

Definition generation for lexical semantic change detection

Mariia Fedorova, Andrey Kutuzov, Yves Scherrer
Language Technology Group
University of Oslo

ACL Findings, 2024



Definitions as senses

- ▶ We use LLM-generated **contextualized word definitions** for **lexical semantic change detection (LSCD)**
 - ▶ ...by comparing vectorized definitions
 - ▶ ...by comparing the distributions of definitions ('definitions-as-senses')
- ▶ Dictionary-like definitions give us **interpretability**:
 - ▶ what **exactly** has changed in the meaning of a word over time?

Definitions as senses

- ▶ We use LLM-generated **contextualized word definitions** for **lexical semantic change detection (LSCD)**
 - ▶ ...by comparing vectorized definitions
 - ▶ ...by comparing the distributions of definitions ('definitions-as-senses')
- ▶ Dictionary-like definitions give us **interpretability**:
 - ▶ what **exactly** has changed in the meaning of a word over time?
- ▶ Still **reasonable performance on three different languages**.

Definitions as senses

- ▶ We use LLM-generated **contextualized word definitions** for **lexical semantic change detection (LSCD)**
 - ▶ ...by comparing vectorized definitions
 - ▶ ...by comparing the distributions of definitions ('definitions-as-senses')
- ▶ Dictionary-like definitions give us **interpretability**:
 - ▶ what **exactly** has changed in the meaning of a word over time?
- ▶ Still **reasonable performance on three different languages**.

Important prior work

- ▶ the general idea of **using definitions as representations** from [Giulianelli et al., 2023]

Definitions as senses

- ▶ We use LLM-generated **contextualized word definitions** for **lexical semantic change detection (LSCD)**
 - ▶ ...by comparing vectorized definitions
 - ▶ ...by comparing the distributions of definitions ('definitions-as-senses')
- ▶ Dictionary-like definitions give us **interpretability**:
 - ▶ what **exactly** has changed in the meaning of a word over time?
- ▶ Still **reasonable performance on three different languages**.

Important prior work

- ▶ the general idea of **using definitions as representations** from [Giulianelli et al., 2023]
- ▶ baseline **sense-based LSCD system** from [Tang et al., 2023]
 - ▶ (they used **pre-defined sense inventories**, no definition generation)

Definitions as senses

- ▶ We use LLM-generated **contextualized word definitions** for **lexical semantic change detection (LSCD)**
 - ▶ ...by comparing vectorized definitions
 - ▶ ...by comparing the distributions of definitions ('definitions-as-senses')
- ▶ Dictionary-like definitions give us **interpretability**:
 - ▶ what **exactly** has changed in the meaning of a word over time?
- ▶ Still **reasonable performance on three different languages**.

Important prior work

- ▶ the general idea of **using definitions as representations** from [Giulianelli et al., 2023]
- ▶ baseline **sense-based LSCD system** from [Tang et al., 2023]
 - ▶ (they used **pre-defined sense inventories**, no definition generation)
- ▶ **fine-tuned definition generation models** from [Kutuzov et al., 2024]

General pipeline

1. The task: to rank a given set of words by the degree of their semantic change in two time periods (represented by two corpora)

General pipeline

1. The task: to rank a given set of words by the degree of their semantic change in two time periods (represented by two corpora)
2. Let's find the senses of all the target word occurrences!
 - ▶ [Tang et al., 2023] did this by choosing a sense from an ontology (WSD)
 - ▶ we instead generate a definition from scratch for each target word occurrence...
 - ▶ ...assuming that different definitions will be generated for different senses.

General pipeline

1. The task: to rank a given set of words by the degree of their semantic change in two time periods (represented by two corpora)
 2. Let's find the senses of all the target word occurrences!
 - ▶ [Tang et al., 2023] did this by choosing a sense from an ontology (WSD)
 - ▶ we instead generate a definition from scratch for each target word occurrence...
 - ▶ ...assuming that different definitions will be generated for different senses.
- 2 sets of generated definitions (for 2 time periods). How to compare them for a change score?

