

---

# A SURVEY ON CONTEXTUALISED SEMANTIC SHIFT DETECTION

---

A PREPRINT

✉ Stefano Montanelli\*

Department of Computer Science  
Università degli Studi di Milano  
Via Celoria 18, 20133 Milan, Italy  
stefano.montanelli@unimi.it

✉ Francesco Periti†

Department of Computer Science  
Università degli Studi di Milano  
Via Celoria 18, 20133 Milan, Italy  
francesco.periti@unimi.it

## ABSTRACT

Semantic Shift Detection (SSD) is the task of identifying, interpreting, and assessing the possible change over time in the meanings of a target word. Traditionally, SSD has been addressed by linguists and social scientists through manual and time-consuming activities. In the recent years, computational approaches based on Natural Language Processing and word embeddings gained increasing attention to automate SSD as much as possible. In particular, over the past three years, significant advancements have been made almost exclusively based on word contextualised embedding models, which can handle the multiple usages/meanings of the words and better capture the related semantic shifts. In this paper, we survey the approaches based on contextualised embeddings for SSD (i.e., CSSDetection) and we propose a classification framework characterised by *meaning representation*, *time-awareness*, and *learning modality* dimensions. The framework is exploited i) to review the measures for shift assessment, ii) to compare the approaches on performance, and iii) to discuss the current issues in terms of scalability, interpretability, and robustness. Open challenges and future research directions about CSSDetection are finally outlined.

**Keywords** Computational Semantics · Contextualised Word Embeddings · Semantic Shift Detection

## 1 Introduction

Word meanings in a language are influenced over time by social practices, events, and political circumstances [Keidar et al., 2022, Castano et al., 2022, Azarbyonad et al., 2017]. Detecting, interpreting, and assessing the possible change of a word over time are usually considered as steps of a broader *Semantic Shift Detection* (SSD) task. Capturing semantic shift requires to arrange testing procedures as well as to define and standardise interviews that are eventually exploited to build large catalogues of word descriptions. All this work is generally addressed by interested scholars, like linguists and social scientists, through manual and time-consuming approaches of “close reading” that keep humans “in-the-loop”.

The growing attention on Computational Semantics issues as well as the recent availability of large digitised diachronic corpora in many different languages, like English [Alatrash et al., 2020], Swedish [Adesam et al., 2019], German [Schlechtweg et al., 2020], Latin [McGillivray and Kilgarriff, 2013], Italian [Basile et al., 2019], Russian [Kutuzov and Pivovarova, 2021a], Spanish [Zamora-Reina et al., 2022], Chinese [Chen et al., 2022], and Norwegian [Kutuzov et al., 2022a], pushed the emergence of a novel family of approaches based on Natural Language Processing (NLP) techniques to automate the SSD task as much as possible.

In this context, distributional word representations (i.e., word embeddings) emerged as an effective solution, based on the idea that semantically related words are close to each other in the embedding space [Mikolov et al., 2013]. The use of static embedding solutions is widely adopted and the main approaches have been reviewed in three survey papers [Tang, 2018, Kutuzov et al., 2018, Tahmasebi et al., 2021]. Typically, static embeddings are used to detect how a word changes

---

\*Authors are listed in alphabetical order

†Corresponding author

in its dominant sense, without considering the possible additional, subordinate senses/meanings that the word can have. However, subordinate senses can change on their own regardless of their dominant sense. For example, considering the word *rock*, the *music* meaning evolved over time to encompass both *music* and a particular lifestyle, while the *stone* meaning remained unchanged [Tahmasebi, 2013, Mitra et al., 2015]. This issue has motivated the recent efforts to enforce SSD by leveraging contextualised embeddings, which are capable of handling the so-called *colexification* phenomena such as homonymy and polysemy. As a result, SSD solutions based on contextualised embeddings have emerged, but a classification framework and a corresponding survey of existing approaches are still missing.

In this paper, we define CSSDetection as a *Contextualised Semantic Shift Detection* task and we survey the main approaches to SSD addressed through the use of contextualised embedding techniques. To this end, we propose a classification framework based on three dimensions of analysis, namely *meaning representation*, *time-awareness*, and *learning modality*, that allows to effectively describe the featuring properties of both *form-* and *sense-oriented* approaches in which CSSDetection solutions are typically distinguished. Assessment methods and metrics used for CSSDetection are also surveyed to discuss how the detected semantic shift of a word is measured and quantified by the considered approaches.

In the recent years, a growing attention has been captured by SSD issues, and a number of events with competitive shared tasks and corpora have been proposed, such as SemEval-20 Task 1 [Schlechtweg et al., 2020], DIACRIta-20 [Basile et al., 2020], RuShiftEval-21 [Kutuzov and Pivovarova, 2021b], and LSCDiscovery-22 [Zamora-Reina et al., 2022]. As a further contribution of our survey, the CSSDetection approaches are compared according to their results in these competitions (where available) with the aim to discuss the related performance issues and possible limitations in real applications.

Unlike the existing surveys on static SSD, the goal of our survey is to focus on contextualised approaches and to highlight the computational perspective, rather than the linguistics one.

The paper is organised as follows. Section 2 presents the CSSDetection problem and the related workflow with related formalisation. The proposed survey framework for approach classification is illustrated in Section 3. The classification of state-of-the-art approaches is discussed in Section 4. A comparative analysis of approach performance is provided in 5; issues about scalability, interpretability, and robustness of CSSDetection approaches are discussed in Section 6. Finally, in Section 7, we outline the open challenges and we give our concluding remarks.

## 2 Problem statement

Consider a diachronic document corpus  $\mathcal{C} = \bigcup_{i=1}^{i=n} C_i$  where  $C_i$  denotes a set of documents (e.g., sentences, paragraphs) of the time  $t_i$ . CSSDetection consists in assessing the change of meaning for a set of target words  $\mathcal{W}$  occurring in  $\mathcal{C}$  across the whole time span  $[1 \dots n]$  by leveraging contextualised embeddings.

Approaches to CSSDetection rely on semantic modeling of words and their foundations lie in the well-known distributional hypothesis: “You shall know a word by the company it keeps” [Firth, 1957], meaning that the semantic representation of a word is determined by analysing the patterns of lexical co-occurrence within a considered document corpus.

As a difference with static models, where words are encoded in a single vector representation, contextualised models generate different word representations according to the context in which they occur. For instance, consider the word *sex*. Different semantic vectors are generated when the word in the input sequence is used with the *fair sex* connotation or with the *sexual activity* meaning.

For the sake of readability, in the following, we consider the problem of CSSDetection on a set of documents  $\mathcal{C} = C_1 \cup C_2$  and we consider to evaluate the shift of a given target word  $w \in \mathcal{W}$  on a single time period from  $t_1$  to  $t_2$ . This simplification enables to review approaches to CSSDetection in a clear and concise fashion, while being easily extendable to the general case. As a matter of fact, when  $\mathcal{C}$  spans more than two time periods, the change is typically measured by re-implementing the approaches across all contiguous pairs of periods [Giulianelli et al., 2020].

The CSSDetection workflow is described in Figure 1.

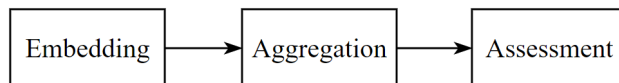


Figure 1: A general workflow of CSSDetection

The initial *embedding stage* has the goal to represent the word occurrences in a multi-dimensional semantic space where the word representations are similar for word occurrences in similar sentences. An optional *aggregation stage* can be enforced to group multiple word representations into a single one for detecting similar usage and/or reducing the computational complexity of the overall CSSDetection task. For example, word occurrences can be aggregated according to a sense-oriented criterion. As another example, multiple word representations can be synthesised into a single prototype representation. The final *assessment stage* consists in the application of a semantic measure to evaluate how the meanings of the word shifted over time.

**Embedding.** Consider the subsets of documents  $C_w^1 \subseteq C_1$  and  $C_w^2 \subseteq C_2$  that contain the word  $w$ . In the embedding stage, a contextualised model  $m$  (e.g., BERT [Devlin et al., 2019], RoBERTa [Liu et al., 2019], ELMo [Peters et al., 2018]) is employed to extract an embedding vector for each occurrence of  $w$  in  $C_w^1$  and  $C_w^2$ . The contextualised embedded representation of the word  $w$  in the  $i$ -th document of a corpus  $C_w^j$  is denoted by  $e_{w,i}^j$  ( $j \in \{1, 2\}$ ). Then, the representation of the word  $w$  in a corpus  $C_w^j$  is defined as:  $\Phi_w^j = \{e_{w,1}^j, \dots, e_{w,z}^j\}$ , with  $z$  being the cardinality of  $C_w^j$ , namely the number of documents in  $C_w^j$ . As a result, we denote as  $\Phi_w^1$  and  $\Phi_w^2$  the sets of embedding vectors generated for the word  $w$  at time  $t_1$  and  $t_2$ , respectively.

**Aggregation.** This stage is optionally executed and it has two main goals: i) to recognise when different word occurrences purport a similar meaning, and ii) to reduce the number of elements to consider for shift detection. To this end, clustering and averaging techniques are proposed for aggregating the word embeddings previously created.

*Clustering.* Clustering techniques are employed to group similar word embeddings in a cluster, each one loosely denoting a specific word meaning. In some approaches, it is assumed that the corpus is *static*, meaning that all the documents in  $C_w^1$  and  $C_w^2$  are available as a whole. Then, a *joint* clustering operation is executed over the embeddings of  $\Phi_w^1 \cup \Phi_w^2$  (e.g. Martinc et al. [2020a]). In other approaches, it is assumed that the corpus is *dynamic*, meaning that documents become available at different time moments and a *separate* clustering operation is performed over the embeddings of  $\Phi_w^1$  and  $\Phi_w^2$ , individually (i.e., one exclusively on  $\Phi_w^1$  and another exclusively on  $\Phi_w^2$  embeddings). When a separate clustering is executed, the resulting clusters need to be aligned in order to recognise similar word meanings at different consecutive times (e.g. Kanjirangat et al. [2020]). To overcome the need for aligning clusters, an *incremental* clustering operation is employed to progressively group the embedding available at the different time steps (e.g. Periti et al. [2022]). The result of clustering is a set of  $k$  clusters where the  $i$ -th cluster is denoted as  $\phi_{w,i}$  and it can fall into one of the following cases (see Figure 2):

- A.  $\phi_{w,i}$  contains only embeddings from  $C_w^1$ ;
- B.  $\phi_{w,i}$  contains a mixture of embeddings from both  $C_w^1$  and  $C_w^2$ ;
- C.  $\phi_{w,i}$  contains only embeddings from  $C_w^2$ .

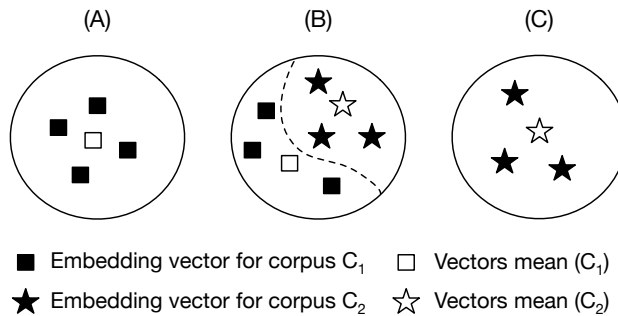


Figure 2: Possible cluster composition (from Periti et al. [2022])

As a result, a cluster  $\phi_{w,i} = \phi_{w,i}^1 \cup \phi_{w,i}^2$  is composed by the union of two partitions  $\phi_{w,i}^1$  and  $\phi_{w,i}^2$  denoting the embeddings from  $\Phi_w^1$  and  $\Phi_w^2$ , respectively. When a *joint* or *incremental* clustering is applied, the resulting clusters can belong to any of the above cases (i.e., A, B, and C). When a *separate* clustering is applied, the resulting clusters can just belong to A and C cases, meaning that  $\phi_{w,i}^2 = \emptyset$  and  $\phi_{w,i}^1 = \emptyset$ , respectively.

*Averaging.* Averaging techniques consist in determining a prototypical representation of the word  $w$ . As an option, a *word-prototype*, can be computed by averaging all its embedding. In this case, *word-prototypes*  $\mu_w^1$  and  $\mu_w^2$  are created as the average embeddings of  $\Phi_w^1$  and  $\Phi_w^2$ , respectively (e.g., Rodina et al. [2020]). As an alternative option, averaging can be executed on top of the results of clustering. For each cluster, averaging is used to create a prototypical representation of all the cluster elements (i.e., the centroid of the cluster). In particular, *sense-prototypes*  $c_{w,i}^1, c_{w,i}^2$  can be created for each cluster  $\phi_{w,i}$  as the average embedding of its cluster partitions  $\phi_{w,i}^1, \phi_{w,i}^2$ , respectively (e.g., Periti et al. [2022]).

**Assessment.** This stage has the goal to measure the shift on the meanings of the word  $w$  across the corpora  $C_1$  and  $C_2$  by considering the sets  $\Phi_w^1$  and  $\Phi_w^2$ . In the literature, a number of functions are proposed for semantic shift assessment. We can distinguish measures that assess the shift by considering the whole set of embedding representations  $\Phi_w^i$ , by those that exploit the prototypical representations  $c_{w,i}^i$  and/or  $\mu_w^i$  generated during the aggregation step through clustering and/or averaging. According to Kutuzov et al. [2018], the definition of a rigorous, formal, mathematical model for representing the assessment functions used in CSSDetection approaches is a challenging issue. In the following, we provide a formal definition of an abstract function  $f$  that has the goal to encompass all the existing assessment measures. The semantic shift assessment  $s_w = f(\cdot, \cdot)$  is defined as follows:

$$f : \{\mathbb{R}^D\}^{(p_1+z_1 \cdot \delta)}, \{\mathbb{R}^D\}^{(p_2+z_2 \cdot \delta)}, c \rightarrow \mathcal{S}$$

where  $D$  is the dimension of the word vectors in  $\Phi_w^1$  and  $\Phi_w^2$ ;  $p_1, p_2$  are the number of prototypical embeddings under consideration for  $C_w^1, C_w^2$ , respectively;  $z_1, z_2$  are the number of vectors in  $\Phi_w^1$  and  $\Phi_w^2$ , respectively;  $\delta \in \{0, 1\}$  is a flag that allows to distinguish the approaches according to the kind of embedding used (i.e., original and/or prototypical);  $c$  is a counting function that determines the normalised number of embeddings in the cluster partitions  $\phi_{w,i}^1$  and  $\phi_{w,i}^2$ , respectively.

The counting function  $c$  is defined as:

$$c : \{\mathbb{R}^D\}^{z_1}, \{\mathbb{R}^D\}^{z_2} \rightarrow \mathbb{R}^k, \mathbb{R}^k$$

where  $k$  denotes the number of  $k$  clusters obtained when a clustering operation is enforced during the aggregation stage. When the clustering operation is not enforced, each embedding is mapped to a singleton group (i.e.,  $k = z_1 + z_2$ ).

The signature of  $f$  depends on the possible execution of an aggregation technique:

- *Clustering.* When the clustering operation is executed, then  $p_1 = p_2 = 0$  and  $\delta = 1$ . This means that all the  $z_1 + z_2$  embeddings in  $\Phi_w^1 \cup \Phi_w^2$  are exploited for semantic shift assessment (e.g., Martinc et al. [2020a]).
- *Averaging.* When the averaging operation is executed, then  $p_1 = p_2 = 1$ . In some approaches,  $\delta = 0$  and this means that the function  $f$  is defined as a distance measure over prototypical representations (e.g., Martinc et al. [2020b]). In some other approaches,  $\delta = 1$  and this means that  $f$  is defined as a distance measure over the original embeddings  $\Phi_w$  and their prototypical representations (e.g., Pömsl and Lyapin [2020]).
- *Clustering + Averaging.* When both clustering and averaging are performed,  $p_1, p_2 > 0$  and  $\delta$  can be both 0 or 1 as in the previous case (e.g., Periti et al. [2022]).

The output  $\mathcal{S}$  is defined according to four different assessment questions.

