

SST-BERT at SemEval-2020 Task 1: Semantic Shift Tracing by Clustering in BERT-based Embedding Spaces

Vani K Sandra Mitrović Alessandro Antonucci Fabio Rinaldi
IDSIA, Lugano, Switzerland
{vanik, sandra, alessandro, fabio.rinaldi}@idsia.ch

Abstract

Lexical semantic change detection (also known as *semantic shift tracing*) is a task of identifying words that have changed their meaning over time. Unsupervised semantic shift tracing, focal point of SemEval2020, is particularly challenging. Given the unsupervised setup, in this work, we propose to identify clusters among different occurrences of each target word, considering these as representatives of different word meanings. As such, disagreements in obtained clusters naturally allow to quantify the level of semantic shift per each target word in four target languages. To leverage this idea, clustering is performed on contextualized (BERT-based) embeddings of word occurrences. The obtained results show that our approach performs well both measured separately (per language) and overall, where we surpass all provided SemEval baselines.

1 Problem Setup

Consider two corpora \mathcal{C}_1 and \mathcal{C}_2 for a same language but associated with different time stamps (say, respectively, t_1 and $t_2 > t_1$). Let \mathcal{W} be a set of *target* words occurring in both corpora. Each target word $w \in \mathcal{W}$ might assume multiple meanings, to be called *senses*, within the two corpora. A pool of experts annotated a representative amount of occurrences with their corresponding senses. The problem we consider is to characterize the semantic shift related to those senses from one corpus to the other without having access to the expert annotations. In particular, we address the two following two subtasks:

- **Subtask 1:** Decide, for each $w \in \mathcal{W}$, whether or not w gained or lost at least a *sense* between t_1 and t_2 . This is a binary decision task. We will denote this subtask as **(S1)**.
- **Subtask 2:** Define, for the elements of \mathcal{W} , a measure of their degree of lexical semantic change between t_1 and t_2 and sort these elements consequently. This is a ranking task. This subtask will be referred to as **(S2)**.

We¹ describe two different methods able to address both subtasks. As both methods require a preprocessing step based on transformers, let us start from this preliminary operation.

2 Preprocessing

For the proposed approach, we used embeddings derived from BERT (Bidirectional Encoder Representations from Transformers) (Devlin et al., 2018) model to represent the text information in the corpora. BERT uses attention mechanism to learn the contextual relations and reads the input bidirectionally. It is an encoder-only model (as the goal is to generate a language model) opposed to transformers (encoder-decoder model)(Vaswani et al., 2017). BERT is trained on two main objectives, masked language model (MLM) and next sentence prediction (NSP).

The corpus data is initially segmented at sentence-level and the BERT Word Piece tokenizer is applied on

This work is licensed under a Creative Commons Attribution 4.0 International License. License details: <http://creativecommons.org/licenses/by/4.0/>.

¹Our team name in SemEval2020 competition is NLP@IDSIA.

these sentences, to get the token-level representations. BERT-base model with twelve transformer layers is used and we derived the final embedding by concatenating the final four layers. If a single word gets split by the tokenizer, we take the average embedding value of the sub-tokens. Thus, for each target word we extract the embeddings from all the sentences associated with it from both corpora, \mathcal{C}_1 and \mathcal{C}_2 . For corpora other than English language, we used multilingual BERT models of respective languages. We used the pre-trained model to generate the embeddings for all experiments, since the task is completely unsupervised in nature.

Let us assume, for each $w \in \mathcal{W}$, that $\mathcal{S}_j := \{s_{ij}^{(w)}\}_{i=1}^{n_j^{(w)}}$ denote the $n_j^{(w)}$ sentences in corpus \mathcal{C}_j where target word w occurred, for each $j = 1, 2$. The transformer maps those sentences into corresponding vectors of a d -dimensional space. Let us perform a single transformation for both \mathcal{S}_1 and \mathcal{S}_2 simultaneously. Denote as $x_{ij}^{(w)} \in \mathbb{R}^d$ the vector associated with $s_{ij}^{(w)}$, for each $i = 1, \dots, n_{w,j}$ and $j = 1, 2$. We regard the relative distances between these vectors as proxy indicators of their semantic similarity. This helps to cluster the similar senses together. While different distances, such as cosine, Manhattan and Euclidean could be used to measure BERT embedding similarities, we adopt the Euclidean distance given the findings in (Hewitt and Manning, 2019; Inui et al., 2019; Reif et al., 2019) which demonstrate that the syntax tree distance between two words in BERT embedding space corresponds to the square of the Euclidean distance. In the next two sections, in order to address subtasks (S1) and (S2), we focus on the vectors associated with the same target word w and use their relative Euclidean distances as indicators of possible semantic shifts from one corpus to the other.