General pipeline

1. The task: to rank a given set of words by the degree of their semantic change in two time periods (represented by two corpora)
 2. Let's find the senses of all the target word occurrences!
 - ▶ [Tang et al., 2023] did this by choosing a sense from an ontology (WSD)
 - ▶ we instead generate a definition from scratch for each target word occurrence...
 - ▶ ...assuming that different definitions will be generated for different senses.
- 2 sets of generated definitions (for 2 time periods). How to compare them for a change score?

Method 1. Definition embeddings

1. Embed the generated definitions into a vector space with, e.g., DistilRoBERTa

General pipeline

1. The task: to rank a given set of words by the degree of their semantic change in two time periods (represented by two corpora)
 2. Let's find the senses of all the target word occurrences!
 - ▶ [Tang et al., 2023] did this by choosing a sense from an ontology (WSD)
 - ▶ we instead generate a definition from scratch for each target word occurrence...
 - ▶ ...assuming that different definitions will be generated for different senses.
- 2 sets of generated definitions (for 2 time periods). How to compare them for a change score?

Method 1. Definition embeddings

1. Embed the generated definitions into a vector space with, e.g., DistilRoBERTa
2. Use well-developed LSCD methods of measuring distances between representations: PRT, APD, PRT/APD [Kutuzov et al., 2022]

General pipeline

1. The task: to rank a given set of words by the degree of their semantic change in two time periods (represented by two corpora)
 2. Let's find the senses of all the target word occurrences!
 - ▶ [Tang et al., 2023] did this by choosing a sense from an ontology (WSD)
 - ▶ we instead generate a definition from scratch for each target word occurrence...
 - ▶ ...assuming that different definitions will be generated for different senses.
- 2 sets of generated definitions (for 2 time periods). How to compare them for a change score?

Method 1. Definition embeddings

1. Embed the generated definitions into a vector space with, e.g., DistilRoBERTa
2. Use well-developed LSCD methods of measuring distances between representations: PRT, APD, PRT/APD [Kutuzov et al., 2022]
3. The only difference: definition embeddings instead of token embeddings

General pipeline

1. The task: to rank a given set of words by the degree of their semantic change in two time periods (represented by two corpora)
 2. Let's find the senses of all the target word occurrences!
 - ▶ [Tang et al., 2023] did this by choosing a sense from an ontology (WSD)
 - ▶ we instead generate a definition from scratch for each target word occurrence...
 - ▶ ...assuming that different definitions will be generated for different senses.
- 2 sets of generated definitions (for 2 time periods). How to compare them for a change score?

Method 1. Definition embeddings

1. Embed the generated definitions into a vector space with, e.g., DistilRoBERTa
2. Use well-developed LSCD methods of measuring distances between representations: PRT, APD, PRT/APD [Kutuzov et al., 2022]
3. The only difference: definition embeddings instead of token embeddings

Good performance, but interpretability is lost after vectorizing definitions.

Method 2. Definitions as strings



- ▶ Let's use **definitions in their textual form**
- ▶ Every definition is a 'sense'
- ▶ **definition distributions** for 2 time periods are compared similarly to how sense distributions are compared in [Tang et al., 2023]

- ▶ Let's use **definitions in their textual form**
- ▶ Every definition is a 'sense'
- ▶ **definition distributions** for 2 time periods are compared similarly to how sense distributions are compared in [Tang et al., 2023]
- ▶ These 'senses' are too granular, because of too many unique definitions. E.g., 'plane':
 - ▶ *'An aircraft, especially one designed for military use'*
 - ▶ *'An aircraft, especially a military aircraft'*, etc

Method 2. Definitions as strings

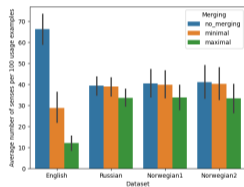


- ▶ Let's use **definitions in their textual form**
- ▶ Every definition is a 'sense'
- ▶ **definition distributions** for 2 time periods are compared similarly to how sense distributions are compared in [Tang et al., 2023]
- ▶ These 'senses' are too granular, because of too many unique definitions. E.g., 'plane':
 - ▶ '*An aircraft, especially one designed for military use*'
 - ▶ '*An aircraft, especially a military aircraft*', etc
- ▶ Hence, we **merge similar definitions** into one:
 - ▶ '**Minimalist merging**': only with the most frequent definition
 - ▶ '**Full-fledged merging**': any pair of definitions can be merged together
- ▶ Merging if Levenshtein edit distance exceeds a pre-defined threshold.