- *Grade Change Detection.* The goal of Grade Change Detection is to quantify the assessment  $s_w$ . Then,  $\mathcal{S} = \mathbb{R}$  represents the extent to which  $w$  shifts between  $C_1$  and  $C_2$  [Schlechtweg et al., 2020].
- *Binary Change Detection.* The goal of Binary Change Detection is to classify  $w$  as stable (without lost or gained sense(s)) or changed (with lost or gained sense(s)). In this case,  $s_w$  is binary, meaning that  $\mathcal{S} = \{0, 1\}$  for stable and changed, respectively [Schlechtweg et al., 2020].
- *Sense Gain Detection.* The goal of Sense Gain Detection is to recognise whether  $w$  gained meanings or not. In this case,  $s_w$  is binary, meaning that  $\mathcal{S} = \{0, 1\}$  for not-gained and gained, respectively [Zamora-Reina et al., 2022].
- *Sense Loss Detection.* The goal of Sense Gain Detection is to recognise whether  $w$  lost meanings or not. In this case,  $s_w$  is binary, meaning that  $\mathcal{S} = \{0, 1\}$  for not-lost and lost, respectively [Zamora-Reina et al., 2022].

In this survey, we focus on approaches that adopt Grade Change Detection since it is the most commonly enforced assessment. As a matter of fact, we note that the approaches based on Grade Change Detection can be transformed into Binary Change Detection by binarising  $s_w$  through a threshold  $\theta$ . We do not address Sense Gain and Sense Loss Detection as they are relatively novel assessment questions.

For the sake of clarity, we summarise the notation used throughout this paper in Table 1.

Notation	Definition
$w$	Target word
$C_j$	Set of documents at time $t_j$
$C_w^j$	Subset of documents of $C_j$ containing the word $w$
$e_{w,i}^j$	Contextualised embedding of the word $w$ in the $i$ -th document of a corpus $C_w^j$
$\Phi_w^j$	Set of the embeddings of $w$ in the corpus $C_w^j$
$\phi_{w,i}$	$i$ -th cluster containing the embeddings of the word $w$
$\phi_{w,i}^j$	Subset of contextualised embeddings $\Phi_w^j$ in the cluster $\phi_{w,i}$
$\mu_w^j$	Prototypical representation of $w$ for $\Phi_w^j$
$\mu_{w,i}^j$	Prototypical representation of $w$ for $\phi_{w,i}^j$

Table 1: A reference table of notations used in the paper

### 3 A classification framework for CSSDetection

A consolidated and widely-accepted classification framework of CSSDetection approaches is not available. A basic framework is focused on the meaning representation of the words by distinguishing between form- and sense-based approaches [Giulianelli et al., 2020, Wenjun Qiu and Xu, 2022]. However, such a distinction is not universally recognised with a unique interpretation. Sometimes, these two categories are referred as *type*- and *token*-based, where averaging and clustering are enforced to aggregate embeddings, respectively [Laicher et al., 2020, Schlechtweg et al., 2020]. More recently, *average*- and *cluster*-based categories have been proposed to rename form and sense ones to highlight the method used for embedding aggregation [Periti et al., 2022].

In the following, we propose a comprehensive classification framework that extends the basic distinction between form- and sense-based approaches by introducing three dimensions of analysis, namely *meaning representation*, *time-awareness*, and *learning modality* (see Table 2).

Meaning representation	Time-awareness	Learning modality
form-based	time-oblivious	supervised
sense-based	time-aware	unsupervised

Table 2: A classification framework for CSSDetection

**Meaning representation.** Borrowing the distinction proposed by Giulianelli et al. [2020], this dimension focuses on the meaning representation of a word. Two categories are defined:

- *form-based*: the meaning representation concerns the high-level properties of the target word  $w$ , such as its degree of polysemy or its dominant sense. When the polysemy is considered, the CSSDetection approaches do not enforce any aggregation stage and the semantic shift of  $w$  is assessed by measuring the degree of change on the embeddings  $\Phi_w^1$  and  $\Phi_w^2$  (i.e., change on the degree of polysemy). When the dominant sense is considered, all the meanings of  $w$  are collapsed into a single one on which the shift is assessed. Typically, the embeddings  $\Phi_w^1$  and  $\Phi_w^2$  are averaged into corresponding word prototypes  $\mu_w^1$  and  $\mu_w^2$ , respectively. In this case, the CSSDetection approaches focus on one meaning of  $w$  that can be considered as an approximation of the *dominant sense* since, generally, it is the most frequent in the corpus, and thus the one most represented in the word prototype. We stress that form-based approaches are not able to represent how minor meanings *compete* and *cooperate* to change the dominant sense [Hu et al., 2019].
- *sense-based*: the meaning representation concerns the low-level properties of the target word  $w$ , such as its different word usages (i.e., its multiple meanings). All the senses of a word  $w$  are represented and considered in the shift assessment, namely the dominant sense and the minor ones. Typically, the embeddings  $\Phi_w^1$  and  $\Phi_w^2$  are aggregated into clusters, each one representing a different usage/meaning of  $w$ . Sense-based approaches allow to capture the changes over the different meanings of  $w$  as well as to interpret the word change (e.g., a new/existing meaning has gained/lost importance).

**Time awareness.** This dimension focuses on how the time information of the documents is considered in the embedding model. Two categories are defined:

- *time-oblivious*: this category is based on the assumption that a document of time  $t$  adopts linguistic patterns that already characterise the language at the time  $t$  by its own. Thus, it is not needed that the embedding model is aware of the time in which a document is inserted in the corpus. A time-oblivious model is based on *the*

*contextual nature of embeddings generated by the model, which by definition are dependent on the context that is always time-specific* [Martinc et al., 2020a].

- *time-aware*: this category is based on the assumption that contextualised embedding models are not capable of *adapting to time and generalising temporally* since they are *usually pre-trained on corpora derived from a snapshot of the web crawled at a specific moment in time* [Rosin et al., 2022]. Thus, it is needed that the embedding model is aware of the time in which a document is inserted in the corpus. As a result, a time-aware model encodes the time information as well as the linguistic context of a document while generating the embeddings.

**Learning modality.** This dimension is about the possible use of external knowledge for describing and learning the word meanings to recognise. Two categories are defined:

- *supervised*: a form of supervision is enforced to inject external knowledge to support the shift assessment. In addition to the text in the corpora  $C_1$  and  $C_2$ , a lexicographic/manual supervision is employed. By lexicographic supervision, we mean that a dictionary/thesaurus is introduced to recognise the meanings of the word  $w$ . This solution can be considered as an alternative to aggregation by clustering for meaning identification. By manual supervision, we mean that a human-annotated dataset with gold labels is provided for training the embedding model.
- *unsupervised*: the shift assessment is exclusively based on the text of the corpora  $C_1, C_2$  without any external knowledge support. As a result, the word meanings to recognise emerge from the corpora and the shift is completely assessed by exploiting unsupervised learning techniques. The use of aggregation by clustering is an example of unsupervised learning for meaning detection.

## 4 Approaches to CSSDetection

In this section, we review the literature about CSSDetection according to the classification framework discussed in Section 3. In particular, the solutions are presented in Sections 4.1 and 4.2 according to the meaning representation of the considered target word, namely *form-* and *sense-* based approaches, respectively. Moreover, in Section 4.3, we describe the so-called *ensemble* approaches, namely approaches that are based on a combination of form-/sense-based solutions.

For the sake of comparison, in each category (i.e., form, sense, ensemble), a summary table is provided to frame the literature papers according to our dimensions of analysis as well as to report additional descriptive features about the following aspects:

- *Language model*: the contextualise language model used (e.g., ELMo, BERT, RoBERTa);
- *Training language*: the language of the dataset used for training the model. The possible options are *monolingual* to denote when training is executed on a single language, or *multilingual* when more than one language is considered.
- *Type of training*: how the model is trained. We distinguish five categories:
  - *trained*: the model is trained from scratch through the typical objective function of the architecture model;
  - *pre-trained*: the model is pre-trained through the typical objective function without further training;
  - fine-tuned for *domain-adaptation*: the model is pre-trained through the typical objective function, then it is fine-tuned on new data through the same objective function;
  - fine-tuned for *incremental domain-adaptation*: the model is fine-tuned on the corpus of the first time period  $C_1$ . Then, it is re-tuned separately on the corpus  $C_2$ . The model at time  $t_2$  is initialised with the weights from the model at time  $t_1$ , so that both models are inherently related the one to the other;
  - *fine-tuned*: the model is pre-trained through the typical objective function, then it is fine-tuned on new data through a different objective function.
- *Layer*: the architecture’s layer(s) from which word representations are extracted;
- *Layer aggregation*: the type of aggregation used to synthesise the word representations extracted from different layers into a single embedding;
- *Clustering algorithm*: the clustering algorithm used in the aggregation stage;
- *Shift function*: the function  $f$  used to detect/assess the semantic shift;
- *Corpus language*: the natural language of the corpus in the considered experiments of shift assessment (e.g., English, Italian, Spanish).

#### 4.1 Form-based approaches

Ref.	Time awareness	Learning modality	Language model	Training language	Type of training	Layer	Layer aggregation	Clustering algorithm	Shift function	Corpus language
Arefyev et al.	time-oblivious	supervised	XLM-R-large	multilingual	fine-tuned	last	-	-	APD	Russian
Beck	time-oblivious	unsupervised	mBERT-base	multilingual	pre-trained	last two	average	K-Means	CD	English, German, Latin, Swedish
Martinc et al.	time-oblivious	unsupervised	BERT-base, mBERT-base	monolingual, multilingual	domain-adaptation	last four	sum	-	CD	English, Slovenian
Horn	time-oblivious	unsupervised	BERT-base, RoBERTa-base	monolingual	domain-adaptation, pre-trained	-	-	-	CD	English
Hofmann et al.	time-aware	unsupervised	BERT-base	monolingual	fine-tuned	last	-	-	CD	English
Zhou and Li	time-aware	unsupervised	BERT-base	monolingual	domain-adaptation	last four	sum	-	CD	English, German, Latin, Swedish
Rosin et al.	time-aware	unsupervised	BERT-base, BERT-tiny	monolingual	fine-tuned	all, last, last four	average	-	CD, TD	English, Latin
Rosin and Radinsky	time-aware	unsupervised	BERT-base, BERT-small, BERT-tiny	monolingual	fine-tuned	all, last, last four, last two	average	-	CD	English, German, Latin
Kutuzov and Giulianelli	time-oblivious	unsupervised	BERT-base, ELMo, mBERT-base	monolingual, multilingual	domain-adaptation, incremental domain-adaptation, pre-trained, trained	all, last, last four	average	-	APD, CD, PRT	English, German, Latin, Swedish
Giulianelli et al.	time-oblivious	unsupervised	BERT-base	monolingual	pre-trained	all	sum	-	APD	English
Keidar et al.	time-oblivious	unsupervised	RoBERTa-base	monolingual	domain-adaptation	all, first, last	sum	-	APD	English
Pömsl and Lyapin	time-aware	unsupervised	BERT-base, mBERT-base	monolingual, multilingual	fine-tuned	last	-	-	APD	English, German, Latin, Swedish
Kudisov and Arefyev	time-oblivious	unsupervised	XLM-R-large	multilingual	pre-trained	-	-	-	APD	Spanish
Laicher et al.	time-oblivious	unsupervised	BERT-base	monolingual	pre-trained	first, first + last, first four, last, last four	average	-	APD-OLD/NEW, CD	English, German, Swedish
Wang et al.	time-oblivious	unsupervised	mBERT-base	multilingual	pre-trained	last	-	-	APD, HD	Italian
Kutuzov	time-oblivious	unsupervised	BERT-base, ELMo, mBERT-base	monolingual, multilingual	domain-adaptation, pre-trained	all, last, last four	average	-	APD, DIV, PRT	English, German, Latin, Swedish, Russian
Ryzhova et al.	time-oblivious	unsupervised	ELMo, RuBERT 2019	multilingual	pre-trained, trained	-	-	-	APD	Russian
Rodina et al.	time-oblivious	unsupervised	ELMo, RuBERT	monolingual, multilingual	domain-adaptation	last	-	-	PRT	Russian
Liu et al.	time-oblivious	unsupervised	BERT-base, LatinBERT 2020	multilingual, monolingual	domain-adaptation	last four	sum	-	CD	English, German, Latin, Swedish
Giulianelli et al.	time-oblivious	unsupervised	XLM-R-base	multilingual	domain-adaptation	all	average	-	APD, PRT	English, German, Italian, Latin, Norwegian, Russian, Swedish
Laicher et al.	time-oblivious	unsupervised	mBERT-base	multilingual	pre-trained	all, last four	average	-	APD	Italian
Wenjun Qiu and Xu	time-oblivious	unsupervised	BERT-base	monolingual	domain-adaptation pre-trained	last four	sum	-	CD	English
Periti et al.	time-oblivious	unsupervised	BERT-base, mBERT-base	monolingual, multilingual	pre-trained	last four	sum	-	CD, DIV	English, Latin
Montariol et al.	time-oblivious	unsupervised	BERT-base, mBERT-base	monolingual, multilingual	domain-adaptation	last four	sum	-	CD	English, German, Latin, Swedish

Table 3: Summary view of form-based approaches. Missing information is denoted with a dash

According to Table 3, we note that most form-based approaches are time-oblivious. A few time-aware approaches have been recently appeared and they are all characterised by the adoption of a specific fine-tuning operation to inject time information into the model. All the papers leverage unsupervised learning modalities with the exception of Arefyev et al. [2021]. The aggregation stage is mostly based on averaging, while clustering is only enforced in Beck [2020] where a cluster represents the dominant sense of the word  $w$ . In particular, in Beck [2020], a word is considered as changing when clustering the embeddings  $\Phi_w^1$  and  $\Phi_w^2$  via K-means with  $k = 2$  generates two groups where one of the two clusters contains at least 90% of the embeddings from one corpus only ( $C_w^1$  or  $C_w^2$ ).

In form-based approaches, the following shift functions are proposed for measuring the semantic shift  $s_w$ .

**Cosine distance (CD).** The shift  $s_w$  is measured as the *cosine distance* (CD) between the word prototypes  $\mu_w^1, \mu_w^2$  as follows:

$$CD(\mu_w^1, \mu_w^2) = 1 - CS(\mu_w^1, \mu_w^2) \quad (1)$$

where  $CS$  is the *cosine similarity* between the prototypes. Intuitively, the greater the  $CD(\mu_w^1, \mu_w^2)$ , the greater the shift in the dominant sense of  $w$ .

Typically, the prototypes  $\mu_w^1$  and  $\mu_w^2$  are determined through aggregation by averaging over  $\Phi_w^1$  and  $\Phi_w^2$ , respectively (e.g., Martinc et al. [2020b]). As a difference, in Horn [2021], the prototype embedding  $\mu_w^2$  at time step  $t = 2$  is computed by updating the prototype embedding  $\mu_w^1$  at time step  $t = 1$  through a weighted running average (e.g., Finch [2009]).

In Martinc et al. [2020b], the CD metric is employed in a multilingual experiment where the shift is measured across a diachronic corpus with texts of different languages. This is the only example of cross-language shift detection.

CD is also used in time-aware approaches. The integration of extra-linguistic information into word embeddings, such as time and social space, has been proposed in previous work based on static models [Rudolph and Blei, 2018, Zeng et al., 2018]. Recently, this integration has been also applied to contextualised embeddings [Huang and Paul, 2019, Röttger and Pierrehumbert, 2021]. In Hofmann et al. [2021], a pre-trained model is fine-tuned to encapsulate time and social space in the embedding model [Hofmann et al., 2021]. Then, the shift  $s_w$  is assessed by computing the CD between embeddings generated by the original pre-trained model and the embeddings generated by the time-aware, fine-tuned model. In particular, in Zhou and Li [2020], a *temporal referencing* mechanism is adopted to encode time-awareness into a pre-trained model. Temporal referencing is a pre-processing step of the documents that tags each occurrence of  $w$  in  $C_w^1$  and  $C_w^2$  with a special marker denoting the corpus/time in which it appears [Ferrari et al., 2017, Dubossarsky et al., 2019]. The embeddings of a tagged word are learned by fine-tuning the model for domain-adaptation. In this case,  $s_w$  is assessed by computing the CD between  $\mu_{w[1]}^1$  and  $\mu_{w[2]}^1$ , where  $w[i]$  denotes  $w$  with the temporal marker  $t_i$ . Similarly to Zhou and Li [2020], a time-aware approach is proposed in Rosin et al. [2022] where a time marker is added to documents instead of words and the model is fine-tuned to predict the injected time information. As an alternative, in Rosin and Radinsky [2022], a *temporal attention* mechanism is adopted to generate the embeddings  $\Phi_w^1$  and  $\Phi_w^2$  for calculating CD.