Both methods will be based on clustering algorithms used to cluster the vectors associated with a given target word of a single corpus or of the union of the two. We adopt the classical k -means clustering algorithm, which forms the clusters by attempting to minimize the intra-cluster variance. So called *silhouette* method is used for the selection of the optimal number of clusters k and the initialization (i.e., the position of the centroids before starting the algorithm) (Rousseeuw, 1987). Accordingly, given a value of k , we compute in the corresponding cluster, the means of both the nearest-inter cluster distance and the nearest-cluster distance. The difference between these two quantities normalized by the maximum of the two is used as a fitness score to be maximized in order to select the optimal value of k . The same approach is used to determine the initial centroids. In this case, we use the optimal k value and run the k -means algorithm for N iterations, to determine the best centroids.

3 Method 1: Joint Clustering Vectors of Both Corpora

Let us focus on a particular target word $w \in \mathcal{W}$. Accordingly, for the sake of readability, denote its vectors in corpus \mathcal{C}_j simply as $\mathcal{X}_j := \{x_{ij}\}_{i=1}^{n_j}$, for each $j = 1, 2$. We cluster the whole set of vectors of the two corpora, say $\mathcal{X} := \mathcal{X}_1 \cup \mathcal{X}_2$, and denote as $\{\mathcal{X}^k\}_{k=1}^m$ the n clusters returned by the algorithm. Note that we cope with *hard* clustering methods, i.e., $\cup_{k=1}^m \mathcal{X}^k = \mathcal{X}$ and $\mathcal{X}^{k_1} \cap \mathcal{X}^{k_2} = \emptyset$ for each $k_1, k_2 = 1, \dots, m$, with $k_1 \neq k_2$. For each cluster \mathcal{X}^k we count how many of its elements belong to \mathcal{X}_1 , say $n_{1,k}$, and to \mathcal{X}_2 , say $n_{2,k}$. We call *impure* a cluster such that both $n_{1,k} > 0$ and $n_{2,k} > 0$. As we regard the clusters as equivalence classes for the abstract notion of *sense*, if all the m clusters are impure it means that no new senses appeared in \mathcal{C}_2 and no new senses have been lost from \mathcal{C}_1 to \mathcal{C}_2 . If this is not the case we might have new senses in \mathcal{C}_2 , i.e., there is at least a k such that $n_{1,k} = 0$, or, vice versa, an old sense has been lost, i.e., there is at least a k such that $n_{2,k} = 0$. Following the guidelines of the SemEval shared task, we might set a lower bound \underline{n} to the number of occurrences of a word in a cluster before deciding to regard it as a new sense. If this is the case the above conditions for the counts equal to zero should be replaced by $n_{j,k} < \underline{n}$. Overall, this procedure corresponds to a sound algorithm to address (S1). We refer to it as MIS1.

Regarding (S2), after the clustering, we might define a random variable S , to be called the *sense* variable, whose m states are in one-to-one correspondence with the clusters. The variable denotes how likely is finding an occurrence of w with sense S in a corpus. Accordingly, we might use the counts $\{n_{j,k}\}_{k=1}^m$ to learn a probability mass function $P_j(S)$ for each $j = 1, 2$. Following a Bayesian approach,

based on a Laplace uniform prior with equivalent sample size $\sigma > 0$ (Gelman et al., 2013), we have:

$$P(s_k) = \frac{n_{j,k} + \frac{\sigma}{m}}{\sum_{k=1}^m n_{j,k} + \sigma}. \quad (1)$$

