# To Word Senses and Beyond: Inducing Concepts with Contextualized Language Models

Bastien Liétard and Pascal Denis and Mikaela Keller

University of Lille, Inria, CNRS, Centrale Lille,

UMR 9189 – CRIStAL, F-59000 Lille, France

`first-name.last-name@inria.fr`

February 8, 2026

**Abstract**

Polysemy and synonymy are two crucial inter-related facets of lexical ambiguity. While both phenomena are widely documented in lexical resources and have been studied extensively in NLP, leading to dedicated systems, they are often being considered independently in practical problems. While many tasks dealing with polysemy (e.g. Word Sense Disambiguation or Induction) highlight the role of word's senses, the study of synonymy is rooted in the study of concepts, i.e. meanings shared across the lexicon. In this paper, we introduce Concept Induction, the unsupervised task of learning a soft clustering among words that defines a set of concepts directly from data. This task generalizes Word Sense Induction. We propose a bi-level approach to Concept Induction that leverages both a local lemma-centric view and a global cross-lexicon view to induce concepts. We evaluate the obtained clustering on SemCor's annotated data and obtain good performance (BCubed F1 above 0.60). We find that the local and the global levels are mutually beneficial to induce concepts and also senses in our setting. Finally, we create static embeddings representing our induced concepts and use them on the Word-in-Context task, obtaining competitive performance with the State-of-the-Art.

# 1 Introduction

A crucial challenge in understanding natural language comes from the fact that the mapping between word forms and lexical meanings is many-to-many, due to polysemy (i.e., the multiplicity of meanings for a

given form)[1] and synonymy (i.e., the multiplicity of forms for expressing a given meaning). Both polysemy and synonymy have been thoroughly studied in NLP, but mostly as independent problems, giving rise to dedicated systems. Thus, Word Sense Disambiguation (WSD) aims at correctly mapping word occurrences to one of its senses [Raganato et al.2017], while Word Sense Induction (WSI), its unsupervised counterpart, aims at clustering word occurrences into latent senses directly from data [Manandhar et al.2010, Jurgens and Klapaftis2013]. More recently, researchers have proposed the task of Word-in-Context (WiC), which consists in classifying pairs of word occurrences depending on whether they realize the same sense or not [Pilehvar and Camacho-Collados2019]. All these works take a word centric view, which aims at identifying or characterizing the different senses of a given word, where these senses are bound to a word. Another line of work, which takes a broader lexicon-wide perspective, is concerned with identifying synonyms, which are equivalence classes over different words that point to the same concept [Zhang et al.2021, Ghanem et al.2023], where concepts are semantic entities that are not bound to a word. In WordNet [Miller1995, Fellbaum1998], concepts are called synsets, defined as sets of synonyms. However, outside of lexical resources, synonymy and polysemy are usually considered as independent problems in the NLP literature. Yet, these two views are complementary. In lexicology, they correspond to two perspectives on the word-meaning mapping: semasiology and onomasiology. The former is the word-to-meanings view, where one can observe polysemy by looking at the different meanings a given word has. The latter is the meaning-to-words view, in which one can study synonymy by looking at the inventory of words that speakers use to express the same meaning.

In this paper, we propose a new task, called Concept Induction, that directly aims at learning concepts in an unsupervised manner from raw text. More precisely, this task aims at learning a soft clustering over a target lexicon (i.e., a set of words), in such a way that each cluster corresponds to a (latent) concept. Thus, this task both addresses polysemy (since polysemous words should appear in multiple clusters) and synonymy (since synonymous words should appear in the same cluster(s)). Inducing concepts can be interesting for many external applications, like building lexical resources for low-resources languages [Velasco et al.2023], and can bring a different perspective in computational studies of meaning, moving the usual word-centric focus to a more meaning-centric state.

Our approach to Concept Induction relies on word occurrences for a target lexicon, represented as word embeddings derived from a Contextualized Language Model (in this case, BERT Large [Devlin et al.2019]),

---

[1]We take polysemy in its most comprehensive definition, also including homonymy.

which are then grouped, using hard clustering algorithms, into concept denoting clusters. While these concept clusters could in principle be obtained directly from word occurrences, we propose a bi-level methodology that leverages both a local, lemma-centric clustering (i.e., operating on only specific word occurrences), and a global, cross-lexicon clustering (i.e., operating on all words occurrences). From this perspective, our approach generalizes, and in fact builds upon classical Word Sense Induction, in that word senses are learned jointly alongside with concepts. We hypothesize that an approach taking both complementary resolutions in account will lead to improved Concept Induction and Word Sense Induction, i.e. that the two objectives can be mutually beneficial.

