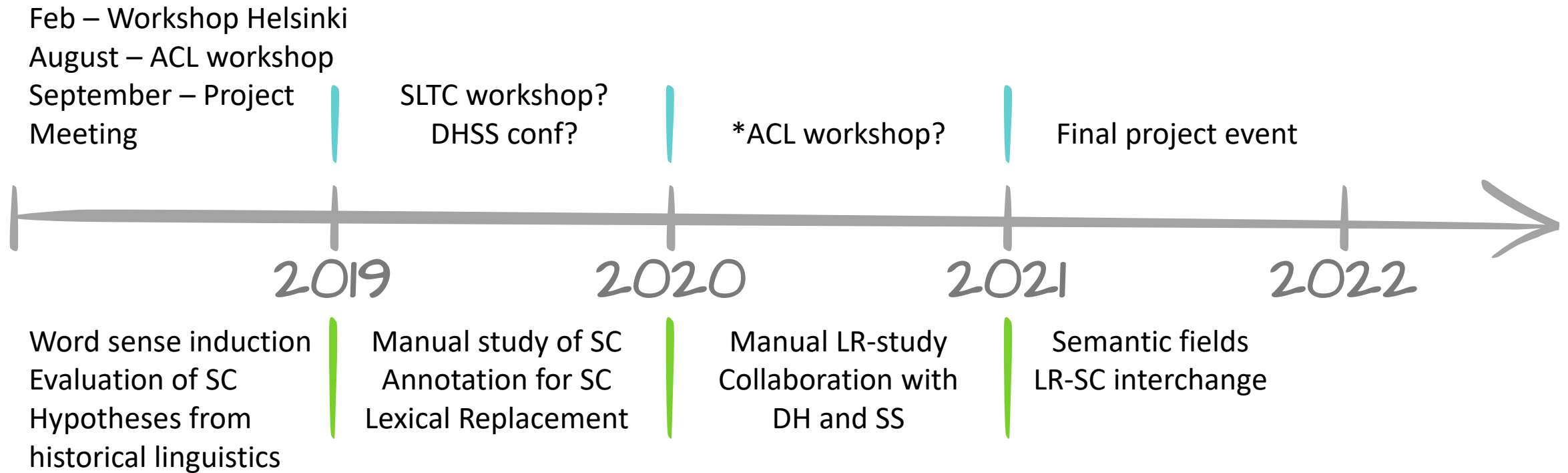
Talks (Stuttgart / Frankfurt spring 2019)



# Project timeline



**SC** = Semantic Change  
**LR** = Lexical Replacement  
**DH** = Digital Humanities  
**SS** = Social Science



# Conclusions



## Complexity in

- Multiple senses
- Many time points

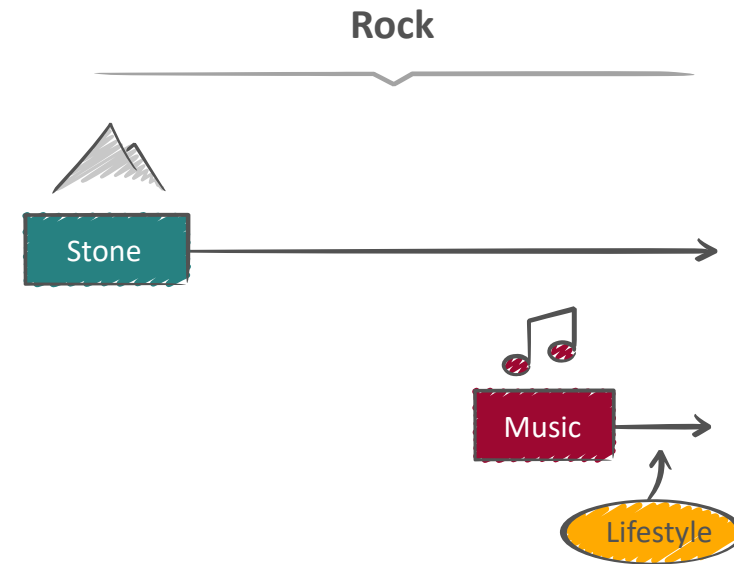


## Not all data are big data!



## Evaluation

- Common datasets and methods!
- What is the result valid for?



# Thank you for listening!



[Nina.tahmasebi@gu.se](mailto:Nina.tahmasebi@gu.se)

[nina@tahmasebi.se](mailto:nina@tahmasebi.se)