In such a probabilistic setup, the semantic shift of the target word w between the two corpora can be therefore described by the dissimilarity between the mass functions $P_1(S)$ and $P_2(S)$. We measure that by the Shannon-Jensen distance SJ , i.e., a symmetrization of the popular Kullback-Leibler divergence. This semantic shift of w corresponds therefore to the distance δ , with $\delta := SJ(P_1, P_2)$ and $SJ(P_1, P_2) := \frac{1}{2}[KL(P_1, P_2) + KL(P_2, P_1)]$ and $KL(P_1, P_2) := \sum_{k=1}^m P_1(s_k) \ln \frac{P_1(s_k)}{P_2(s_k)}$. Note that with the Bayesian smoothing in Equation (1), we cannot have zero probabilities and degenerate values in the computation of the distance. The overall procedure gives an algorithm to address subtask S2, as this corresponds to sort the elements of \mathcal{W} with respect to their value δ . We refer to this procedure as M1S2.

4 Method 2: Separate Clustering of the two Corpora

In this section, while still focusing on a given target word $w \in \mathcal{W}$, we consider a different approach based on the separate clustering of the two set of vectors \mathcal{X}_1 and \mathcal{X}_2 . Let $\{\mathcal{X}_1^{k_1}\}_{k_1=1}^{m_1}$ and $\{\mathcal{X}_2^{k_2}\}_{k_2=1}^{m_2}$ denote these two sets of clusters. As in the previous method we regard each cluster as a representative model of a sense. Yet, unlike the previous case, here we need to define a map between the clusters of the first corpus and those of the second. As discussed before we adopt the Euclidean distance between the vectors as a proxy indicator of semantic similarity. In order to cope with single numerical values, for the sake of simplicity, we represent each cluster with its center of mass. Let $\tilde{x}_1^{k_j}$ denote the center of mass of $\mathcal{X}_j^{k_j}$ for each $k_j = 1, \dots, m_j$ and $j = 1, 2$. If $m_1 = m_2$, i.e., the two corpora have the same number of clusters, we can reduce the identification of the map between the clusters of the two corpora to a minimum weight matching in a complete bipartite graph, whose nodes are associated with the two sets of clusters and whose weights are the Euclidean distances between the centers of mass. If this is not the case and, for instance, $m_1 > m_2$, we add $m_1 - m_2$ *dummy* clusters to the second corpus and set to zero the weights for all the arcs connecting these elements. We similarly proceed if $m_2 > m_1$.

The optimal matching minimizing the sum of the weights can be computed in cubic time with the classical Hungarian algorithm (Kuhn, 1955; Jonker and Volgenant, 1987) and the results is a one-to-one correspondence between the clusters, no matter whether proper or dummy, of the two corpora. As a dummy cluster in a corpus has zero distance from all the clusters of the other corpus, the matching returned by the Hungarian algorithm is properly minimizing the distance between the proper clusters. Two proper clusters in the two corpora matched by the algorithm are intended as representative of the same sense. Proper clusters of a corpus pointing to dummy cluster are regarded instead as a new sense appeared in the second corpus only, or old sense occurred in the first corpus only.

After the matching, we define a single clustering with $m := \max\{m_1, m_2\}$ clusters and proceed exactly as in the previous section. In practice, the vectors of two clusters matched by the Hungarian algorithm are assigned to a single, impure, cluster, while those linked to dummy clusters produce pure clusters. We term M2S1 and M2S2 the two algorithms corresponding to the approach discussed in this section to address the two subtasks. Next section describes the experimental analysis and evaluation results.

5 Method 3: An alternative approach for Subtask 2 (S2)

As an alternative to the previously explained procedure for handling (S2), based on Bayesian approach and Shannon-Jensen divergence, we consider another approach, exploiting only the number of word occurrences per cluster and corpora. More precisely, assuming that we have \mathcal{K} clusters in total and $\forall k \in \{1, \dots, \mathcal{K}\}$ already calculated $n_{1,k}$ and $n_{2,k}$ from each corpora (regardless whether clusters come from single clustering in M1 or after performing optimal cluster matching in M2), we define the coefficient of

semantic change of the word (ranking in (S2) terminology), as:

$$\frac{1}{2S_1S_2} \sum_{k=1}^{\mathcal{K}} |S_2 \cdot n_{1,k} - S_1 \cdot n_{2,k}|$$