To validate our approach, we carried out experiments on the SemCor dataset, which provides a set of concepts (taking the form of WordNet synsets) related to word occurrences. We found that our bi-level clustering approach accurately learn concepts, achieving F1 scores above 0.60 on the task of Concept Induction compared to WordNet's synsets, outperforming competing approaches that use only local and global views. This demonstrates the benefits of our bi-level approach, and its ability to leverage both local and global views when inducing concepts. Interestingly, we show that the benefits go both ways: our proposed approach outperforms lemma-centric approaches when evaluated for WSI. Finally, we show that concept-aware static embeddings derived from our approach are also competitive with state-of-the-art approaches on the Word-in-Context task, while using less training data. Through the new task of concept induction, we also contribute in a new way to the ongoing debate regarding the ability to align vector representations extracted from Contextualized Language Models to the semantic representations posited by (psycho-)linguists. In this vein, we conduct a qualitative evaluation of obtained clusters to ensure they indeed reflect concepts and gather synonyms. The source code we used for experiments is available at
`https://github.com/blietard/concept-induction`.

## 2 Related Work

### 2.1 Lexical resources for concepts

Princeton's WordNet (PWN) [Miller1995, Fellbaum1998] is a lexical database that has been the most widely used as a reference for most word-sense-related tasks for many years. In WordNet, the entry corresponding to a lemma has different word senses, each of them mapping to a synset. Synsets are WordNet's equivalents of our concepts. Lemmas whose word senses belong to the same synset are synonymous. Word-

Net 3.0 contains 117,659 synsets and is built from the work of psycholinguists and lexicographers, that not only describes synonymy but also other lexical relations such as hypernymy/hyponymy, antonymy, meronymy/holonymy, etc. But the amount of resources needed to create such lexical databases with human experts is considerable, making them a very rare and precious resource. They are not available for a large number of active languages, and even more rare for dead languages [Bizzoni et al.2014, Khan et al.2022].

## 2.2 Word senses with Language Models

With the recent development of neural Contextualized Language Models (CLM), several work use their hidden-layers to extract vector representations of word usages and retrieve word senses. These representations are fed to a classification (for WSD) or a clustering (in the case of WSI) algorithm to distinguish the word's senses [Scarlini et al.2020, Nair et al.2020, Saidi and Jarray2023]. These embeddings-based approaches have applications in other fields: [Kutuzov and Giulianelli2020] and [Martinc et al.2020] use sense clusters found using CLM embeddings to study the change in meaning of words, and [Chronis and Erk2020] propose a many-Kmeans method to investigate semantic similarity and relatedness. Another line of work uses list of substitute tokens sampled from the CLM head to infer senses [Amrami and Goldberg2019, Eyal et al.2022] and are successful on WSI benchmarks like [Manandhar et al.2010] and [Jurgens and Klapaftis2013].

## 2.3 Structures of Meaning in CLM

Recent research probes neural CLMs for alignments between representations from their latent spaces and semantic patterns and relations. Section 7.2 of [Haber and Poesio2024] summarizes findings about polysemy in contextualized CLMs, showing that these models were able to detect polysemy and in some cases distinguish actual polysemy from homonymy. They report that representations from different senses may however overlap. [Hanna and Marecek2021] shows that pretrained BERT embed knowledge of hypernymy but is limited to the more common hyponyms.

[Velasco et al.2023] build on top of WSI techniques in an attempt to automatically construct a WordNet for Filipino, thus proposing a modeling of synonymy in this language. However, the evaluation of the synsets they obtained is limited by the lack of sense-annotated data for Filipino, and they could not evaluate the impact of their methodology on the two levels (senses and concepts).

Works like [Ethayarajh2019] and [Chronis and Erk2020] study the kind of information that was dis-

tributed across layers. The former concludes that syntactic and word-order information are distributed in the first layers while in deeper layers, representations are heavily influenced by contexts. The latter demonstrates, with a multi-prototypes embedding approach, that semantic similarity is best found in moderately late layers, while relatedness is best found in last layers.

# 3   Concept Induction

Our main motivation behind Concept Induction is to present a view of the mapping between words and their meaning(s).[2] This view is systemic, meaning that it should not be defined for individual words neither for individual concepts, but rather acknowledging these as a whole with interactions and relations. This extends beyond the primary objective of WSI, which defines word senses as pertaining to individual words only and does not explore relations between lemmas or concepts.