Method 2. Definitions as strings

- ▶ Let's use **definitions in their textual form**
- ▶ Every definition is a 'sense'
- ▶ **definition distributions** for 2 time periods are compared similarly to how sense distributions are compared in [Tang et al., 2023]
- ▶ These 'senses' are too granular, because of too many unique definitions. E.g., 'plane':
 - ▶ *'An aircraft, especially one designed for military use'*
 - ▶ *'An aircraft, especially a military aircraft', etc*
- ▶ Hence, we **merge similar definitions** into one:
 - ▶ **'Minimalist merging'**: only with the most frequent definition
 - ▶ **'Full-fledged merging'**: any pair of definitions can be merged together
- ▶ Merging if Levenshtein edit distance exceeds a pre-defined threshold.

Merging reduces the number of unique definitions: more realistic as senses.



Average number of senses per 100 usages before and after merging, calculated across all datasets for each language.

(should be done cautiously: merging too much can degrade the performance)

One English example



3 most frequent definitions/senses per time period for *ball* (high predicted change rate):

	Period 1 (1810-1860)	Period 2 (1960-2010)
<i>ball</i> Jensen-Shannon divergence (change score): 0.83	<i>A spherical object especially one that is round in shape (82%)</i>	<i>The object hit in a game (80%)</i>
	<i>A party (6%)</i>	<i>The object used in various sports especially in soccer tennis basketball etc (<1%)</i>
	<i>A wedding (<1%)</i>	<i>The object used in various sports especially in soccer basketball and other games which is thrown or kicked (<1%)</i>

One English example



3 most frequent definitions/senses per time period for *ball* (high predicted change rate):

	Period 1 (1810-1860)	Period 2 (1960-2010)
<i>ball</i>	<i>A spherical object especially one that is round in shape (82%)</i>	<i>The object hit in a game (80%)</i>
Jensen-Shannon divergence (change score): 0.83	<i>A party (6%)</i>	<i>The object used in various sports especially in soccer tennis basketball etc (<1%)</i>
	<i>A wedding (<1%)</i>	<i>The object used in various sports especially in soccer basketball and other games which is thrown or kicked (<1%)</i>

More examples, data and code at

<https://github.com/ltgoslo/Definition-generation-for-LSCD>

Performance across different generation and merging strategies



	English		Norwegian-1		Norwegian-2		Russian-1		Russian-2		Russian-3	
	Cosine	JS	Cosine	JS	Cosine	JS	Cosine	JS	Cosine	JS	Cosine	JS
No merging:												
Greedy	0.461	0.405	<i>0.303</i>	0.332	<i>0.211</i>	<i>0.232</i>	0.299	0.390	0.337	0.427	0.383	0.469
Beam	0.457	0.476	<i>0.268</i>	<i>0.238</i>	<i>0.216</i>	<i>0.201</i>	0.304	0.368	0.297	0.403	0.317	0.417
Diverse	0.449	0.382	<i>0.241</i>	<i>0.280</i>	<i>0.069</i>	<i>0.164</i>	0.301	0.345	0.310	0.389	0.348	0.421
Minimalist merging:												
Greedy	0.564	0.565	<i>0.251</i>	<i>0.280</i>	<i>0.192</i>	<i>0.197</i>	0.271	0.391	0.233	0.431	0.325	0.491
Beam	0.510	0.500	<i>0.297</i>	<i>0.240</i>	<i>0.112</i>	<i>0.189</i>	0.298	0.366	0.252	0.383	0.301	0.409
Diverse	0.478	0.434	<i>0.325</i>	<i>0.296</i>	<i>0.162</i>	<i>0.215</i>	0.265	0.354	0.268	0.406	0.287	0.443
Full-fledged merging:												
Greedy	0.417	0.418	<i>0.261</i>	0.362	<i>0.193</i>	0.260	0.286	0.391	0.250	0.416	0.360	0.476
Beam	0.492	0.493	<i>0.265</i>	<i>0.215</i>	<i>0.186</i>	<i>0.226</i>	0.304	0.360	0.250	0.347	0.327	0.420
Diverse	<i>0.312</i>	<i>0.301</i>	<i>0.209</i>	<i>0.315</i>	<i>0.202</i>	<i>0.221</i>	0.236	0.301	0.217	0.379	0.262	0.411
Threshold	50		10		10		10		10		10	