where $S_1 = \sum_{k=1}^{\mathcal{K}} n_{1,k}$ and $S_2 = \sum_{k=1}^{\mathcal{K}} n_{2,k}$. Let us assume that word w has p occurrences in both corpora. It is trivial to see that in the case with $\mathcal{K} = 2$ and clear cut between corpora (e.g. all occurrences in cluster 1 belong to \mathcal{C}_2 and all occurrences in cluster 2 belong to \mathcal{C}_1 , i.e. $n_{2,1} = n_{1,2} = p, n_{1,1} = n_{2,2} = 0$), our coefficient equals 1, which indicates complete change of sense. Likewise, if the distribution of occurrences is uniform ($n_{*,*} = p/2$), it yields 0, meaning no sense change. We denote these two new procedures for (S2) for M1 and M2 as NM1 and NM2, respectively.

6 Experimental Analysis

Experimental analysis is performed according to the rules posed by SemEval2020 challenge² organizers, using provided corpora (2) and baselines (3). Corpora are provided in four languages: English (Alatrash et al., 2020), Latin (McGillivray and Kilgarriff, 2013), German (Textarchiv, 2018) and Swedish (Adesam et al., 2019). Table 1 provides brief statistics for given corpora stating the number of target words (NTW), the total and average number of sentences containing target words (NSTW) per each language. The three baselines provided are: normalized frequency difference (FD), count vectors with column intersection and cosine distance (CNT+CI+CD) and a random baseline always predicting a majority class (RND/MC) - for details see (Schlechtweg et al., 2019).

For evaluation purposes (as instructed by SemEval guidelines), accuracy is exploited for (S1), while Spearman coefficient, taking values between -1 (corresponding to negative correlation) and 1 (perfect correlation), was used for (S2). More details can be found in the system description paper (Schlechtweg et al., 2020).

It is worth mentioning that the upper bound for the number of clusters for K-means which could be retrieved by the silhouette score was set to 10.

As explained before, for (S1), the idea was to compare the number of elements of each cluster \mathcal{X}^k coming from different corpora, say $n_{1,k}$ and $n_{2,k}$, and claim a change in senses if $\exists k: n_{1,k} = 0 \vee n_{2,k} = 0$. This indeed was the procedure applied for Latin corpora. For other languages (with larger sizes of corpora), following the guidelines of the SemEval shared task, additional restrictions in terms of lower (\underline{n}) and upper bounds (\bar{n}) were set, with the following purpose: word is considered as gaining a new sense, if $\exists k: n_{1,k} \leq \bar{n} \wedge n_{2,k} \geq \underline{n}$ (and vice versa for losing a sense). Additionally, suggested values for these bounds were set to $\underline{n} = 5$ and $\bar{n} = 2$.

The code³ is implemented in Python using Scikit (Buitinck et al., 2013) and Transformers (Wolf et al., 2019) library.

Language	NTW	\mathcal{C}_1			\mathcal{C}_2		
		total (NSTW)	average (NSTW)	span (years)	total (NSTW)	average (NSTW)	span (years)
English	37	24917	673.43	1810-1860	29088	786.16	1960-2010
Latin	40	26912	672.80	-200-0	126081	3152.03	0-2000
German	48	78845	1642.60	1800-1899	68621	1429.60	1946-1990
Swedish	31	83703	2700.10	1790-1830	241525	7791.13	1895-1903

Table 1: Summary statistics of corpora with respect to target words for different languages

The experimental results on the corpora over the two subtasks (S1) and (S2) are reported using the methods M1 and M2, in Table 2. Results are shown for the four target languages and overall, as well as

²https://competitions.codalab.org/competitions/20948#learn_the_details-overview

³The source code is available at: https://github.com/vanikanjirangat/SST_BERT_SEMEVAL_TASK1

compared with provided SemEval baselines. Best results per language and subtask are underlined. Overall best results per subtask are denoted in boldface. As can be seen, except for the Latin, the proposed methods are outperforming all baselines on (S1) with the procedure M1S1 being the best for German and Swedish and M2S1 for English. On the other hand, for (S2), results are quite corpus/language dependent, M1S2 scores best for English, M2S2 for Swedish, while baseline 2 (CNT+CI+CD) wins over all the others for Latin and German. Overall, Method 2 outperforms its competitors on both subtasks. The performance with the contextualized embeddings is actually comparable with the baseline approaches in many cases. This could be the fact that pre-trained embeddings from BERT may not be completely suitable for representing meaningful sentence vectors for clustering (Reimers and Gurevych, 2019). These factors have to be investigated in the future.