## 3.1   Basic notions

Consider a set of target words (or lemmas) and for each lemma, we have a set of occurrences of this word in a context (e.g., a sentence or a phrase). The set of target lemmas is referred to as the lexicon, while the corpus is the set of all occurrences. Our goal is to study the meaning of target words as they are used in the corpus.

In this study we call *sense* of a word its usage to refer to a concept. A polysemous word has multiple senses, each of them referring to a distinct concept. Two words are said to be synonyms for a given concept when each of them has one of their senses referring to this shared concept. Senses are defined "locally", i.e. bound to an individual word of the lexicon, as opposed to concepts which are defined "globally", i.e. across the whole lexicon. An occurrence of a word $w$ realizes one of its senses.

Consider the words "test" and "trial" and the following corpus: (A) the jury found them guilty in a fair trial. (B) candidates competed in a trial of skill. (C) the hero underwent a test of strength. The corpus is composed of two occurrences of "trial" and one occurrence of "test." In the corpus, "trial" is polysemous. Its first sense, illustrated in A, refers to a process of law. Its second sense, in B, refers to the concept of the act of undergoing testing. The sense of "test" in sentence C also corresponds to this concept: it's a case

---

[2]This mapping is called *patterns of lexification* by [François2022]; see also *coexpression* and *synexpression* in the terminology proposed by [Haspelmath2023].

where "test" and "trial" are synonymous. Shifting the focus from senses to concepts, we will say that B and C instantiate the same concept, while A is an instance of a different concept.

## 3.2 Task definition

The goal of Concept Induction (CI) is to automatically learn a set of concepts directly from the data, i.e. learning a soft clustering $\mathcal{C}_W$ in the set of target words $W$ that should correspond to the multiple concepts instantiated by occurrences of the corpus. $\mathcal{C}_W$ is a soft clustering because a word can be assigned to several clusters (when it is polysemous). Using a different perspective than WSI, the framework of Concept Induction provides a more complete view on meaning across the lexicon. Both WSI and CI capture polysemy, but CI also reveals synonymy across the lexicon. Like WSI, Concept Induction does not require a pre-defined set of concepts.

## 3.3 Formal framework

Let $W$ be the lexicon. For all word $w$ in $W$, we denote $o_i^w$ the $i$-th occurrence of $w$ in the corpus. We define $\mathcal{O}_w = \{o_i^w\}_{i \leq n_w}$ the set of $n_w$ occurrences of $w$. The corpus, denoted $\mathcal{O}$, is the union of all $\mathcal{O}_w$.

For a given word $w \in W$, the set $\mathcal{O}_w$ can be partitioned according to its different senses. We denote $s_j^w$ the part of occurrences of $w$ in the corpus corresponding to the $j$-th sense of $w$. We refer to these groups of occurrences as the *sense clusters* of $w$. The set $\mathcal{S}_w = \{s_j^w\}_{j \leq n_w^s}$ forms a partition of $\mathcal{O}_w$, and we call $\mathcal{S}$ the set of all sense clusters of all words, i.e., $\mathcal{S} = \bigcup_{w \in W} \mathcal{S}_w$. $\mathcal{S}$ is a "local" (lemma-centric) partition of the whole $\mathcal{O}$. The task of Word Sense Induction aims at learning the partition $\mathcal{S}$ given a corpus $\mathcal{O}$.

In this work, we aim at dividing the corpus into concepts instead of senses. We denote $c_k$, the group of occurrences of words corresponding to the concept indexed by $k$, and $\mathcal{C} = \{c_k\}_{k \leq p}$ the partition of $\mathcal{O}$ in $p$ concept clusters. Unlike sense clusters of $\mathcal{S}$, a concept cluster $c_k \in \mathcal{C}$ can gather occurrences of different words: $\mathcal{C}$ is a "global" partition. Each occurrence $o_i^w$ of a word $w \in W$ is associated to a sense cluster $s_j^w$ and a concept cluster $c_h \in \mathcal{C}$. We can say that a concept corresponding to $c_h$ is instantiated by occurrence $o_i^w$ through the sense corresponding to $s_j^w$, or conversely that $o_i^w$ uses the sense reflected in $s_j^w$ to mean the concept described by concept cluster $c_h$. All occurrences of sense cluster $s_j^w \in \mathcal{S}$ appear in the same concept cluster $c_p \in \mathcal{C}$.