**Inverted similarity over word prototype (PRT).** This measure is proposed as an alternative to CD for improving the effectiveness of the shift detection [Kutuzov and Giulianelli, 2020]. The *inverted similarity over word prototypes* (PRT) measure is defined as:

$$PRT(\mu_w^1, \mu_w^2) = \frac{1}{CS(\mu_w^1, \mu_w^2)}. \quad (2)$$

**Time-diff (TD).** This measure is designed for time-aware approaches and it works on analysing the change of polysemy of a word along time. It is based on the model capability to predict the time of a document and it calculates the shift  $s_w$  by considering the probability distribution of the predicted times [Rosin et al., 2022]. Intuitively, a uniform distribution means that the association document-time is not strong enough to clearly entail a shift. Instead, a non-uniform distribution means that there is an evidence to predict the time of a document. Consider a document  $d_w$ , let  $p_t(d_w)$  be the probability of  $d_w$  to belong to the time  $t$ . The function *time diff* (TD) is defined as the average difference of the predicted time probabilities:

$$TD(C_w^1, C_w^2) = \frac{1}{|C_w|} \sum_{d_w^1 \in C_w^1, d_w^2 \in C_w^2} |p_1(d_w^1) - p_2(d_w^2)|. \quad (3)$$

The experiments performed in Rosin et al. [2022] show that TD outperforms CD on short-term semantic shift. On the contrary, CD outperforms TD over long-term semantic shift. We argue that the time-diff measure is more effective on long-term periods since major differences in writing style emerge and the prediction of document-time associations is more reliable.

**Average pairwise distance (APD).** This measure exploits the variance of the contextualised representations  $\Phi_w^1, \Phi_w^2$  to compute the semantic shift assessment (i.e., variance on the word polysemy). As a difference with the previous measures, APD directly works on word-occurrence embeddings without requiring any aggregation stage, namely clustering nor averaging. The *average pairwise distance* (APD) is defined as follows:

$$APD(\Phi_w^1, \Phi_w^2) = \frac{1}{|\Phi_w^1| |\Phi_w^2|} \cdot \sum_{e_{w,i}^1 \in \Phi_w^1, e_{w,i}^2 \in \Phi_w^2} d(e_{w,i}^1, e_{w,i}^2), \quad (4)$$

where  $d$  is an arbitrary distance measure (e.g., cosine distance, Euclidean distance, Canberra distance). According to the experiments performed in Giulianelli et al. [2020], APD better performs when the Euclidean distance is employed as  $d$ . In Keidar et al. [2022], APD is used over the embeddings  $\Phi_w^1$  and  $\Phi_w^2$  by applying a dimensionality reduction through the Principal Component Analysis (PCA). In Keidar et al. [2022], experiments on both slang and non-slang words are



performed through causal analysis to study how distributional factors (e.g., polysemy, frequency shift) influence the shift  $s_w$ . The results show that slang words experience fewer semantic shifts than non-slang words.

In Kudisov and Arefyev [2022], lexical substitutes are used to assess  $s_w$ . Given a word  $w$ , lexical substitutes of  $w$  are those words that can replace  $w$  in a text fragment without introducing grammatical errors or significantly changing its meaning. A set of lexical substitutes is generated by leveraging a masked language model (e.g., XLM-R). In this case, the word representations  $\Phi_w^1$ , and  $\Phi_w^2$  are defined as the bag-of-words vectors computed over the substitutes (i.e., *bag-of-substitutes*) through Tf-Idf. Then, APD is finally computed over  $\Phi_w^1$ , and  $\Phi_w^2$  to assess  $s_w$ .

APD is also used in a time-aware approach described in Pömsl and Lyapin [2020], where a pre-trained BERT model is fine-tuned to predict the time period of a sentence. APD is finally used to measure the shift between the embeddings extracted from the fine-tuned model.

In Arefyev et al. [2021], ADP is employed to measure the shift  $s_w$  over the embeddings  $\Phi_w^1$  and  $\Phi_w^2$  extracted from a supervised Word-in-Context model (WiC) [Pilehvar and Camacho-Collados, 2019]. This model is trained to reproduce the behavior of human annotators when they are asked to evaluate the similarity of the meaning of a word  $w$  in a pair of given sentences from  $C_w^1$  and  $C_w^2$ , respectively. The embeddings  $\Phi_w^1$  and  $\Phi_w^2$  are extracted from the trained WiC model for calculating the final APD measure.

**Average of average inner distances (APD-OLD/NEW).** The APD-OLD/NEW measure is presented in Laicher et al. [2021] as an extension of APD and it estimates the shift  $s_w$  as the average degree of polysemy of  $w$  in the corpora  $C_w^1$  and  $C_w^2$ , respectively. The *average of average inner distances* (APD-OLD/NEW) is defined as:

$$APD-OLD/NEW(\Phi_w^1, \Phi_w^2) = \frac{AID(\Phi_w^1) + AID(\Phi_w^2)}{2}. \quad (5)$$

where AID is the *average inner distance* and it measures the degree of polysemy of  $w$  in a specific time frame by relying on the APD measure, namely  $AID(\Phi_w^1) = APD(\Phi_w^1, \Phi_w^1)$  and  $AID(\Phi_w^2) = APD(\Phi_w^2, \Phi_w^2)$ , respectively.

**Hausdorff distance (HD).** The shift  $s_w$  is measured as the *Hausdorff distance* (HD) between the word embeddings  $\Phi_w^1$  and  $\Phi_w^2$ . HD relies on the Euclidean distance  $d$  to measure the difference between the embeddings of  $w$  in  $C_w^1$  and  $C_w^2$  and it returns the greatest of all the distances  $d$  from one embedding  $e_w^1 \in \Phi_w^1$  to the closest embedding  $e_w^2 \in \Phi_w^2$ , or vice-versa. The HD measure is defined as follows:

$$HD(\Phi_w^1, \Phi_w^2) = \max \left( \sup_{e_w^1 \in \Phi_w^1} \inf_{e_w^2 \in \Phi_w^2} d(e_w^1, e_w^2), \sup_{e_w^2 \in \Phi_w^2} \inf_{e_w^1 \in \Phi_w^1} d(e_w^2, e_w^1) \right). \quad (6)$$

The experiments performed in Wang et al. [2020] show that HD is sensitive to outliers since it is based on infimum and supremum, thus an outlier embedding may largely affect the final  $s_w$  value.

**Difference between token embedding diversities (DIV).** Similar to APD, this measure assesses the shift  $s_w$  by exploiting the variance of the contextualised representation  $\Phi_w^1$  and  $\Phi_w^2$ . As a difference with APD, the *difference between token embedding diversities* (DIV) leverage a coefficient of variation calculated as the average of the cosine distances  $d$  between the embeddings  $\Phi_w^1$  and  $\Phi_w^2$ , and their prototypical embeddings  $\mu_w^1$  and  $\mu_w^2$ , respectively Kutuzov [2020]. The intuition is that when  $w$  is used in just one sense, its embeddings tend to be close to each other yielding a low coefficient of variation. On the opposite, when  $w$  is used many different senses, its embeddings are distant to each other yielding to a high coefficient of variation. DIV is defined as the absolute difference between the coefficient of variation in  $C_w^1$  and  $C_w^2$ :

$$DIV(\Phi_w^1, \Phi_w^2) = \left| \frac{\sum_{e_w^1 \in \Phi_w^1} d(e_w^1, \mu_w^1)}{|\Phi_w^1|} - \frac{\sum_{e_w^2 \in \Phi_w^2} d(e_w^2, \mu_w^2)}{|\Phi_w^2|} \right| \quad (7)$$

In Kutuzov [2020], the experiments show that when the coefficient of variation is low, the prototypical embeddings  $\mu_w^1$  and  $\mu_w^2$  successfully represent the meanings of the given word  $w$ . On the opposite, when the coefficient of variation is high, the prototypical embeddings  $\mu_w^1$  and  $\mu_w^2$  do not provide a relevant representation of the  $w$  meanings.

## 4.2 Sense-based approaches

According to Table 4, we note that all the sense-based approaches are time-oblivious and that fine-tuning is sometimes adopted, but mainly for domain-adaptation purposes. Most papers leverage unsupervised learning modalities. Only

Ref.	Time awareness	Learning modality	Language model	Training language	Type of training	Layer	Layer aggregation	Clustering algorithm	Shift function	Corpus language
Hu et al.	time-oblivious	supervised	BERT-base	monolingual	pre-trained	last	-	-	MNS	English
Rachinskiy and Arefyev	time-oblivious	supervised	XLM-R-base	multilingual	fine-tuned, pre-trained	-	-	-	APD	Russian
Rachinskiy and Arefyev	time-oblivious	supervised	XLM-R-base	multilingual	fine-tuned, pre-trained	last	-	-	APD, JSD	Spanish
Periti et al.	time-oblivious	unsupervised	BERT-base, mBERT-base	monolingual, multilingual	pre-trained	last four	sum	AP, APP, IAPNA	JSD, PDIS, PDIV	English, Latin
Montariol et al.	time-oblivious	unsupervised	BERT-base, mBERT-base	monolingual, multilingual	domain-adaptation	last four	sum	K-Means, AP	JSD, WD	English, German, Latin, Swedish
Rodina et al.	time-oblivious	unsupervised	mBERT-base, ELMo	monolingual, multilingual	domain-adaptation	last	-	K-Means, AP	JSD, MS	Russian
Kanjirang et al.	time-oblivious	unsupervised	mBERT-base	multilingual	pre-trained	last four	concatenation	K-Means	CSC, JSD	English, German, Latin, Swedish
Giulianelli et al.	time-oblivious	unsupervised	BERT-base	monolingual	pre-trained	all	sum	K-Means	ED, JSD	English
Arefyev and Zhikov	time-oblivious	unsupervised	XLM-R-base	multilingual	domain-adaptation	-	-	AGG	CDCD	English, German, Latin, Swedish
Kashleva et al.	time-oblivious	unsupervised	BERT-base	monolingual	domain-adaptation	all	sum	K-Means	APDP	Spanish
Martinc et al.	time-oblivious	unsupervised	BERT-base, mBERT-base	monolingual, multilingual	domain-adaptation	last four	sum	K-Means, AP	JSD	English, German, Latin, Swedish
Kutuzov and Giulianelli	time-oblivious	unsupervised	BERT-base, ELMo, mBERT-base	monolingual, multilingual	domain-adaptation, incremental domain-adaptation, pre-trained	all, last, last four	average	AP	JSD	English, German, Latin, Swedish
Giulianelli et al.	time-oblivious	unsupervised	XLM-R-base	multilingual	domain-adaptation	all	average	AP	JSD	English, German, Italian, Latin, Norwegian, Russian, Swedish
Wang et al.	time-oblivious	unsupervised	mBERT-base	multilingual	domain-adaptation	last	-	GMMs, K-Means	JSD	Italian
Keidar et al.	time-oblivious	unsupervised	RoBERTa-base	monolingual	domain-adaptation	all, first, last	sum	AP, K-Means, GMMs	ED, JSD	English
Karnysheva and Schwarz	time-oblivious	unsupervised	ELMo, mELMo	monolingual, multilingual	pre-trained	all	-	K-Means, DBSCAN	JSD	English, German, Latin, Swedish
Cuba Gyllensten et al.	time-oblivious	unsupervised	XLM-R-base	multilingual	pre-trained	last	-	K-Means	JSD	English, German, Latin, Swedish
Rother et al.	time-oblivious	unsupervised	mBERT-base, XLM-R-base	multilingual	fine-tuned	last	-	BIRCH, DBSCAN, GMMs, HDBSCAN	JSD	English, German, Latin, Swedish

Table 4: Summary view of sense-based approaches. Missing information is denoted with a dash

a few exceptions employ a lexicographic supervision (i.e., Hu et al. [2019], Rachinskiy and Arefyev [2021, 2022]). As a difference with form-based, sense-based approaches usually enforce clustering in the aggregation stage. The aggregation by averaging is only exploited in Periti et al. [2022], Hu et al. [2019], Montariol et al. [2021] where sense prototypes are computed on top of the results of a clustering operation.

When clustering is adopted, the function  $f$  that calculates the shift  $s_w$  can be directly defined over the embeddings  $\Phi_w^1$  and  $\Phi_w^2$ . As an alternative, the function  $f$  can be defined over the distribution of the embeddings in the resulting clusters (i.e., *cluster distribution*). In this case, as a result of the clustering operation, a counting function  $c$  is used to determine two cluster distributions  $p_w^1$  and  $p_w^2$  that represent the normalised number of embeddings in the cluster partitions  $\phi_{w,i}^1$  and  $\phi_{w,i}^2$ , respectively (see Section 2). The  $i$ -th value  $p_{w,i}^j$  in  $p_w^j$  (with  $j \in \{1, 2\}$ ) represents the number of embeddings of  $\phi_{w,i}^j$  in the  $i$ -th cluster, namely:

$$p_{w,i}^j = \frac{|\phi_{w,i}^j|}{|\Phi_w^j|}. \quad (8)$$

Finally, the function  $f$  is defined as a compound function  $f = g \circ c$ , where the result of the  $c$  function is exploited by a shift function  $g$  which works on the cluster distributions  $p_w^1$  and  $p_w^2$ .

In sense-based approaches, the following shift functions are proposed for measuring the semantic shift  $s_w$ .

**Maximum novelty score (MNS).** This measure exploits the cluster distributions  $p_w^1$  and  $p_w^2$  by leveraging the idea that the higher is the ratio between the number of embeddings  $\Phi_w^1$  and  $\Phi_w^2$  in a cluster, the higher is the semantic shift of the

considered word  $w$ . The *maximum novelty score* (MNS) is defined as:

$$MNS(p_w^1, p_w^2) = \max\{NS(p_{w,1}^1, p_{w,1}^2), \dots, NS(p_{w,k}^1, p_{w,k}^2)\}, \quad (9)$$

where  $NS(p_{w,i}^1, p_{w,i}^2) = p_{w,i}^1/p_{w,i}^2$  is the *novelty score* proposed in Cook et al. [2014], and  $k$  is the number of clusters produced as a result of the aggregation stage.

In Hu et al. [2019], MNS is employed as a shift measure in a supervised learning approach. In particular, a lexicographic supervision (i.e., the Oxford English dictionary) is employed to provide the meanings of the target word  $w$ . Each word occurrence in  $\Phi_w^1$  and  $\Phi_w^2$  is associated with the closest meaning of the dictionary according to the cosine distance. As a result, for each word/dictionary meaning, a cluster of word embeddings is defined and MNS is exploited to calculate the overall shift.

**Maximum square (MS).** This measure is an alternative to MNS to assess the shift of  $s_w$ . The intuition of MS is that slight changes in cluster distributions  $p_w^1$  and  $p_w^2$  may occur due to noise and do not represent a real semantic shift [Rodina et al., 2020]. The *maximum square* (MS) aims at identifying strong changes in the cluster distributions. As a difference with MNS, the square difference between  $p_{w,i}^1$  and  $p_{w,i}^2$  is used to capture the degree of shift instead of the novelty score (NS):

$$MS(p_w^1, p_w^2) = \max_i (p_{w,i}^1 - p_{w,i}^2)^2 \quad (10)$$

**Jensen-Shannon divergence (JSD).** This measure extends the Kullback-Leibler (KL) divergence, which calculates how one probability distribution is different from another. The *Jensen-Shannon divergence* (JSD) calculates the shift  $s_w$  as the symmetrical KL score of the cluster distributions  $p_w^1$  from  $p_w^2$ , namely:

$$JSD(p_w^1, p_w^2) = \frac{1}{2} (KL(p_w^1 || M) + KL(p_w^2 || M)) , \quad (11)$$

where KL is the Kullback-Leibler divergence and  $M = (p_w^1 + p_w^2)/2$ .

JSD is also used in approaches where aggregation by clustering is performed separately over the embeddings  $\Phi_w^1$  and  $\Phi_w^2$  [Kanjirangat et al., 2020]. As a result, the clusters need to be aligned to determine the distributions  $p_w^1$  and  $p_w^2$  before the JSD calculation. As a difference with Kanjirangat et al. [2020], an evolutionary clustering algorithm is employed in Periti et al. [2022] to apply the JSD measure without requiring any alignment step over the resulting clusters.

As a final remark, JSD can be employed to measure the shift  $s_w$  over more than two time periods. However, the experiments in Giulianelli et al. [2020] show that the JSD effectiveness over a single time period outperforms the version over more time periods since JSD is insensitive to the order of the temporal intervals.