Language	Method 1		Method 2		SemEval Baselines					
					FD		CNT+CI+CD		MC	
	(S1)	(S2)	(S1)	(S2)	(S1)	(S2)	(S1)	(S2)	(S1)	(S2)
English	0.541	<u>0.028</u>	<u>0.622</u>	-0.008	0.432	-0.217	0.595	0.022	0.568	/
Latin	0.375	0.253	0.625	0.253	<u>0.650</u>	0.020	0.525	<u>0.359</u>	0.350	/
German	<u>0.708</u>	0.176	0.625	0.176	0.417	0.014	0.688	<u>0.216</u>	0.646	/
Swedish	<u>0.742</u>	0.275	0.677	<u>0.321</u>	0.258	-0.15	0.645	-0.022	<u>0.742</u>	/
Overall	0.591	0.183	0.637	0.185	0.439	-0.083	0.613	0.144	0.576	/

Table 2: Performance of the proposed methods and different SemEval baselines. The baselines correspond to normalized frequency difference (FD), count vectors with column intersection and cosine distance (CNT+CI+CD) and majority class(MC)

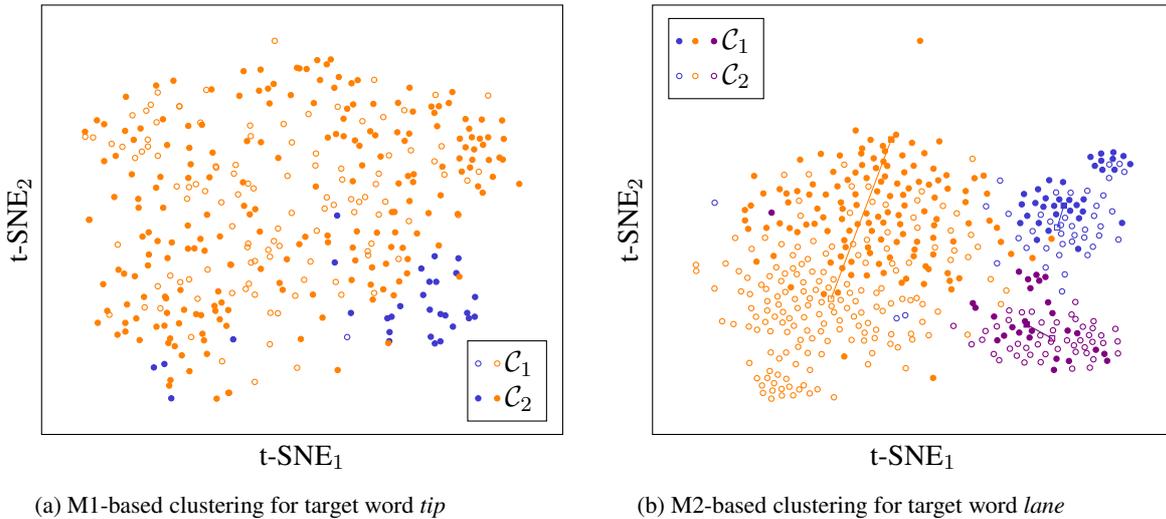


Figure 1: Obtained clusterings for target words

Figure 1a shows an example of the application of method M1 for the English target word *tip*, based on 2D t-SNE (Maaten and Hinton, 2008) projections. The method produces two clusters, each denoted with a different color. As it can be seen, one of the clusters (orange, $k = 1$) is remarkably larger than the other (blue, $k = 2$). Additionally, it is also quite impure containing word occurrences from both corpora (more precisely, $n_{1,1} = 112$ and $n_{2,1} = 211$), while the other cluster contains only 31 instance whose distribution is $n_{1,2} = 1$ and $n_{2,2} = 30$. Given that $n_{1,2} = 1 \leq 2 = \bar{n}$ and $n_{2,2} = 30 \geq 5 = \underline{n}$, method M1S1 correctly detects the sense change for the word *tip*.

Figure 1b shows an example of the application of method M2 for the English target word *lane* (2D t-SNE projections). The clustering algorithms produce three clusters for each corpus, each denoted with a different color, and the matching algorithm detect the correspondence between clusters minimizing the

distances between the centers of mass, depicted as squares in the figure. Matching clusters of the two corpora are depicted with the same color.