LSCD performance (Spearman's ρ) with different generation, merging and distance calculation strategies.

Threshold: Levenshtein edit distance threshold for merging definitions.

Results (Spearman's ρ for graded LSCD) compared to the baselines



Method	English	Norwegian-1	Norwegian-2	Russian-1	Russian-2	Russian-3
Non-interpretable methods:						
XLM-R token embeddings	0.514 \diamond	0.394 \diamond	0.387 \diamond	0.376 \diamond	0.480\diamond	0.457 \diamond
XL-LEXEME+APD (WiC-based)	0.886	0.659	0.640	0.796	0.820	0.863
Definition embeddings (ours)	0.637	0.496	0.565	0.488	0.462	0.504
Interpretable methods:						
Lesk without PoS	0.423 \clubsuit	<i>0.178</i>	0.500	0.294	0.279	0.286
Lesk with PoS	0.587	<i>0.150</i>	0.474	—	—	—
ARES sense embeddings	0.529 \clubsuit	—	—	—	—	—
LMMS sense embeddings	0.589 \clubsuit	—	—	—	—	—
Definitions as senses (ours)	0.565	0.362	<i>0.260</i>	0.391	0.431	0.491

(\diamond : best results from [Giulianelli et al., 2022], \clubsuit : [Tang et al., 2023], XL-LEXEME results from [Periti and Tahmasebi, 2024])

Definition embeddings yield better performance, but definitions as senses are interpretable and explainable.



Giulianelli, M., Kutuzov, A., and Pivovarova, L. (2022).

Do not fire the linguist: Grammatical profiles help language models detect semantic change.



In Tahmasebi, N., Montariol, S., Kutuzov, A., Hengchen, S., Dubossarsky, H., and Borin, L., editors, Proceedings of the 3rd Workshop on Computational Approaches to Historical Language Change, pages 54–67, Dublin, Ireland. Association for Computational Linguistics.



Giulianelli, M., Luden, I., Fernandez, R., and Kutuzov, A. (2023).

Interpretable word sense representations via definition generation: The case of semantic change analysis.

In Rogers, A., Boyd-Graber, J., and Okazaki, N., editors, Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 3130–3148, Toronto, Canada. Association for Computational Linguistics.

-  Kutuzov, A., Fedorova, M., Schlechtweg, D., and Arefyev, N. (2024). Enriching word usage graphs with cluster definitions. In Calzolari, N., Kan, M.-Y., Hoste, V., Lenci, A., Sakti, S., and Xue, N., editors, Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 6189–6198, Torino, Italia. ELRA and ICCL.
-  Kutuzov, A., Velldal, E., and Øvrelid, L. (2022). Contextualized embeddings for semantic change detection: Lessons learned. In Derczynski, L., editor, Northern European Journal of Language Technology, Volume 8, Copenhagen, Denmark. Northern European Association of Language Technology.



Periti, F. and Tahmasebi, N. (2024).

A systematic comparison of contextualized word embeddings for lexical semantic change. In Duh, K., Gomez, H., and Bethard, S., editors, Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 4262–4282, Mexico City, Mexico. Association for Computational Linguistics.



Tang, X., Zhou, Y., Aida, T., Sen, P., and Bollegala, D. (2023).

Can word sense distribution detect semantic changes of words?

In Bouamor, H., Pino, J., and Bali, K., editors, Findings of the Association for Computational Linguistics: EMNLP 2023, pages 3575–3590, Singapore. Association for Computational Linguistics.