**Coefficient of semantic change (CSC).** This measure is proposed as an alternative to JSD where the difference over the weighted number of elements in  $\phi_{w,i}^1$  and  $\phi_{w,i}^2$  for each cluster  $i$  is employed to replace KL in measuring the shift [Kanjirangat et al., 2020]. The *coefficient of semantic change* (CSC) is defined as follows:

$$CSC(p_w^1, p_w^2) = \frac{1}{P_w^1 \cdot P_w^2} \sum_{k=1}^K |P_w^2 \cdot p_{w,k}^1 - P_w^1 \cdot p_{w,k}^2| , \quad (12)$$

where  $P_w^j = \sum_{i=1}^k p_{w,i}^j$  is the weight of each cluster distribution and  $k$  is the number of clusters.

**Cosine distance between cluster distributions (CDCD).** As a further alternative of JSD, this measure assess the shift  $s_w$  by considering the cluster distributions  $p_w^1$  and  $p_w^2$  as vectors and by applying the cosine distance over them to assess the semantic shift  $s_w$ . The *cosine distance between cluster distributions* (CDCD) is defined as follows:

$$CDCD(p_w^1, p_w^2) = 1 - \frac{P_w^1 \cdot P_w^2}{\|p_w^1\| \times \|p_w^2\|} \quad (13)$$

In Arefyev and Zhikov [2020], CDCD is calculated between the cluster distributions  $p_w^1$  and  $p_w^2$  obtained by enforcing clustering over bag-of-substitutes (see the description of Arefyev and Zhikov [2020] in Section 4.1).

**Entropy difference (ED).** This measure is based on the idea that the higher is the uncertainty in the interpretation of a word occurrence due to the  $w$  polysemy in  $C_w^1$  and  $C_w^2$ , the higher is the semantic shift  $s_w$ . The intuition is that high

values of ED are associated with the broadening of a word’s interpretation, while negative values indicate a narrowing interpretation [Giulianelli et al., 2020]. The *entropy difference* (ED) is defined as follows:

$$ED(p_w^1, p_w^2) = \eta(p_w^1) - \eta(p_w^2), \quad (14)$$

where  $\eta(p_w^j)$  is the degree of polysemy of  $w$  in the corpus  $C_j$ , which is calculated as the normalised entropy of its cluster distribution  $p_w^j$ :

$$\eta(p_w^j) = \log_{K_w} \left( \prod_{k=1}^{K_w} p_{w,i}^j \cdot^{-p_{w,i}^j} \right).$$

As shown in Giulianelli et al. [2020], ED is not capable of properly assessing  $s_w$  when new usage types of  $w$  emerge, while old ones become obsolescent at the same time, since it may lead to no entropy reduction.

**Cosine distance between word prototypes (PDIS).** This measure is presented in Periti et al. [2022] as an extension of the CD measure adopted by form-oriented approaches. The idea of PDIS is that the aggregation by averaging over cluster prototypes can be employed to produce summary descriptions of the cluster contents (i.e., *semantic prototypes*). The *cosine distance between word prototypes* (PDIS) is defined as the CD between  $\bar{c}_w^1, \bar{c}_w^2$ , that is:

$$PDIS(\bar{c}_w^1, \bar{c}_w^2) = 1 - \frac{\bar{c}_w^1 \cdot \bar{c}_w^2}{\|\bar{c}_w^1\| \times \|\bar{c}_w^2\|} \quad (15)$$

where  $\bar{c}_w^1$  and  $\bar{c}_w^2$  are semantic prototypes defined as the average embeddings of all the sense prototypes  $c_{w,i}^1$  and  $c_{w,i}^2$ , respectively.

**Difference between prototype embedding diversities (PDIV).** This measure is presented in Periti et al. [2022] as an extension of the DIV measure adopted by form-oriented approaches. PDIV leverages the same intuition of PDIS, namely the semantic prototypes can be employed to calculate the coefficient of ambiguity of  $w$  by measuring the difference between a semantic prototype  $\bar{c}_w^j$  and each sense prototype  $c_{w,i}^j$ . The *difference between prototype embedding diversities* (PDIV) is defined as the absolute difference between these ambiguity coefficients:

$$PDIV(\Psi_w^1, \Psi_w^2) = \left| \frac{\sum_{c_{w,k}^1 \in \Psi_w^1} d(c_{w,k}^1, \bar{c}_w^1)}{|\Psi_w^1|} - \frac{\sum_{c_{w,k}^2 \in \Psi_w^2} d(c_{w,k}^2, \bar{c}_w^2)}{|\Psi_w^2|} \right|, \quad (16)$$

where  $\Psi_w^1$  and  $\Psi_w^2$  denote the set of sense prototypes of  $c_{w,i}^1$  and  $c_{w,i}^2$ , respectively.

**Average pairwise distance (APD).** In addition to form-based approaches (see Section 4.1), the APD measure is exploited to assess  $s_w$  also in sense-based approaches. In Rachinskiy and Arefyev [2021, 2022], APD is applied to the contextualised embeddings  $\Phi_w^1$  and  $\Phi_w^2$  extracted from a fine-tuned XLM-R model. In particular, an English corpus is used to fine-tune the pre-trained model to select the most appropriate WordNet’s definition for each word occurrence [Blevins and Zettlemoyer, 2020]. As a result of the fine-tuning, both WordNet’s definitions and word occurrences are embedded in the same vector space and the meaning of any word occurrence can be induced by selecting the closest definition in the vector space. In Rachinskiy and Arefyev [2021], the zero-shot, cross-lingual transferability property of XLM-R is exploited to obtain word representations for Russian language and APD is finally applied [Ming-Wei et al., 2008, Choi et al., 2021]. The authors of Rachinskiy and Arefyev [2021] claim that the approach is useful to overstep the lack of lexicographic supervision for low-resource languages and that most concept definitions in English also hold in other languages, such as Russian. However, this claim is not completely satisfied, since some words can drastically change their meaning across languages. For example, the Russian word "сЛОВО" (pioneer, scout) is strongly connected to the Communist ideology in the Soviet Period, but it isn’t in the English language.

**Average pairwise distance between sense prototypes (APDP).** This measure is an extension of APD and it considers all the pairs of sense prototypes  $c_{w,i}^1$  and  $c_{w,i}^2$  instead of all the original embeddings in  $\Phi_w^1$  and  $\Phi_w^2$  [Kashleva et al., 2022]. The *average pairwise distance between sense prototypes* (APDP) is defined as:

$$APDP(\Psi_w^1, \Psi_w^2) = \frac{1}{|\Psi_w^1| |\Psi_w^2|} \cdot \sum_{c_{w,k}^1 \in \Psi_w^1, c_{w,k}^2 \in \Psi_w^2} d(c_{w,k}^1, c_{w,k}^2) \quad (17)$$

**Wassertein distance (WD).** This measure models the shift assessment as an *optimal transport problem* and it is exploited as an alternative to cluster alignment when aggregation by clustering is performed separately over the

embeddings  $\Phi_w^1$  and  $\Phi_w^2$  [Montariol et al., 2021]. WD quantifies the effort of re-configuring the cluster distribution of  $p_w^1$  into  $p_w^2$ , namely minimising the cost of moving one unit of mass (i.e., a sense prototype) from  $\Psi_w^1$  to  $\Psi_w^2$ . The *Wassertein distance* (WD) is defined as:

$$WD(p_w^1, p_w^2) = \min_{\gamma} \sum_i^{k_1} \sum_j^{k_2} CD(c_{w,i}^1, c_{w,j}^2) \gamma_{c_{w,i}^1 \rightarrow c_{w,j}^2} \quad (18)$$

$$\begin{aligned} \text{such that: } \quad & \gamma_{c_{w,i}^1 \rightarrow c_{w,j}^2} \geq 0 \\ & \sum_i \gamma_{c_{w,i}^1 \rightarrow c_{w,j}^2} = p_w^1 \\ & \sum_j \gamma_{c_{w,i}^1 \rightarrow c_{w,j}^2} = p_w^2 \end{aligned}$$

where all  $\gamma_{c_{w,i}^1 \rightarrow c_{w,j}^2}$  represents the (unknown) effort required to reconfigure the mass distribution  $p_w^1$  into  $p_w^2$ ;  $k_1$  and  $k_2$  are the number of clusters obtained by clustering  $\Phi_w^1$  and  $\Phi_w^2$ , respectively;  $CD$  is the cosine distance computed over the sense prototypes  $c_{w,i}^1 \in \Psi_w^1$  and  $c_{w,j}^2 \in \Psi_w^2$  [Bonneel et al., 2011].

### 4.3 Ensemble-based approaches

In this section, we review the CSSDetection approaches that rely on an *ensemble mechanism*, namely the combination of two or more assessment functions to determine the semantic shift score. Ensembling can mean that more than one form- and/or sense-based measure is adopted in a given approach. Ensembling can also mean that a disciplined use of both static and contextualised embedding models is used. A final semantic shift score is then returned by the whole ensemble process.

Ref.	Time awareness	Learning modality	Language model	Training language	Type of training	Layer	Layer aggregation	Clustering algorithm	Shift function	Corpus language
Pömsl and Lyapin	time-aware	unsupervised	BERT-base, mBERT-base	monolingual, multilingual	fine-tuned	last	-	-	APD	English, German, Latin, Swedish
Teodorescu et al.	time-oblivious	unsupervised	XLM-large	multilingual	trained	last four	sum	-	APD	Spanish
Martinc et al.	time-oblivious	unsupervised	BERT-base, mBERT-base	monolingual, multilingual	domain-adaptation	last four	sum	AP	CD, JSD	English, German, Latin, Swedish
Wang et al.	time-oblivious	unsupervised	mBERT-base	multilingual	pre-trained	last	-	GMMs, K-Means	APD, HD, JSD	Italian
Giulianelli et al.	time-oblivious	unsupervised	XLM-R-base	multilingual	domain-adaptation	all	average	-	APD, PRT	English, German, Italian, Latin, Norwegian, Russian, Swedish
Ryzhova et al.	time-oblivious	unsupervised	ELMo, RuBERT	monolingual, multilingual	pre-trained trained	-	-	-	APD	Russian
Kutuzov et al.	time-oblivious	unsupervised	BERT-base, ELMo	monolingual, multilingual	domain adaptation	last	-	-	APD, PRT	English, German, Latin, Swedish
Rachinskiy and Arefyev	time-oblivious	supervised	XLM-R-base	multilingual	fine-tuned, pre-trained	-	-	-	APD	Russian
Rosin and Radinsky	time-aware	unsupervised	BERT-base	monolingual	fine-tuned	-	-	-	CD	English, Latin, German

Table 5: Summary view of ensemble approaches. Missing information is denoted with a dash

According to Table 5, we note that all the ensemble approaches are time-oblivious with the exception of Pömsl and Lyapin [2020] and Rosin and Radinsky [2022]. We also note that unsupervised learning modalities are adopted with the exception of Rachinskiy and Arefyev [2021]. As a further remark, most of the ensemble solutions exploit models trained over different languages.

Some ensemble approaches combine form-based and sense-based measures to improve the quality of results. On the one hand, form-based measures are exploited to better capture the dominant sense of the target word  $w$ . On the other hand, sense-based measures are exploited to represent all the meanings of  $w$ , including the minor ones. The combination of CD (see form-based approaches in Section 4.1) and JSD (see sense-based approaches in Section 4.2) is proposed in Martinc et al. [2020c]. As a further ensemble experiment, the results of combining APD, HD, and JSD are discussed

in Wang et al. [2020]. The APD measure is also considered in Rachinskiy and Arefyev [2021], where multiple shift scores are calculated by using different distance metrics (e.g., Manhattan distance, CD, euclidean distance) and these scores are exploited to train a regression model as an ensemble.

Ensemble approaches based on two form-based measures are also proposed. For instance, in Giulianelli et al. [2022], the final semantic shift  $s_w$  is obtained by averaging APD and PRT scores. This is motivated by experimental results where sometimes APD outperforms PRT, while some other times PRT outperforms APD [Kutuzov and Giulianelli, 2020].

Some other ensemble approaches are based on the idea to combine static and contextualised embeddings. The intuition is that static embeddings can capture the dominant sense of the target word  $w$ , better than form-based, contextualised embeddings. In Pömsl and Lyapin [2020], Teodorescu et al. [2022], the semantic shift  $s_w$  is assessed by leveraging both static and contextualised embeddings. In particular,  $s_w$  is determined by the linear combination of the scores obtained by two approaches: i) the APD measure over contextualised embeddings (see form-based approaches in Section 4.1); ii) the CD measure over static embeddings aligned according to the approach described in Hamilton et al. [2016]. Similarly, in Martinc et al. [2020c], instead of directly using the APD measure, JSD is exploited over clusters of contextualised embeddings (see sense-based approaches in Section 4.2). As a further difference, the scores obtained by static and contextualised approaches are combined by multiplication. The intuition is that, since the score distributions of the two approaches are unknown, multiplication prevents an approach from contributing more than the other one in the final score.

CSSDetection approaches can be also combined with grammatical profiles under the intuition that grammatical changes are slow and gradual, while lexical contexts can change very quickly Kutuzov et al. [2021], Giulianelli et al. [2022]. Grammatical profile vectors  $gp_w^1$  and  $gp_w^2$  are associated with the times  $t_1$  and  $t_2$ , respectively, to represent morphological and syntactical features of the considered language in the time period. In Ryzhova et al. [2021], the contextualised embeddings of the word  $w$  occurrences are combined with the grammatical vectors. A linear regression model with regularisation is trained by using as features the cosine similarities over  $\Phi_w^1$  and  $\Phi_w^2$ , and over the grammatical vectors  $gp_w^1$  and  $gp_w^2$ .

As a further ensemble approach, the combination of the time-aware techniques presented in Rosin and Radinsky [2022] and Rosin et al. [2022] (see form-based approaches in Section 4.1) is proposed in order to better inject time into word embeddings.

#### 4.4 Discussion

According to Section 4.1, 4.2, and 4.3, we note that form-based approaches are more popular than sense-based ones. Most papers are characterised by time-oblivious approaches and only a few time-aware approaches have recently appeared (e.g., Rosin and Radinsky [2022]). All approaches leverage unsupervised learning modalities with few exceptions (e.g., Hu et al. [2019]). We argue that the motivation is due to the recent introduction of a reference evaluation framework for semantic shift assessment proposed at SemEval Shared Task 1, where participants were asked to adopt an unsupervised configuration [Schlechtweg et al., 2020].

All papers are featured by contextualised word embeddings extracted from BERT-like models. Regardless of their version (i.e., tiny, small, base, large), BERT and XLM-R are the most frequently used models, and only a few experiments rely on ELMo and RoBERTa. As a matter of fact, the size of data needed to train or fine-tune an XLM-R model is several orders of magnitude greater than BERT. Moreover, even if less frequently employed than BERT, ELMo seems to be promising for CSSDetection and outperform BERT, while being much faster in training and inference [Kutuzov and Giulianelli, 2020]. As a further interesting remark, the use of static *document* embeddings extracted from a Doc2Vec model has been proposed to provide pseudo-contextualised *word* embeddings as an alternative to BERT [Periti et al., 2022].

Monolingual and multilingual language models are both popular. The BERT models are the most frequently used monolingual models. XLM-R models are generally preferred to mBERT (i.e., multilingual BERT) models, since the former are trained on a larger amount of data and languages, thus the intuition is that they can better encode the language usages. Multilingual models are used both in multilingual settings, where corpora of different languages are considered (e.g., Martinc et al. [2020b]), and monolingual settings, where just corpora of one language are given (e.g., in Giulianelli et al. [2022]). In a monolingual setting, the use of a multilingual model is motivated by two reasons: i) a model pre-trained on a specific language is not available (e.g., Kutuzov and Giulianelli [2020]), ii) multilingual models are employed to exploit their cross-lingual transferability property (e.g., Rachinskiy and Arefyev [2021]).

About the type of training, most of the papers directly use pre-trained models or fine-tune them for domain adaptation. Only a few papers propose to exploit a specific fine-tuning (e.g., Pömsl and Lyapin [2020]) or to incrementally fine-tune

a pre-trained model (e.g., Kutuzov and Giulianelli [2020]). Experiments indicate that fine-tuning a pre-trained model for domain adaptation consistently boosts the quality of results when compared against pre-trained models (e.g., Wenjun Qiu and Xu [2022]). The impact of fine-tuning on performance is analysed in Martinc et al. [2020a], where it is shown that optimal results are achieved by fine-tuning a pre-trained model for five epochs and that, after five epochs, performance decreases due to over-fitting. However, we argue that the fine-tuning effectiveness strictly depends on the size and domain of the considered corpora. In many papers, a different number of epochs is proposed with varying results (e.g., Kutuzov and Giulianelli [2020]).