The results of the alternative method M3 for subtask (S2) with respect to M1 and M2 and two baselines (FD and CNT+CI+CD), are provided in Table 3. We can see that NM1 improves results on English and German languages, and overall.

	M1	M2	SemEval Baselines		NM1	NM2
			FD	CNT+CI+CD		
English	0.028	-0.008	-0.217	0.022	<u>0.159</u>	0.037
Latin	0.253	0.253	0.020	<u>0.359</u>	0.231	0.333
German	0.176	0.176	0.014	0.216	<u>0.525</u>	0.062
Swedish	0.275	<u>0.321</u>	-0.15	-0.022	0.141	0.095
Overall	0.183	0.185	-0.083	0.144	0.264	0.132

Table 3: Performance of the semantic change coefficient methods NM1 & NM2 for (S2)

7 Related Work

The task of identifying words whose meaning has changed over time is well-known and the related literature is, therefore, resourceful with many recent advancements. Albeit, there are still quite some issues to be resolved, primary regarding the methodology and the respective ground truth (missing semantic change annotations).

As for the latter, the first step is to decide whether to aim for a binary response (equivalent of SemEval2020 subtask 1) or to provide graded ratings of a sense change (equivalent of SemEval2020 subtask 2). Given that in both cases, but particularly with graded rating, inter-annotator agreement rates vary greatly, as evidenced in (Erk et al., 2009), establishing a definition of a standard test set is extremely difficult. In (Schlechtweg et al., 2018) a unifying evaluation framework for unsupervised lexical semantic change detection was proposed based on changes in relatedness of word use pairs in each time period.

Regarding the former, many different approaches for unsupervised lexical semantic change detection have been suggested. A detailed survey of studies can be found in (Kutuzov et al., 2018; Tahmasebi et al., 2018). Most notably, several works (Baroni et al., 2014; Kim et al., 2014; Hamilton et al., 2016) showed the benefits of using dense word representations for semantic shift detection. Furthermore, (Kulkarni et al., 2015) showcased that these outperform the frequency-based methods. However, unlike our work, none of these works exploits clustering. More close to our approach, that is, considering clusters as representative semantic areas, are the works of (Mitra et al., 2014) and (Dubossarsky et al., 2015). The main differences are however, that in (Mitra et al., 2014), clustering is performed on the level of the ego-network of each word, where the network is constructed based on word co-occurrences, while we perform clustering of the word embeddings itself. Additionally, in contrast to (Dubossarsky et al., 2015) where the authors consider incremental learning of word embeddings in yearly chunks and vary the number of clusters from 500 to 5000, we use silhouette scores to determine the optimal number of clusters per each target word.

8 Conclusion and Future Work

A word can have a different meaning (sense) in different contexts and/or different time periods. Despite being quite extensively studied, the problem of identifying words that have changed their meaning over time, particularly in an unsupervised way, still challenges researchers.

In this work, we propose two approaches, both of which combine contextualized word embeddings (obtained by BERT) and clustering, differing thus only in the way the clustering has been performed. Considering obtained clusters as proxies for word meanings allows us to quantify the level of change per each target word in four target languages. The obtained results, especially looking overall, across all

target languages, where we are outperforming all provided baselines, demonstrate the usefulness of the suggested approach.

As potential directions for future work we plan to investigate various strategies, including different clustering methods and time-wise comparison of target words nearest-neighbours, in an attempt to identify actual word senses more accurately. Additionally, we would like to further scrutinize how the linguistic particularities of different corpora might have contributed to the variability of the results.