When a transformer-based model is used, contextualised word embeddings are typically extracted from the last one or the last four layers of the model. Experiments show that the semantic features of text are mainly encoded in the last four encoder layers of BERT [Jawahar et al., 2019, Devlin et al., 2019]. In some papers, contextualised embeddings are extracted by aggregating the output of the first and the last encoded layers. In this case, the idea is to combine *surface* features (i.e., phrase-level information [Jawahar et al., 2019]) encoded in the first layer with the semantic features from the last one. Only in Laicher et al. [2021], the standalone use of lower layers of BERT is proposed. Middle layers of BERT are usually excluded since they mainly encode syntactic features Jawahar et al. [2019]. When contextualised embeddings are extracted from more than one layer, they are generally aggregated by average or sum (e.g., Periti et al. [2022]). As an alternative, the use of concatenation is proposed in Kanjirangat et al. [2020].

As a further note, when a BERT-like model is used, some words may be split into word pieces by a subword-based tokenisation algorithm [Sennrich et al., 2016, Wu et al., 2016]. In this case, word piece representations are generally synthesised into a single word representation  $e_{w,k}^j$  through averaging (e.g., Martinc et al. [2020b]), or concatenating (e.g., Martinc et al. [2020c]). As alternative to avoid such problem, the pre-trained vocabulary associated with the model can be extended by adding some words of interest. Then, a fine-tuning step is performed in order to learn the weights associated with the added words (e.g., Rosin et al. [2022]).

Clustering operations are typically exploited in sense-based approaches to perform Word Sense Induction [Lau et al., 2012]. The only form-based approach that relies on clustering is presented in Beck [2020] (see Section 4.1 for details). The clustering algorithms that are most frequently employed are K-Means and affinity propagation (AP). Further considered clustering algorithms are Gaussian Mixture Models (GMMs) (e.g., Rother et al. [2020]), agglomerative clustering (AGG) (e.g., Arefyev and Zhikov [2020]), DBSCAN (e.g., Karnysheva and Schwarz [2020]), HDBSCAN (e.g., Rother et al. [2020]), Balanced Iterative Reducing and Clustering using Hierarchies (BIRCH) (e.g., Rother et al. [2020]), A-Posteriori affinity Propagation (APP) (e.g., Periti et al. [2022]), and Incremental Affinity Propagation based on Nearest neighbor Assignment (IAPNA) (e.g., Periti et al. [2022]). Since K-Means, GMMs, and AGG require to define the number of clusters in advance, the use of a silhouette score is generally employed to determine the optimal number of clusters. As an alternative, the AP algorithm is employed to let emerge the number of clusters without prefixing it. DBSCAN is proposed due to its capability of reducing noise by specifying i) the minimum number of embeddings of each cluster, and ii) the maximum distance  $\epsilon$  between two embeddings in a cluster. HDBSCAN is the hierarchical version of DBSCAN and it can manage clusters of different sizes. As a difference with DBSCAN, HDBSCAN can detect noise without the  $\epsilon$  parameter. APP and IAPNA are incremental extensions of AP, and their use is proposed for semantic shift detection when more than one time interval is considered. In Rother et al. [2020], different clustering algorithms are compared and the experiments show that i) DBSCAN is very sensitive to scale since  $\epsilon$  is prefixed, and ii) BIRCH tends to find a lot of small clusters that are marginal with respect to word meanings.

About the shift functions, a detailed presentation of possible alternatives has been provided in Sections 4.1 and 4.2. As a final remark, we note that CD and APD are frequently exploited in form-based approaches, while JSD is commonly employed in sense-based approaches.

Finally, as for the language of considered corpora, most papers consider the shared benchmark datasets taken from competitive evaluation campaigns (e.g., LSCDiscovery [Zamora-Reina et al., 2022]). Common considered languages are English, German, Latin, and Swedish that appeared in 2020 at SemEval Task 1. Russian appeared in 2021 at RuShiftEval. Spanish appeared in 2022 at LSCDiscovery. The Italian language was introduced in 2020 at DIACRIta, even if it did not receive much attention. The approach described in Martinc et al. [2020b] represents a novel attempt to consider a diachronic corpus containing texts of different languages, namely English and Slovenian.

## 5 Comparison of approaches on performances

In this section, we propose a comparison of the CSSDetection approaches based on their performance obtained by considering the evaluation framework adopted in SSD tasks of shared competitions. The framework is based on a reference benchmark which contains a diachronic textual corpus in a given language. The framework is also

Ref.	SemEval English $C_1 - C_2$	SemEval German $C_1 - C_2$	SemEval Latin $C_1 - C_2$	SemEval Swedish $C_1 - C_2$	GEMS English $C_1 - C_2$	LivFC English $C_1 - C_2$	COHA English $C_1 - C_2$	LSCD Spanish $C_1 - C_2$	DURel German $C_1 - C_2$	SURel German $C_1 - C_2$	$C_1 - C_2$	RSE Russian $C_2 - C_3$	$C_1 - C_3$	NOR Norwegian $C_1 - C_2$ $C_2 - C_3$
Teodorescu et al.	-	-	-	-	-	-	-	ensemble APD 0.573	-	-	-	-	-	-
Zhou and Li	form-based CD 0.392	form-based CD 0.392	form-based CD 0.392	form-based CD 0.392	-	-	-	-	-	-	-	-	-	-
Montariol et al.	sense-based AP + WD 0.456	sense-based AP + JSD 0.583	form-based CD 0.496	sense-based K-Means + WD 0.332	<b>sense-based AP + JSD 0.510</b>	-	-	-	sense-based AP + JSD 0.712	-	-	-	-	-
Peri et al.	sense-based AP + JSD 0.514*	-	sense-based APP + JSD 0.512*	-	-	-	-	-	-	-	-	-	-	-
Pnsl and Lyapin	ensemble APD 0.246	ensemble APD 0.725	ensemble APD 0.463	ensemble APD 0.546	-	-	-	-	ensemble APD <b>0.802</b>	ensemble APD <b>0.723</b>	-	-	-	-
Rachinskiy and Arefyev	-	-	-	-	-	-	-	-	-	-	ensemble APD 0.781	ensemble APD 0.803	ensemble APD 0.822	-
Rachinskiy and Arefyev	-	-	-	-	-	-	-	sense APDP <b>0.745</b>	-	-	-	-	-	-
Rodina et al.	-	-	-	-	-	-	-	-	-	-	form-based PRT 0.557	sense-based AP + JSD 0.406	-	-
Rosin et al.	form-based CD 0.467	form-based CD 0.512	form-based CD 0.512	-	-	<b>form-based TD 0.620</b>	-	-	-	-	-	-	-	-
Rosin and Radinsky	<b>form-based CD 0.627</b>	<b>form-based CD 0.763</b>	<b>form-based CD 0.565</b>	-	-	-	-	-	-	-	-	-	-	-
Rother et al.	sense-based HDBSCAN 0.512	sense-based GMMs 0.605	sense-based GMMs 0.321	sense-based HDBSCAN 0.308	-	-	-	-	-	-	-	-	-	-
Ryzhova et al.	-	-	-	-	-	-	-	-	-	-	ensemble regression 0.480*	ensemble regression 0.487*	ensemble regression 0.560*	-
Kudisov and Arefyev	-	-	-	-	-	-	-	form-based APD 0.637	-	-	-	-	-	-
Kutuzov	form-based APD 0.605	form-based PRT 0.740	form-based PRT 0.561	<b>form-based APD 0.610</b>	sense-based AP + JSD 0.456*	-	-	-	-	-	-	-	-	-
Laicher et al.	form-based APD 0.571*	form-based CD 0.755*	-	form-based APD 0.602*	-	-	-	-	-	-	-	-	-	-
Liu et al.	form-based CD 0.341	form-based CD 0.512	form-based CD 0.304	form-based CD 0.304	form-based CD 0.286	form-based CD 0.561	-	-	-	-	-	-	-	-
Martinc et al.	ensemble AP + JSD 0.361	ensemble AP + JSD 0.642	form-based CD 0.496	ensemble AP + JSD 0.343	-	-	-	-	-	-	-	-	-	-
Giulianelli et al.	-	-	-	-	form-based APD 0.285*	-	-	-	-	-	-	-	-	-
Giulianelli et al.	form-based APD 0.514	ensemble PRT 0.354	ensemble PRT 0.572	ensemble APD 0.397	-	-	-	-	-	-	ensemble APD + PRT 0.376	form-based APD 0.480	form-based APD 0.457	ensemble APD + PRT <b>0.394</b> <b>0.503</b>
Hu et al.	-	-	-	-	-	-	sense-based MNS <b>0.428*</b>	-	-	-	-	-	-	-
Kanjirang et al.	sense-based K-Means + JSD 0.028*	sense-based K-Means + JSD 0.173*	sense-based K-Means + JSD 0.253*	sense-based K-Means + CSC 0.321*	-	-	-	-	-	-	-	-	-	-
Karnysheva and Schwarz	sense-based K-Means + JSD -0.155*	sense-based DBSCAN + JSD 0.388*	sense-based DBSCAN + JSD 0.177*	sense-based K-Means + JSD -0.062*	-	-	-	-	-	-	-	-	-	-
Kashleva et al.	-	-	-	-	-	-	-	sense-based APDP 0.553	-	-	-	-	-	-
Keidar et al.	form-based APD 0.489	-	-	-	-	-	-	-	-	-	-	-	-	-
Arefyev et al.	-	-	-	-	-	-	-	-	-	-	form-based APD <b>0.825</b>	form-based APD <b>0.821</b>	form-based APD <b>0.823</b>	-
Arefyev and Zhikov	sense-based AGG + CD 0.299	sense-based AGG + CD 0.094	sense-based AGG + CD -0.134	sense-based AGG + CD 0.274	-	-	-	-	-	-	-	-	-	-
Beck	form-based CD 0.293*	form-based CD 0.414*	form-based CD 0.343*	form-based CD 0.300*	-	-	-	-	-	-	-	-	-	-
Cuba Gyllensten et al.	form-based CD 0.209*	form-based CD 0.656*	form-based CD 0.399*	form-based CD 0.234*	-	-	-	-	-	-	-	-	-	-
Kutuzov et al.	form-based APD 0.605	form-based PRT 0.740	form-based PRT 0.561	form-based APD 0.569	form-based APD 0.394	-	-	-	-	-	-	-	-	-

Table 6: The Spearman’s correlation score of CSSDetection approaches in selected experiments over corpora of different languages. For each corpus, the top performance is reported in bold. Asterisks denote experiments based on a pre-trained model

characterised by a test-set of target words, where each word is associated with a continuous shift score (i.e., *gold score*) calculated on the basis of manual annotation. Different metrics are also defined in the framework to evaluate the performance of the approaches according to the kind of assessment question that the task aims to address, namely *Grade/Binary Change*, *Sense Gain/Loss* (see Section 2).

In Table 6, we compare the CSSDetection approaches by considering the experiments on *Grade Change Detection* task performed and reported in the corresponding literature papers. In such a kind of task, the Spearman’s correlation score is typically employed for assessing the performance of a given experiment by measuring the correlation between the predicted shift scores and the gold scores<sup>3</sup>. The Spearman’s correlation evaluates the monotonic relationship between the rank-order of the predicted scores and the gold ones. When multiple experiments are discussed in a paper, in Table 6, we report the best Spearman’s correlation score obtained.

<sup>3</sup>In Montariol et al. [2021], as an alternative to the Spearman’s correlation score, the *Discount Cumulative Gain* is proposed Montariol et al. [2021]. However, most papers still use Spearman’s, since it is currently employed in competitive shared tasks.



In the comparison, twelve diachronic corpora are exploited. In particular, we consider: i) the four SemEval datasets for English (SemEval English), German (SemEval German), Latin (SemEval Latin), and Swedish (SemEval Swedish) [Schlechtweg et al., 2020]; ii) the English dataset proposed in Gulordava and Baroni [2011] (GEMS English); iii) the English LiverpoolFC dataset proposed in Del Tredici et al. [2019] (LivFC English); iv) the COHA English dataset (COHA English); v) the LSCDiscovery dataset for Spanish (LSCD Spanish) [Zamora-Reina et al., 2022]; vi) the DUREl dataset for German (DUREl German) [Schlechtweg et al., 2018]; vii) the SUREl dataset for German (SUREl German) [Hätty et al., 2019]; viii) the RuShiftEval dataset for Russian (RSE Russian) [Kutuzov and Pivovarova, 2021b]; and ix) the NorDiaChange dataset for Norwegian (NOR Norwegian) [Kutuzov et al., 2022a]. In Table 6, for each corpus, we highlight when a single time interval  $C_1 - C_2$  or two consecutive time intervals  $C_1 - C_2$  and  $C_2 - C_3$  are considered, respectively. As a further remark, we note that the RSE Russian corpus is the only case where a test-set for the time interval  $C_1 - C_3$  as a whole is provided.

For the sake of readability, the performance according to the Spearman’s correlation scores shown in Table 6 are labeled with the semantic shift function of the considered CSSDetection approach and the corresponding framing with respect to form-based, sense-based, and ensemble-based categories (see Section 4).

As a general remark, we cannot find an approach outperforming all the others on all the considered corpora. This can suggest that an approach is language-dependant, namely it works well on one language and it is not appropriate for others. By relying on the experiments presented in Kutuzov and Giulianelli [2020], we claim that the approaches are not language-dependant and the performance of an approach is influenced by the employed assessment measure in relation with the distribution of the gold scores in the considered test-set. The experiments in Kutuzov and Giulianelli [2020] show that when the distribution of the gold scores is skewed, namely some words are highly shifted and some others are barely shifted, the APD measure achieves better performance on Spearman’s correlation than the PRT measure. On the contrary, when the distribution of the gold scores is almost uniform, namely most of the words are similarly shifted, the PRT measure achieves better performance than the APD measure.

As a further remark, we note that the approaches characterised by fine-tuning achieve greater performance. This is also confirmed in the experiments of Martinc et al. [2020a] where fine-tuning a model boosts the performance when the model is not affected by under or over-fitting.

On average, form-based approaches outperform sense-based approaches in Grade Change Detection tasks. We argue that such a result is motivated by the structure of the test-sets, where just one semantic shift score is provided for each target word. Form-based approaches benefit from this structure since they work on measuring the shift over one general word property (i.e., the dominant sense, or the degree of polysemy). On the opposite, sense-based approaches are disadvantaged by this structure since they work on measuring the shift over multiple word meanings and they need to produce a single, comprehensive shift value that summarises all the single-meaning shifts for the comparison against the gold score. As a result, capturing some (minor) meanings can negatively affect the comprehensive shift value, and to address this issue, small clusters are usually considered as possible noise and filtered out [Martinc et al., 2020c].

Table 6 shows that form-based approaches based on APD, CD, or PRT measures tend to obtain higher performance than sense- and ensemble-based approaches. GEMS English, COHA English, and LSCD Spanish are the only benchmarks where sense-based approaches outperform form-based ones. This can be motivated by the small number of experiments performed. Indeed, for COHA English experiments with form-based approaches have not been tested [Hu et al., 2019], while only a few experiments and a limited number of configurations with form-based approaches have been tested on GEMS English. For LSCD Spanish, the top performance is 0.745 and the corresponding approach leverages the APDP measure, which is an extension of APD characterised by the use of an average-of-average operation. This result is in line with the intuition presented in Periti et al. [2022], where the use of averaging on top of clustering contributes to reduce the noise in the contextualised embeddings of the target word.

We also note that ensemble approaches are on average characterised by high performance. In particular, top performances are provided by ensemble approaches on DUREl German (0.802), SUREl German (0.723), and NOR Norwegian (0.394 and 0.503). It is interesting to observe that the performance on DUREl and SUREl German are obtained through an approach combining static and contextualised word embeddings, thus highlighting that such a kind of combination can be effective. For NOR Norwegian in the time interval  $C_1 - C_2$ , the best approach exploits both APD and PRT; this is a further confirmation that APD and PRT are top-performing measures in semantic shift detection. For the subsequent time interval  $C_2 - C_3$ , the best result on NOR Norwegian is obtained with a combination of APD with grammatical profiles. This is a confirmation of the intuition presented in Giulianelli et al. [2022], which suggests that ensembling grammatical profiles with contextualised embeddings can enhance performance by incorporating morphological and syntactic features not fully captured by contextualised models.

For SemEval English, SemEval German, and SemEval Latin, the top performance are 0.627, 0.763, and 0.565, respectively, and they are obtained by the time-aware approach proposed in Rosin and Radinsky [2022]. Also for

LivFC English (0.620), the top performance is obtained by leveraging a time-aware approach [Rosin et al., 2022]. We argue that extra-linguistic information (e.g., time information) can have a positive impact on performance. The injection of extra-linguistic information can contribute to increase the performance also when small-size language models are employed, since they are less affected by noise than larger models. As a confirmation, in contrast to the widespread belief that the larger the models the higher the performance, the best result for SemEval English is obtained by exploiting contextualised embeddings extracted from a BERT-tiny model [Turc et al., 2019, Rosin and Radinsky, 2022]. This is also true for SemEval Swedish (0.610), where the top performance is obtained by calculating the APD measure over contextualised embeddings extracted from an ELMo model [Kutuzov, 2020], which is far smaller than BERT-like models.

Finally, we note that the use of supervised learning modalities contributes to achieve high performance. As an example, the top performances for RSE Russian are 0.825 on  $C_1 - C_2$ , 0.821 on  $C_2 - C_3$ , and 0.823 on  $C_1 - C_3$  and they are obtained by a form-based, supervised approach [Arefyev et al., 2021].

## 6 Scalability, interpretability, and robustness issues

In this section, we analyse the CSSDetection approaches by considering possible scalability, interpretability, and reliability issues.

### 6.1 Scalability issues

In the CSSDetection approaches, any occurrence of the target word considered for shift assessment is represented by a specific embedding. As a basic implementation, all the contextualised embeddings are stored in memory for processing. The higher the number of occurrences of a target word, the higher the number of embeddings to manage. As a result, when the size of the diachronic corpus grows, possible issues arise both in terms of memory and computation time. Similar issues occur when multiple target words are considered for shift assessment. In this case, a possible workaround for addressing the memory issue is to process one target word at a time. However, in this way, the memory issue *shifts* to a computation time issue. For feasibility convenience, most experiments work on a small set of target words. This kind of limitations inhibits the possibility to address tasks like the detection of the most changed word in a corpus. The need to work on solutions capable of dealing with such a kind of scalability issues has recently been promoted in LSCDiscovery, where participants were asked to assess the semantic shift on all the words of the dictionary [Zamora-Reina et al., 2022].

Some possible solutions to the scalability issues have been proposed in literature. For instance, approaches based on measures that enforce aggregation by averaging (e.g., CD, PRT) are time-scalable, since only the prototypes are considered for shift assessment instead of the whole set of embeddings. Also approaches based on APD or JSD measures can be adjusted to become time-scalable. In particular, the number of embeddings to store and process can be reduced by random sampling the occurrences of the target word  $w$ . This means that i) a smaller number of similarity scores needs to be calculated with APD (e.g., Ryzhova et al. [2021]), and ii) JSD works on top of clustering algorithms that converge faster (e.g., Rodina et al. [2020]). As an alternative to random sampling, an online *aggregation by summing* method is proposed in Montariol et al. [2021], where a pre-fixed number of contextualised embeddings  $n$  is stored in memory. An embedding  $e_w$  is stored when the number of embeddings in memory is less than  $n$  and  $e_w$  is strongly dissimilar from all the other embeddings previously stored. If  $e_w$  is not stored, it is aggregated to the most similar embedding stored in memory through sum.

The dimensionality reduction of the embeddings is proposed as a further alternative to enforce scalability. For example, in Rother et al. [2020], the embedding dimensionality is reduced to 10 (from 768) by combining an autoencoder with the UMAP (Uniform Manifold Approximation and Projection) algorithm [McInnes et al., 2018]. In Keidar et al. [2022], UMAP and PCA are used to project contextualised embedding into  $h \in \{2, 5, 10, 20, 50, 100\}$  dimensions. With respect to this solution, we argue that, although it can improve the memory scalability, time scalability is negatively affected since dimensionality reduction takes time. However, in Rother et al. [2020], it is shown that the dimensionality reduction can still contribute to time scalability when the goal is to test and compare the effectiveness of different clustering algorithms and the reduced embeddings are saved and re-used. As a further option, the use of small language models, such as TinyBert or ELMo, is gaining more and more attention since the dimension of the generated embeddings is far lower (e.g., Rosin and Radinsky [2022]).

Scalability issues can also arise when the shift needs to be assessed on a corpus  $C = \bigcup_i^n C_i$  defined over more than one time interval ( $n > 2$ ). Typically, CSSDetection approaches calculate the shift score  $s_w$  over each pair of time intervals  $(t_i, t_{i+1})$  by iteratively re-applying the same assessment workflow. As a difference, an incremental approach based on a clustering algorithm called *A Posteriori affinity Propagation* (APP) is proposed in Periti et al. [2022] to speed up

the aggregation stage. In each time interval, clustering is incrementally executed by considering the prototypes of the previous time period (i.e., aggregation by averaging) and the incoming embeddings of the current time period.

## 6.2 Interpretability issues

Interpretability issues arise when it is not possible to determine which meaning(s) have changed among all the meanings of a target word, namely the meaning(s) that mainly caused the shift score assessed by a considered approach. Definitely, form-based approaches are affected by such a kind of issues, since they model the shift as the change in the dominant sense or in the degree of polysemy of a word, without considering the possible multiple meanings. On the opposite, sense-based approaches aim at providing an interpretation of the word change, since they attempt to model the shift by considering the multiple word senses. However, interpretability issues can arise also when sense-based approaches are employed due to three main motivations.

*Word meaning representation.* Sense-based approaches mostly rely on clustering techniques to represent word meanings. The K-Means and the AP clustering algorithms are usually employed to this end. K-Means requires that the number of target clusters is prefixed, and this can be inappropriate to effectively represent the meanings of a target word that are not known beforehand. AP lets the number of target clusters emerge, but experimental results show that the association of a cluster with a word meaning can be imprecise. We argue that this can be due to the distributional nature of contextualised models that tends to capture changes in contextual variance (i.e., word usages) rather than changes in lexicographic senses (i.e., word meanings) [Kutuzov et al., 2022b]. As an example, sometimes AP produces more than 100 clusters, which is rather unrealistic if we assume that a cluster represents a word meaning [Periti et al., 2022]. As a matter of fact, a word may completely change its context without changing its meaning [Martinc et al., 2020a].

*Word meaning description.* Each cluster obtained during the aggregation stage of a sense-based approach needs to be associated with a description that denotes the corresponding word meaning. This can be done by human experts on the basis of the cluster contents. However, this is time-consuming, given that a cluster can consist of several hundreds/thousands of elements. As an alternative, clustering analysis techniques have been proposed to label clusters by summarising their contents. As a possible option, a cluster description can be extracted from the content by considering the top featuring keywords based on lexical occurrences (e.g., Tf-Idf) [Kellert and Mahmud Uz Zaman, 2022, Montariol et al., 2021]. In Giulianelli et al. [2020], the sense-prototype of a cluster is proposed as a cluster exemplar and the corpus sentences that are closest to the prototype are adopted as cluster/meaning description. However, when a cluster contains outliers, these sentences could not provide an effective description.

*Word meaning evolution.* When a corpus  $C = \bigcup_i^n$  defined over more than one time interval is considered, the clusters defined at a time step  $t_i$  need to be linked to the clusters of the previous time-step  $t_{i-1}$  to trace the evolution of the corresponding meaning over time (i.e., cluster/meaning history). Since the clustering executions at each time-step are independent, the capability of recognising corresponding clusters/meanings at different time-steps can be challenging. As a possible solution, alignment techniques can be employed to link similar word meanings in different, consecutive time periods [Kanjirangat et al., 2020, Montariol et al., 2021]. As a further option, evolutionary clustering algorithms can be exploited without requiring any alignment mechanism across time periods [Periti et al., 2022].

## 6.3 Robustness issues

Robustness issues arise when the assessment score is not reliable due to data imbalance, model stability, and model bias.

*Data imbalance.* The diachronic corpus  $C$  must equally reflect the presence of the target word  $w$  in both the time steps  $t_1$  and  $t_2$ . This means that the frequency of  $w$  must not strongly change in the considered time period. However, in common scenarios, more documents are available for the most recent time step  $t_2$  and “it may not be possible to achieve balance in the sense expected from a modern corpus” [Tahmasebi et al., 2021]. As a consequence, the frequency of  $w$  can be strongly higher in  $t_2$  than in  $t_1$  and the embeddings  $\Phi_w$  can produce a distorted representation of the target word when the model is trained/fine-tuned [Zhou et al., 2021, Wendlandt et al., 2018]. As a further remark, data imbalance issues can occur when some word meanings are more frequent than others. For instance, the dominant sense is usually more represented than other senses in the corpus  $C$ . As a result, when a sense-based approach is adopted, the embedding distributions  $p_w^1, p_w^2$  can be skewed, meaning that a larger number of embeddings is associated with the dominant sense rather than with the other minor senses. In sense-based approaches, the word meanings are represented by clusters, and *the number of clusters consistently reflects word frequency* [Kutuzov, 2020]. When a meaning is associated with a few embeddings/clusters, its contribution to the overall assessment score is marginally leading to an inflated or underestimated assessment score. In this respect, a qualitative analysis of “potentially erroneous” outputs of CSSDetection approaches is presented in Kutuzov et al. [2022b]. Some examples of potentially erroneous assessment scores occur when i) a word with strongly context-dependent meanings is considered, whose embeddings are mutually different; ii) a word is frequently used in a very specific context in only one time step  $t_1$  or  $t_2$ ; iii) a word is affected by a

*syntactic change*, not a semantic one. In Liu et al. [2021], a solution is proposed to reduce the false discovery rate and to improve the precision of the shift assessment by leveraging permutation-based statistical test and term-frequency thresholding.

**Model stability.** Contextualised pre-trained models are usually trained on modern text sources. For example, the original English BERT model is pre-trained on Wikipedia and BooksCorpus [Zhu et al., 2015]. As a result, pre-trained models are prone to represent words from a modern perspective, and thus they tend to ignore the temporal information of a considered corpus. This way, when historical corpora are considered, the possible obsolete word usages cannot be properly represented. This problem has been investigated in the literature by comparing the performance of pre-trained against fine-tuned models [Kutuzov and Giulianelli, 2020, Wenjun Qiu and Xu, 2022]. In line with the considerations of Section 4.4, the results show that fine-tuning the model on the whole diachronic corpus improves the quality of word representations for historical texts. Since fine-tuning the model can be expensive in terms of time and computational resources, a measure for estimating the model effectiveness for historical sources is presented in Ishihara et al. [2022]. In particular, in Ishihara et al. [2022], this measure is used to decide whether a model should be re-trained or fine-tuned.

**Model bias** The contextualised embeddings can possibly be affected by biases on the encoded information. For instance, a possible bias can arise from orthographic information, such as the word form and the position of a word in a sentence, since they influence the output of the top BERT layers [Laicher et al., 2021]. Text pre-processing techniques are proposed as a solution to reduce the influence of orthography in the embeddings, thus increasing the robustness of encoded semantic information. To this end, lower-casing the corpus text is a commonly-employed solution. However, *the lower-casing of words often conflates parts of speech*, thus another possible bias can raise. For example, the proper noun `Apple` and the common noun `apple` become identical after lower-casing [Hengchen et al., 2021]. The possible bias introduced by Named Entities and proper nouns is investigated in Laicher et al. [2021], Martinc et al. [2020c]. In Wenjun Qiu and Xu [2022], text normalisation techniques are proposed based on the removal of accent markers. In some languages, such a kind of normalisation can introduce a bias since different words can be conflated. For example, `papà` (e.g., the Italian word for dad) and `papa` (e.g., the Italian word for pope) cannot be distinguished after the accent removal. Further text pre-processing techniques can be employed to reduce the possible bias due to orthographic information. In Schlechtweg et al. [2020], lemmatisation and punctuation removal are proposed. Experimental results on lemmatisation for reducing the model bias on BERT embeddings are presented in Laicher et al. [2021]. Further experiments show that lemmatising the target word alone is more beneficial than lemmatising the whole corpus [Laicher et al., 2021]. Filtering out unimportant words, such as stop words and low-frequency words, can be also beneficial [Zhou and Li, 2020]. As an alternative solution to reduce word-form biases, the embedding of a word occurrence can be computed by averaging its original embedding and the embeddings of its nearest words in the input sentence [Zhou and Li, 2020].

When aggregation by clustering is enforced, the possible word-form biases can affect the clustering result [Laicher et al., 2021]. As a solution, clustering refinement techniques have been proposed. As an option, the removal of the clusters containing only one or two instances is adopted, since they are not considered significant [Martinc et al., 2020c]. As a further option, in Martinc et al. [2020a], clusters with less than two members are considered as weak clusters and they are merged with the closest strong cluster, i.e. cluster with more than two members. In Periti et al. [2022], clusters containing less than 5 percent of the whole set of embeddings are assumed to be poorly informative and are thus dropped. However, we argue that the use of clustering refinement techniques must be carefully considered since also small clusters can be important when the corpus is unbalanced in the number of meanings of a word.

## 7 Challenges and concluding remarks

In this survey, we analysed the CSSDetection task by providing a formal definition of the problem and a reference classification framework based on meaning representation, time awareness, and learning modality dimensions. The literature approaches are surveyed according to the given framework by considering the assessment function, the language model, the achieved performance, and the possible scalability/interpretability/robustness issues.

In Hengchen et al. [2021], an overview of open challenges about computational SSD is presented. In the following, we extend such an overview by focusing on those challenges that are specific to CSSDetection in relation with the issues discussed in Section 6.

**Scalability.** The trend in CSSDetection is to adopt increasingly larger models with the idea that they better represent language features. As a consequence, scalability issues arise and they are being addressed as discussed in Section 6.1. However, contrary to this trend, we argue that the use of small-size models, such as those introduced in Rosin and Radinsky [2022], Rosin et al. [2022], needs to be further explored since they are competitive in terms of performance.

**Word meaning representation.** In Section 5, we show that form-based approaches outperform sense-based approaches in the Grade Change Detection assessment. However, we argue that sense-based approaches are promising since they focus on encoding word senses and they can enrich the mere degree of semantic shift of a word  $w$  with the information about the specific meaning of  $w$  that changed. In this direction, the SSD should be considered as a temporal/diachronic extension of other problems such as Word Sense Induction [Alsulaimani et al., 2020], Word Meaning Disambiguation [Godbole et al., 2022], and Word-in-Context [Loureiro et al., 2022]. So far, word senses have been represented through aggregation by clustering under the idea that each cluster represents a specific word meaning. However, according to the interpretability issues of Section 6, clustering techniques are often affected by noise and they are typically capable of representing word usages rather than word meanings. Thus, further investigations are required to represent lexicographic meanings in a more faithful way. The possible integration of the linguistic theory should be also considered for sense modeling. For example, the linguistics theory could contribute to tailor the behavior of CSSDetection approaches by focusing on a specific type of word sense to consider (e.g., standard word meaning, topic use, pragmatics, connotation) [Hengchen et al., 2021]. As a further example, the linguistics theory can clarify the specific word meaning captured by a given cluster representation.

**Word meaning description.** According to Section 6, current solutions to meaning description are focused on determining a representative label taken from the cluster contents (e.g., Tf-Idf, sentence(s) featuring the sense-prototype). Such solutions are mostly oriented to highlight the lexical features of the cluster/meaning without considering any element that reflects the cluster’s semantics. As a consequence, open challenges are based on the need of comprehensive description techniques capable of capturing both lexical and semantic aspects such as position in text, semantics, or co-occurrences across different documents.

**Word meaning evolution.** In shared competitions, the reference evaluation framework for CSSDetection is based on one/two time periods that are considered for shift detection. The extension of the evaluation framework to consider more time periods is an open challenge. In particular, methods and practices of CSSDetection approaches need to be tested/extended for detecting both short- and long-term semantic shifts and for promoting the design of incremental techniques able to handle dynamic corpora (i.e., corpora that become progressively available).

In this context, a further challenge is about the capability to trace the change of a meaning over multiple time steps (i.e., meaning evolution). As mentioned in Section 2, alignment techniques can be used to link similar word meanings in different, consecutive time periods. However, such a solution is not completely satisfactory due to possible limitations (e.g., scalability, robustness of alignment) and further research work is needed to better track the meaning evolution over time (e.g., Periti et al. [2022]).

**Model stability.** Most of the approaches surveyed in this paper are time-oblivious and face the problem of model stability through fine-tuning. Since this practice can be expensive in terms of time and resources, we argue that further research on the development of time-aware approaches is needed, in that, they do not suffer the model stability problem.