References

- Yvonne Adesam, Dana Dannélls, and Nina Tahmasebi. 2019. Exploring the quality of the digital historical newspaper archive kubhist. In *DHN*, pages 9–17.
- Reem Alatrash, Dominik Schlechtweg, Jonas Kuhn, and Sabine Schulte im Walde. 2020. Ccoha: Clean corpus of historical american english. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 6958–6966.
- Marco Baroni, Georgiana Dinu, and Germán Kruszewski. 2014. Don’t count, predict! a systematic comparison of context-counting vs. context-predicting semantic vectors. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 238–247.
- Lars Buitinck, Gilles Louppe, Mathieu Blondel, Fabian Pedregosa, Andreas Mueller, Olivier Grisel, Vlad Niculae, Peter Prettenhofer, Alexandre Gramfort, Jaques Grobler, Robert Layton, Jake VanderPlas, Arnaud Joly, Brian Holt, and Gaël Varoquaux. 2013. API design for machine learning software: experiences from the scikit-learn project. In *ECML PKDD Workshop: Languages for Data Mining and Machine Learning*, pages 108–122.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Haim Dubossarsky, Yulia Tsvetkov, Chris Dyer, and Eitan Grossman. 2015. A bottom up approach to category mapping and meaning change. In *NetWordS*, pages 66–70.
- Katrin Erk, Diana McCarthy, and Nicholas Gaylord. 2009. Investigations on word senses and word usages. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 1-Volume 1*, pages 10–18. Association for Computational Linguistics.
- Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin. 2013. *Bayesian data analysis*. CRC press.
- William L Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Diachronic word embeddings reveal statistical laws of semantic change. *arXiv preprint arXiv:1605.09096*.
- John Hewitt and Christopher D Manning. 2019. A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4129–4138.
- Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan. 2019. Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (emnlp-ijcnlp). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*.
- Roy Jonker and Anton Volgenant. 1987. A shortest augmenting path algorithm for dense and sparse linear assignment problems. *Computing*, 38(4):325–340.
- Yoon Kim, Yi-I Chiu, Kentaro Hanaki, Darshan Hegde, and Slav Petrov. 2014. Temporal analysis of language through neural language models. *arXiv preprint arXiv:1405.3515*.
- Harold W Kuhn. 1955. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97.
- Vivek Kulkarni, Rami Al-Rfou, Bryan Perozzi, and Steven Skiena. 2015. Statistically significant detection of linguistic change. In *Proceedings of the 24th International Conference on World Wide Web*, pages 625–635.
- Andrey Kutuzov, Lilja Øvrelid, Terrence Szymanski, and Erik Velldal. 2018. Diachronic word embeddings and semantic shifts: a survey. *arXiv preprint arXiv:1806.03537*.

- Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605.
- Barbara McGillivray and Adam Kilgarriff. 2013. Tools for historical corpus research, and a corpus of latin. *New Methods in Historical Corpus Linguistics*, (3):247–257.
- Sunny Mitra, Ritwik Mitra, Martin Riedl, Chris Biemann, Animesh Mukherjee, and Pawan Goyal. 2014. That’s sick dude!: Automatic identification of word sense change across different timescales. *arXiv preprint arXiv:1405.4392*.
- Emily Reif, Ann Yuan, Martin Wattenberg, Fernanda B Viegas, Andy Coenen, Adam Pearce, and Been Kim. 2019. Visualizing and measuring the geometry of bert. In *Advances in Neural Information Processing Systems*, pages 8592–8600.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3973–3983.
- Peter J Rousseeuw. 1987. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65.
- Dominik Schlechtweg, Sabine Schulte im Walde, and Stefanie Eckmann. 2018. Diachronic usage relatedness (durel): A framework for the annotation of lexical semantic change. *arXiv preprint arXiv:1804.06517*.
- Dominik Schlechtweg, Anna Hättü, Marco Del Tredici, and Sabine Schulte im Walde. 2019. A wind of change: Detecting and evaluating lexical semantic change across times and domains. *arXiv preprint arXiv:1906.02979*.
- Dominik Schlechtweg, Barbara McGillivray, Simon Hengchen, Haim Dubossarsky, and Nina Tahmasebi. 2020. SemEval-2020 Task 1: Unsupervised Lexical Semantic Change Detection. In *To appear in Proceedings of the 14th International Workshop on Semantic Evaluation*, Barcelona, Spain. Association for Computational Linguistics.
- Nina Tahmasebi, Lars Borin, and Adam Jatowt. 2018. Survey of computational approaches to lexical semantic change. *arXiv preprint arXiv:1811.06278*.
- Deutsches Textarchiv. 2018. Grundlage für ein referenzkorpus der neuhochdeutschen sprache. *Herausgegeben von der Berlin-Brandenburgischen Akademie der Wissenschaften, Berlin*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, R’emi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *ArXiv*, abs/1910.03771.