**Model bias.** The solutions to model bias issues presented in Section 6 are language-dependent and they are mainly exploited in approaches based on monolingual contextualised model. Further research work is needed to test the effectiveness of existing solutions also in approaches based on multilingual contextualised models. In addition, we argue that future work should concern the application of denoising and debiasing techniques to both monolingual and multilingual embedding models (e.g., Kaneko and Bollegala [2021]) with the aim to improve CSSDetection performance by reducing orthographic biases regardless of the language(s) on which the models were trained.

**Further challenges** not strictly related to the issues of Section 6 are the following:

*Semantic Shift Interpretation.* Most of the literature papers do not investigate the nature of the detected shifts, meaning that they do not classify the semantic shifts according to the existing linguistic theory (e.g., amelioration, pejoration, broadening, narrowing, metaphorisation, metonymisation, and metonymy) [Campbell, 2013, Hock and Joseph, 2019]. Further studies on the causes and types of semantic changes are needed. These studies could be crucial to detect “laws” of semantic shift that describe the condition under which the meanings of words are prone to change. For example, some laws are hypothesised in Xu and Kemp [2015], Dubossarsky et al. [2015], Hamilton et al. [2016], but later the validity of some of them has been questioned [Dubossarsky et al., 2017]. Contextualised embeddings could contribute to test the validity of current laws and to propose new ones. To the best of our knowledge, some steps in this direction are only moved in Hu et al. [2019] for modeling the word change from an ecological viewpoint.

*Computational models of meaning change.* Almost all experiments on CSSDetection are based on BERT embeddings. Although there are open questions about how to maximise the effectiveness of BERT embeddings in different language setups, the effectiveness of BERT for tracing semantic shifts has been extensively investigated. We believe that CSSDetection should be extended by considering a wider range of contextualised embedding models. Some work explored the effectiveness of ELMo [Kutuzov and Giulianelli, 2020, Rodina et al., 2020]. However, the performance of ELMo in different contexts and setups should be analysed in more detail. Furthermore, it might be worth investigating smaller versions of BERT, like ALBERT [Lan et al., 2019] and DistilBERT [Sanh et al., 2019]. Further models can also be considered like seq2seq and generative models, which recently showed interesting results in the field of temporal Word-in-Context problem [Lyu et al., 2022].

*Multilingual models.* In past shared competitions on SSD, monolingual models have generally been preferred to multilingual ones. We believe that a systematic comparison of monolingual vs. multilingual models is required to determine scenarios and conditions where the former type of models provides better performance than the latter type or vice-versa. Multilingual embeddings can also contribute to CSSDetection since they could enable a language-independent semantic shift assessment, meaning that the gold-scores of different languages can be exploited as a whole for the evaluation of a given approach.

*Cross-language shift detection.* As introduced in Martinc et al. [2020b], further investigations are required to address the problem of cross-language shift detection. We argue that solutions to such a kind of problem can be also useful for CSSDetection since they can detect semantic change of *cognates* and *borrowings* (e.g., [Fourrier and Montariol, 2022]), as well as *contact-induced* semantic shifts (e.g., [Miletic et al., 2021])<sup>4</sup>.

*Use cases.* So far, detecting semantic shifts through contextualised embeddings is still a theoretical problem not yet integrated in real application scenarios like historical information retrieval, lexicography, linguistic research, or social-analysis. For this reason, further use cases and experiences must be developed and shared.

*Context Shift over different domains.* The attention gained by diachronic semantic shift detection through the use of word embeddings paved the way for modeling other linguistics issues such as the identification of diatopic lexical variation [Seifart, 2019], the detection of semantic shifts of grammatical constructions [Fonteyn et al., 2020], or the comparison of how speakers who disagree on a subject use the same words [Garí Soler et al., 2022]. The CSSDetection approaches can be tested and possibly extended to cope with such a kind of linguistics issues.

## Acknowledgments

We would like to thank Nina Tahmasebi for her valuable comments and constructive feedback on the manuscript.

## References

- Daphna Keidar, Andreas Opedal, Zhijing Jin, and Mrinmaya Sachan. Slangvolution: A Causal Analysis of Semantic Change and Frequency Dynamics in Slang. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1422–1442, Dublin, Ireland, May 2022. Association for Computational Linguistics (ACL).
- Silvana Castano, Alfio Ferrara, Stefano Montanelli, and Francesco Periti. Semantic Shift Detection in Vatican Publications: a Case Study from Leo XIII to Francis. In *Proc. of SEBD*, pages 231–243, Pisa, Italy, June 2022. CEUR-WS.
- Hosein Azarbondy, Mostafa Dehghani, Kaspar Beelen, Alexandra Arkut, Maarten Marx, and Jaap Kamps. Words Are Malleable: Computing Semantic Shifts in Political and Media Discourse. In *Proc. of the ACM International Conference on Information and Knowledge Management (CIKM)*, page 1509–1518, New York, NY, USA, 2017. Association for Computing Machinery (ACM).

---

<sup>4</sup>In linguistics, cognates are sets of words in different languages that have been inherited in direct descent from an etymological ancestor in a common parent language. Borrowings (or loanwords) are words adopted by the speakers of one language from a different language. Contact-induced semantic shifts are diachronic changes within a recipient language that are traceable to languages other than the direct ancestor of the recipient language and that have spread and are conventionalised within a community speaking the recipient language.

- Reem Alatrash, Dominik Schlechtweg, Jonas Kuhn, and Sabine Schulte im Walde. CCOHA: Clean Corpus of Historical American English. In *Proc. of the Language Resources and Evaluation Conference (LREC)*, pages 6958–6966, Marseille, France, May 2020. European Language Resources Association (ELRA).
- Yvonne Adesam, Dana Dannélls, and Nina Tahmasebi. Exploring the Quality of the Digital Historical Newspaper Archive KubHist. In *Proc. of the Digital Humanities Conference (DHN)*, pages 9–17, Copenhagen, Denmark, 2019.
- Dominik Schlechtweg, Barbara McGillivray, Simon Hengchen, Haim Dubossarsky, and Nina Tahmasebi. SemEval-2020 Task 1: Unsupervised Lexical Semantic Change Detection. In *Proc. of the Workshop on Semantic Evaluation (SemEval)*, pages 1–23, Barcelona (online), December 2020. International Committee for Computational Linguistics (ICCL).
- Barbara McGillivray and Adam Kilgarriff. Tools for historical corpus research, and a corpus of Latin. *New Methods in Historical Corpus Linguistics*, 1(3):247–257, 2013.
- Pierpaolo Basile, Giovanni Semeraro, and A. Caputo. Kronos-it: a dataset for the italian semantic change detection task. In *CLiC-it*, 2019.
- Andrey Kutuzov and Lidia Pivovarova. Three-part Diachronic Semantic Change Dataset for Russian. In *In Proc. of the International Workshop on Computational Approaches to Historical Language Change (LChange)*, pages 7–13, (online), August 2021a. Association for Computational Linguistics (ACL).
- Frank D. Zamora-Reina, Felipe Bravo-Marquez, and Dominik Schlechtweg. LSCDiscovery: A Shared Task on Semantic Change Discovery and Detection in Spanish. In *Proc. of the Workshop on Computational Approaches to Historical Language Change (LChange)*, pages 149–164, Dublin, Ireland, May 2022. Association for Computational Linguistics (ACL).
- Jing Chen, Emmanuele Chersoni, and Chu-ren Huang. Lexicon of Changes: Towards the Evaluation of Diachronic Semantic Shift in Chinese. In *Proc. of the Workshop on Computational Approaches to Historical Language Change (LChange)*, pages 113–118, Dublin, Ireland, May 2022. Association for Computational Linguistics (ACL).
- Andrey Kutuzov, Samia Touileb, Petter MÅihlum, Tita Enstad, and Alexandra Wittemann. NorDiaChange: Diachronic Semantic Change Dataset for Norwegian. In *Proc. of the Language Resources and Evaluation Conference*, pages 2563–2572, Marseille, France, June 2022a. European Language Resources Association (ELRA).
- Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient Estimation of Word Representations in Vector Space. In Yoshua Bengio and Yann LeCun, editors, *Proc. of the International Conference on Learning Representations (ICLR)*, Scottsdale, Arizona, 2013.
- Xuri Tang. A State-of-the-Art of Semantic Change Computation, 2018.
- Andrey Kutuzov, Lilja Øvrelid, Terrence Szymanski, and Erik Velldal. Diachronic Word Embeddings and Semantic Shifts: a Survey. In *Proc. of the International Conference on Computational Linguistics (COLING)*, pages 1384–1397, Santa Fe, New Mexico, USA, August 2018. Association for Computational Linguistics (ACL).
- Nina Tahmasebi, Lars Borin, and Adam Jatowt. Survey of Computational Approaches to Lexical Semantic Change Detection. In *Computational approaches to semantic change*. Language Science Press, June 2021.
- Nina N. Tahmasebi. *Models and Algorithms for Automatic Detection of Language Evolution*. PhD thesis, Gottfried Wilhelm Leibniz Universität Hannover, 2013.
- Sunny Mitra, Ritwik Mitra, Suman Kalyan Maity, Martin Riedl, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. An Automatic Approach to Identify Word Sense Changes in Text Media Across Timescales. *Natural Language Engineering*, 21(5):773–798, 2015.
- Pierpaolo Basile, Annalina Caputo, Tommaso Caselli, Pierluigi Cassotti, and Rossella Varvara. DIACR-Ita@ EVALITA2020: Overview of the EVALITA2020 DiachronicLexical Semantics (DIACR-Ita) Task. In *Proc. of the Evaluation Campaign of Natural Language Processing and Speech Tools for Italian (EVALITA)*, Online, December 2020. CEUR-WS.
- Andrey Kutuzov and Lidia Pivovarova. RuShiftEval: A Shared Task on Semantic Shift Detection for Russian. In *Proc. of the International Conference on Computational Linguistics and Intellectual Technologies (Dialogue)*, number 20, (online), 2021b. Redkollegija sbornika.
- John R. Firth. A Synopsis of Linguistic Theory. *Studies in linguistic analysis*, 1957.
- Mario Giulianelli, Marco Del Tredici, and Raquel Fernández. Analysing Lexical Semantic Change with Contextualised Word Representations. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 3960–3973, (online), July 2020. Association for Computational Linguistics (ACL).

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proc. of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics (ACL).
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A Robustly Optimized BERT Pretraining Approach, 2019.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep Contextualized Word Representations. In *Proc. of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 2227–2237, New Orleans, Louisiana, June 2018. Association for Computational Linguistics (ACL).
- Matej Martinc, Syrielle Montariol, Elaine Zosa, and Lidia Pivovarov. *Capturing Evolution in Word Usage: Just Add More Clusters?*, page 343–349. Association for Computing Machinery (ACM), New York, NY, USA, 2020a.
- Vani Kanjirang, Sandra Mitrovic, Alessandro Antonucci, and Fabio Rinaldi. SST-BERT at SemEval-2020 Task 1: Semantic Shift Tracing by Clustering in BERT-based Embedding Spaces. In *Proc. of the Workshop on Semantic Evaluation (SemEval)*, pages 214–221, Barcelona (online), December 2020. International Committee for Computational Linguistics (ICCL).
- Francesco Periti, Alfio Ferrara, Stefano Montanelli, and Martin Ruskov. What is Done is Done: an Incremental Approach to Semantic Shift Detection. In *Proc. of the Workshop on Computational Approaches to Historical Language Change (LChange)*, pages 33–43, Dublin, Ireland, May 2022. Association for Computational Linguistics (ACL).
- Julia Rodina, Yuliya Trofimova, Andrey Kutuzov, and Ekaterina Artemova. ELMo and BERT in Semantic Change Detection for Russian, October 2020.
- Matej Martinc, Petra Kralj Novak, and Senja Pollak. Leveraging Contextual Embeddings for Detecting Diachronic Semantic Shift. In *Proc. of the Language Resources and Evaluation Conference (LREC)*, pages 4811–4819, Marseille, France, May 2020b. European Language Resources Association (ELRA).
- Martin Pömsl and Roman Lyapin. CIRCE at SemEval-2020 Task 1: Ensembling Context-Free and Context-Dependent Word Representations. In *Proc. of the Workshop on Semantic Evaluation (SemEval)*, pages 180–186, Barcelona (online), December 2020. International Committee for Computational Linguistics (ICCL).
- Wenjun Qiu and Yang Xu. HistBERT: A Pre-trained Language Model for Diachronic Lexical Semantic Analysis, 2022.
- Severin Laicher, Gioia Baldissin, Enrique Castañeda, Dominik Schlechtweg, and Sabine Schulte. CL-IMS @ DIACR-Ita: Volente o Nolente: BERT does not Outperform SGNS on Semantic Change Detection. In *Proc. of the Evaluation Campaign of Natural Language Processing and Speech Tools for Italian (EVALITA)*, pages 438–443, Marrakech, Morocco, December 2020. CEUR-WS.
- Renfen Hu, Shen Li, and Shichen Liang. Diachronic Sense Modeling with Deep Contextualized Word Embeddings: An Ecological View. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 3899–3908, Florence, Italy, July 2019. Association for Computational Linguistics (ACL).
- Guy D. Rosin, Ido Guy, and Kira Radinsky. Time Masking for Temporal Language Models. In *Proc. of the ACM International Conference on Web Search and Data Mining (WSDM)*, page 833–841, (online), 2022. Association for Computing Machinery (ACM).
- Nikolay Arefyev, Maksim Fedoseev, Vitaly Protastov, Daniil Homiskiy, Adis Davletov, and Alexander Panchenko. DeepMistake: Which Senses are Hard to Distinguish for a Word-in-Context Model. In *Proc. of the Conference on Computational Linguistics and Intellectual Technologies (Dialogue)*, (online), 2021.
- Christin Beck. DiaSense at SemEval-2020 Task 1: Modeling Sense Change via Pre-trained BERT Embeddings. In *Proc. of the Workshop on Semantic Evaluation (SemEval)*, pages 50–58, Barcelona (online), December 2020. International Committee for Computational Linguistics (ICCL).
- Franziska Horn. Exploring Word Usage Change with Continuously Evolving Embeddings. In *Proc. of the Annual Meeting of the Association for Computational Linguistics and the International Joint Conference on Natural Language Processing (ACL-IJCNLP 2021)*, pages 290–297, (online), August 2021. Association for Computational Linguistics (ACL).
- Valentin Hofmann, Janet Pierrehumbert, and Hinrich Schütze. Dynamic Contextualized Word Embeddings. In *Proc. of the Annual Meeting of the Association for Computational Linguistics and the International Joint Conference on Natural Language Processing (ACL-IJCNLP 2021)*, pages 6970–6984, (online), August 2021. Association for Computational Linguistics (ACL).



- Jinan Zhou and Jiaxin Li. TemporalTeller at SemEval-2020 Task 1: Unsupervised Lexical Semantic Change Detection with Temporal Referencing. In *Proc. of the Workshop on Semantic Evaluation (SemEval)*, pages 222–231, Barcelona (online), December 2020. International Committee for Computational Linguistics (ICCL).
- Guy D. Rosin and Kira Radinsky. Temporal Attention for Language Models. In *Findings of the Association for Computational Linguistics (NAACL 2022)*, pages 1498–1508, Seattle, United States, July 2022. Association for Computational Linguistics (ACL).
- Andrey Kutuzov and Mario Giulianelli. UiO-UvA at SemEval-2020 Task 1: Contextualised Embeddings for Lexical Semantic Change Detection. In *Proc. of the Workshop on Semantic Evaluation (SemEval)*, pages 126–134, Barcelona (online), December 2020. International Committee for Computational Linguistics (ICCL).
- Artem Kудisov and Nikolay Arefyev. BOS at LSCDiscovery: Lexical Substitution for Interpretable Lexical Semantic Change Detection. In *Proc. of the Workshop on Computational Approaches to Historical Language Change (LChange)*, pages 165–172, Dublin, Ireland, May 2022. Association for Computational Linguistics (ACL).
- Severin Laicher, Sinan Kurtyigit, Dominik Schlechtweg, Jonas Kuhn, and Sabine Schulte im Walde. Explaining and Improving BERT Performance on Lexical Semantic Change Detection. In *Proc. of the Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop (EACL)*, pages 192–202, (online), April 2021. Association for Computational Linguistics (ACL).
- Benyou Wang, Emanuele Di Buccio, and Massimo Melucci. University of Padova @ DIACR-Ita. In *Proc. of the Evaluation Campaign of Natural Language Processing and Speech Tools for Italian (EVALITA 2020)*, Marrakech, Morocco, December 2020. CEUR-WS.
- Andrey Kutuzov. *Distributional Word Embeddings in Modeling Diachronic Semantic Change*. PhD thesis, University of Oslo, 2020.
- Anastasiia Ryzhova, Daria Ryzhova, and Ilya Sochenkov. Detection of Semantic Changes in Russian Nouns with Distributional Models and Grammatical Features. In *Proc. of the Conference on Computational Linguistics and Intellectual Technologies (Dialogue)*, (online), 2021.
- Yuri Kuratov and Mikhail Arkhipov. Adaptation of Deep Bidirectional Multilingual Transformers for Russian Language, 2019.
- Yang Liu, Alan Medlar, and Dorota Glowacka. Statistically Significant Detection of Semantic Shifts using Contextual Word Embeddings. In *Proc. of the Workshop on Evaluation and Comparison of NLP Systems (Eval4NLP)*, pages 104–113, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics (ACL).
- David Bamman and Patrick J. Burns. Latin BERT: A Contextual Language Model for Classical Philology, 2020.
- Mario Giulianelli, Andrey Kutuzov, and Lidia Pivovarova. Do Not Fire the Linguist: Grammatical Profiles Help Language Models Detect Semantic Change. In *Proc. of the Workshop on Computational Approaches to Historical Language Change (LChange)*, pages 54–67, Dublin, Ireland, May 2022. Association for Computational Linguistics (ACL).
- Syrielle Montariol, Matej Martinc, and Lidia Pivovarova. Scalable and Interpretable Semantic Change Detection. In *Proc. of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 4642–4652, (online), June 2021. Association for Computational Linguistics (ACL).
- Tony Finch. Incremental Calculation of Weighted Mean and Variance. *University of Cambridge*, 4(11-5):41–42, 2009.
- Maja Rudolph and David Blei. Dynamic Embeddings for Language Evolution. In *Proc. of the World Wide Web Conference (WWW)*, page 1003–1011, Lyon, France, 2018. International World Wide Web Conferences Steering Committee (IW3C2).
- Ziqian Zeng, Xin Liu, and Yangqiu Song. Biased Random Walk Based Social Regularization for Word Embeddings. In *Proc. of the International Joint Conference on Artificial Intelligence (IJCAI)*, page 4560–4566, Stockholm, Sweden, 2018. AAAI Press.
- Xiaolei Huang and Michael J. Paul. Neural Temporality Adaptation for Document Classification: Diachronic Word Embeddings and Domain Adaptation Models. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 4113–4123, Florence, Italy, July 2019. Association for Computational Linguistics (ACL).
- Paul Röttger and Janet Pierrehumbert. Temporal Adaptation of BERT and Performance on Downstream Document Classification: Insights from Social Media. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2400–2412, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics (ACL).

- Alessio Ferrari, Beatrice Donati, and Stefania Gnesi. Detecting Domain-Specific Ambiguities: An NLP Approach Based on Wikipedia Crawling and Word Embeddings. In *Proc. of the IEEE International Requirements Engineering Conference Workshops (REW)*, pages 393–399. IEEE, 2017.
- Haim Dubossarsky, Simon Hengchen, Nina Tahmasebi, and Dominik Schlechtweg. Time-Out: Temporal Referencing for Robust Modeling of Lexical Semantic Change. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 457–470, Florence, Italy, July 2019. Association for Computational Linguistics (ACL).
- Mohammad Taher Pilehvar and Jose Camacho-Collados. WiC: the Word-in-Context Dataset for Evaluating Context-Sensitive Meaning Representations. In *Proc. of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 1267–1273, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics (ACL).
- Maxim Rachinskiy and Nikolay Arefyev. Zeroshot Crosslingual Transfer of a Gloss Language Model for Semantic Change Detection. In *Proc. of the Conference on Computational Linguistics and Intellectual Technologies (Dialogue)*, (online), 2021.
- Maxim Rachinskiy and Nikolay Arefyev. GlossReader at LSCDiscovery: Train to Select a Proper Gloss in English – Discover Lexical Semantic Change in Spanish. In *Proc. of the Workshop on Computational Approaches to Historical Language Change (LChange)*, pages 198–203, Dublin, Ireland, May 2022. Association for Computational Linguistics (ACL).
- Nikolay Arefyev and Vasily Zhikov. BOS at SemEval-2020 Task 1: Word Sense Induction via Lexical Substitution for Lexical Semantic Change Detection. In *Proc. of the Workshop on Semantic Evaluation (SemEval)*, pages 171–179, Barcelona (online), December 2020. International Committee for Computational Linguistics (ICCL).
- Kseniia Kashleva, Alexander Shein, Elizaveta Tukhtina, and Svetlana Vydrina. HSE at LSCDiscovery in Spanish: Clustering and Profiling for Lexical Semantic Change Discovery. In *Proc. of the Workshop on Computational Approaches to Historical Language Change (LChange)*, pages 193–197, Dublin, Ireland, May 2022. Association for Computational Linguistics (ACL).
- Matej Martinc, Syrielle Montariol, Elaine Zosa, and Lidia Pivovarov. Discovery Team at SemEval-2020 Task 1: Context-sensitive Embeddings Not Always Better than Static for Semantic Change Detection. In *Proc. of the Workshop on Semantic Evaluation (SemEval)*, pages 67–73, Barcelona (online), December 2020c. International Committee for Computational Linguistics (ICCL).
- Anna Karnysheva and Pia Schwarz. TUE at SemEval-2020 Task 1: Detecting Semantic Change by Clustering Contextual Word Embeddings. In *Proc. of the Workshop on Semantic Evaluation (SemEval)*, pages 232–238, Barcelona (online), December 2020. International Committee for Computational Linguistics (ICCL).
- Amaru Cuba Gyllensten, Evangelia Gogoulou, Ariel Ekgren, and Magnus Sahlgren. SenseCluster at SemEval-2020 Task 1: Unsupervised Lexical Semantic Change Detection. In *Proc. of the Workshop on Semantic Evaluation (SemEval)*, pages 112–118, Barcelona (online), December 2020. International Committee for Computational Linguistics (ICCL).
- David Rother, Thomas Haider, and Steffen Eger. CMCE at SemEval-2020 Task 1: Clustering on Manifolds of Contextualized Embeddings to Detect Historical Meaning Shifts. In *Proc. of the Workshop on Semantic Evaluation (SemEval)*, pages 187–193, Barcelona (online), December 2020. International Committee for Computational Linguistics (ICCL).
- Paul Cook, Jey Han Lau, Diana McCarthy, and Timothy Baldwin. Novel Word-sense Identification. In *Proc. of the International Conference on Computational Linguistics (COLING)*, pages 1624–1635, Dublin, Ireland, August 2014. Association for Computational Linguistics (ACL).
- Terra Blevins and Luke Zettlemoyer. Moving Down the Long Tail of Word Sense Disambiguation with Gloss Informed Bi-encoders. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1006–1017, (online), July 2020. Association for Computational Linguistics (ACL).
- Chang Ming-Wei, Lev Ratinov, Dan Roth, and Vivek Srikumar. Importance of Semantic Representation: Dataless Classification. In *Proc. of the National Conference on Artificial Intelligence (AAAI)*, volume 2, 2008.
- Hyunjin Choi, Judong Kim, Seongho Joe, Seungjai Min, and Youngjune Gwon. Analyzing Zero-shot Cross-lingual Transfer in Supervised NLP Tasks. In *Proc. of the International Conference on Pattern Recognition (ICPR)*, pages 9608–9613, 2021.
- Nicolas Bonneel, Michiel van de Panne, Sylvain Paris, and Wolfgang Heidrich. Displacement Interpolation Using Lagrangian Mass Transport. In *Proc. of the SIGGRAPH Asia Conference*, New York, NY, USA, 2011. Association for Computing Machinery (ACM).
- Daniela Teodorescu, Spencer von der Ohe, and Grzegorz Kondrak. UAlberta at LSCDiscovery: Lexical Semantic Change Detection via Word Sense Disambiguation. In *Proc. of the Workshop on Computational Approaches to*

- Historical Language Change (LChange)*, pages 180–186, Dublin, Ireland, May 2022. Association for Computational Linguistics (ACL).
- Andrey Kutuzov, Erik Velldal, and Lilja Øvrelid. Contextualized embeddings for semantic change detection: Lessons learned. In *Proc. of the Northern European Journal of Language Technology (NEJLT)*, volume 8, 2022b.
- William L. Hamilton, Jure Leskovec, and Dan Jurafsky. Diachronic Word Embeddings Reveal Statistical Laws of Semantic Change. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1489–1501, Berlin, Germany, August 2016. Association for Computational Linguistics (ACL).
- Andrey Kutuzov, Lidia Pivovarov, and Mario Giulianelli. Grammatical Profiling for Semantic Change Detection. In *Proc. of the Conference on Computational Natural Language Learning (CoNLL)*, pages 423–434, (online), November 2021. Association for Computational Linguistics (ACL).
- Ganesh Jawahar, Benoît Sagot, and Djamel Seddah. What Does BERT Learn about the Structure of Language? In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 3651–3657, Florence, Italy, July 2019. Association for Computational Linguistics (ACL).
- Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural Machine Translation of Rare Words with Subword Units. In *Proc. of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1715–1725, Berlin, Germany, August 2016. Association for Computational Linguistics (ACL).
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation, 2016.
- Jey Han Lau, Paul Cook, Diana McCarthy, David Newman, and Timothy Baldwin. Word Sense Induction for Novel Sense Detection. In *Proc. of the Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 591–601, Avignon, France, April 2012. Association for Computational Linguistics. URL <https://aclanthology.org/E12-1060>.
- Kristina Gulordava and Marco Baroni. A Distributional Similarity Approach to the Detection of Semantic Change in the Google Books Ngram Corpus. In *Proc. of the Workshop on Geometrical Models of Natural Language Semantics (GEMS)*, pages 67–71, Edinburgh, UK, July 2011. Association for Computational Linguistics (ACL).
- Marco Del Tredici, Raquel Fernández, and Gemma Boleda. Short-Term Meaning Shift: A Distributional Exploration. In *Proc. of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2069–2075, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics (ACL).
- Dominik Schlechtweg, Sabine Schulte im Walde, and Stefanie Eckmann. Diachronic Usage Relatedness DUREl: A Framework for the Annotation of Lexical Semantic Change. In *Proc. of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 169–174, New Orleans, Louisiana, June 2018. Association for Computational Linguistics (ACL).
- Anna Häty, Dominik Schlechtweg, and Sabine Schulte im Walde. SUREl: A Gold Standard for Incorporating Meaning Shifts into Term Extraction. In *Proc. of the Joint Conference on Lexical and Computational Semantics (SEM)*, pages 1–8, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics (ACL).
- Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Well-Read Students Learn Better: The Impact of Student Initialization on Knowledge Distillation. 2019.
- Leland McInnes, John Healy, and James Melville. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction, 2018.
- Olga Kellert and Md Mahmud Uz Zaman. Using Neural Topic Models to Track Context Shifts of Words: a Case Study of COVID-related Terms before and after the Lockdown in April 2020. In *Proc. of the Workshop on Computational Approaches to Historical Language Change (LChange)*, pages 131–139, Dublin, Ireland, May 2022. Association for Computational Linguistics (ACL).
- Kaitlyn Zhou, Kawin Ethayarajh, and Dan Jurafsky. Frequency-based Distortions in Contextualized Word Embeddings, 2021.
- Laura Wendlandt, Jonathan K. Kummerfeld, and Rada Mihalcea. Factors Influencing the Surprising Instability of Word Embeddings. In *Proc. of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 2092–2102, New Orleans, Louisiana, June 2018. Association for Computational Linguistics (ACL).

- Yukun Zhu, Ryan Kiros, Richard Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Aligning Books and Movies: Towards Story-like Visual Explanations by Watching Movies and Reading Books, 2015.
- Shotaro Ishihara, Hiromu Takahashi, and Hono Shirai. Semantic Shift Stability: Efficient Way to Detect Performance Degradation of Word Embeddings and Pre-trained Language Models. In *Proc. of the Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 205–216, Online only, November 2022. Association for Computational Linguistics (ACL).
- Simon Hengchen, Nina Tahmasebi, Dominik Schlechtweg, and Haim Dubossarsky. Challenges for computational lexical semantic change, June 2021.
- Ashjan Alsulaimani, Erwan Moreau, and Carl Vogel. An Evaluation Method for Diachronic Word Sense Induction. In *Findings of the Association for Computational Linguistics (EMNLP)*, pages 3171–3180, (online), November 2020. Association for Computational Linguistics (ACL).
- Mihir Godbole, Parth Dandavate, and Aditya Kane. Temporal Word Meaning Disambiguation using TimeLMs, 2022.
- Daniel Loureiro, Aminette D’Souza, Areej Nasser Muhajab, Isabella A. White, Gabriel Wong, Luis Espinosa-Anke, Leonardo Neves, Francesco Barbieri, and Jose Camacho-Collados. TempoWiC: An Evaluation Benchmark for Detecting Meaning Shift in Social Media. In *Proc. of the International Conference on Computational Linguistics (COLING)*, pages 3353–3359, Gyeongju, Republic of Korea, October 2022. International Committee on Computational Linguistics (ICCL).
- Masahiro Kaneko and Danushka Bollegala. Debiasing Pre-trained Contextualised Embeddings. In *Proc. of the Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 1256–1266, (online), April 2021. Association for Computational Linguistics (ACL).
- Lyle Campbell. *Historical linguistics*. Edinburgh University Press, 2013.
- Hans Henrich Hock and Brian D Joseph. Language history, language change, and language relationship. In *Language History, Language Change, and Language Relationship*. De Gruyter Mouton, 2019.
- Yang Xu and Charles Kemp. A Computational Evaluation of Two Laws of Semantic Change. In David C. Noelle, Rick Dale, Anne S. Warlaumont, Jeff Yoshimi, Teenie Matlock, Carolyn D. Jennings, and Paul P. Maglio, editors, *Proc. of the Annual Meeting of the Cognitive Science Society (CogSci)*, Pasadena, California, USA, 2015.
- Haim Dubossarsky, Yulia Tsvetkov, Chris Dyer, and Eitan Grossman. A Bottom Up Approach to Category Mapping and Meaning Change. In *NetWordS*, pages 66–70. CEUR-WS, 2015.
- Haim Dubossarsky, Daphna Weinshall, and Eitan Grossman. Outta Control: Laws of Semantic Change and Inherent Biases in Word Representation Models. In *Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1136–1145, Copenhagen, Denmark, September 2017. Association for Computational Linguistics (ACL).
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. ALBERT: A Lite BERT for Self-Supervised Learning of Language Representations, 2019.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. DistilBERT, a Distilled Version of BERT: Smaller, Faster, Cheaper and Lighter, 2019.
- Chenyang Lyu, Yongxin Zhou, and Tianbo Ji. MLLabs-LIG at TempoWiC 2022: A Generative Approach for Examining Temporal Meaning Shift. In *Proc. of the EMNLP 2022*, 2022.
- Clémentine Fourier and Syrielle Montariol. Caveats of Measuring Semantic Change of Cognates and Borrowings using Multilingual Word Embeddings. In *Proc. of the Workshop on Computational Approaches to Historical Language Change (LChange)*, pages 97–112, Dublin, Ireland, May 2022. Association for Computational Linguistics (ACL).
- Filip Miletic, Anne Przewozny-Desriaux, and Ludovic Tanguy. Detecting Contact-Induced Semantic Shifts: What Can Embedding-Based Methods Do in Practice? In *Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 10852–10865, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics (ACL).
- Frank Seifart. *Contact-induced change*, pages 13–23. De Gruyter Mouton, Berlin, Boston, 2019.
- Lauren Fonteyn, F Karsdorp, B McGillivray, A Nerghens, and M Wevers. What About Grammar? Using BERT Embeddings to Explore Functional-Semantic Shifts of Semi-Lexical and Grammatical Constructions. In *Proc. of the Workshop on Computational Humanities Research (CHR)*, pages 257–268, Amsterdam, the Netherlands, 2020. CEUR-WS.

---

Aina Garí Soler, Matthieu Labeau, and Chloé Clavel. One Word, Two Sides: Traces of Stance in Contextualized Word Representations. In *Proc. of the International Conference on Computational Linguistics (COLING)*, pages 3950–3959, Gyeongju, Republic of Korea, October 2022. International Committee on Computational Linguistics (ICCL).