In summary, $\mathcal{S}$ and $\mathcal{C}$ are partitions of $\mathcal{O}$ and are naturally constrained as follows:

IMAGE NOT PROVIDED
Illustration of our framework. The words "trial" is polysemous and has two senses corresponding to two different concepts, and is synonym with "test" for this second meaning.

Figure 1: Illustration of our framework. The words "trial" is polysemous and has two senses corresponding to two different concepts, and is synonym with "test" for this second meaning.

1. By definition, a sense in $\mathcal{S}$ is associated to one and only one word $w \in W$.

2. An occurrence $o_i^w$ realizes exactly one sense $s_j^w \in \mathcal{S}$.

3. An occurrence $o_i^w$ instantiates exactly one concept $c_p \in \mathcal{C}$.

4. In a given sense $s_j^w \in \mathcal{S}$, all occurrences are assigned to the same concept $c_p \in \mathcal{C}$.

5. All $s_j^w \in \mathcal{S}_w$ (i.e. same word) refer to distinct concepts.

From the partition $\mathcal{C}$ on occurrences, one can derive $\mathcal{C}_W$, a clustering of the set of words $W$ into concepts. To each concept cluster $c_k \in \mathcal{C}$ we associate a cluster in $\mathcal{C}_W$ that contains all lemmas of $W$ whose occurrences were assigned to $c_k$. In $\mathcal{C}_W$, a polysemous word with $n$ senses appears in $n$ distinct clusters (one per sense), and synonyms appear in at least one common cluster (one per shared concept).

We denote $\widehat{\mathcal{C}_W}$ the word-level soft-clustering and $\widehat{\mathcal{C}}$ the partition of occurrences that we learned on the data.

In Figure 1 we illustrate this framework, using a corpus of occurrences of the words "test" and "trial". In this scenario, $W = \{$"test", "trial"$\}$ and two concepts are instantiated: a process of law to determine someone's guilt and a challenge to evaluate a skill. The lemma "trial" exhibits two senses as it has occurrences corresponding to both concepts: "trial" is polysemous. The second concept is also instantiated by occurrences of "test", therefore "trial" and "test" show synonymy in this case. This toy example also follows all constraints formulated above.

## 4  Methodology

In this section we describe the methods we propose and evaluate for Concept Induction. We learn a clustering $\widehat{\mathcal{C}_W}$ drawing inspiration from the relations between $\mathcal{O}$, $\mathcal{S}$, $\mathcal{C}$ and $\mathcal{C}_W$. In particular, the overall objective of our methodology consist in finding $\mathcal{C}$ (i.e. partition occurrences into concept clusters) to derive $\mathcal{C}_W$. Section 3.3 highlighted that there are two levels of partitions: a local level (senses) and a global one (concepts). The

proposed approaches rely on both levels and the use of a Contextualized Language Model (CLM) to gather representations of occurrences influenced by the context.

## 4.1 Proposed Bi-level Method

**Local (lemma-centric) clustering** Firstly, we propose to learn a word-sense partition for each target words individually. Using the CLM hidden layers, we extract a vector representation (the occurrence embedding) of every occurrence $o_i^w$. We then learn a partition $\widehat{\mathcal{S}}_w$ of each $\mathcal{O}_w$ using a clustering algorithm on the embeddings. Each $\widehat{\mathcal{S}}_w$ describes the locally estimated sense clusters of word $w$. Jointly considering these partitions for all $w \in W$, we obtain a partition $\widehat{\mathcal{S}}$ of the whole set of occurrences $\mathcal{O}$. This partition is local in the sense that each word has its occurrences clustered independently from other words.

**Global (cross-lexicon) clustering** Once we have a local clustering $\widehat{\mathcal{S}}$, we turn from considering words independently to consider all words together. In this step, we learn a global clustering by merging local clusters of occurrences. To do so, we average embeddings of all occurrences in the same local cluster to get a single embedding representing each local cluster. Then we run a second clustering algorithm, this time using the averaged embeddings of local clusters. This global clustering defines a new partition $\widehat{\mathcal{C}}$ of the the corpus $\mathcal{O}$: when two local clusters $s_j^1$ and $s_j^2$ are merged into the same global cluster $c_k$ (because their embeddings were clustered together), all their occurrences are assigned to global cluster $c_k$. From this global occurrence partition $\widehat{\mathcal{C}}$ we can easily extract $\widehat{\mathcal{C}}_W$, a word-level soft-clustering of lemmas whose occurrences appear in the same $c_k$.

This Bi-level method directly implements the system of constraints described in Section 3.3. Only constraint 5 is not enforced by design. Indeed, our local clusters being learned and not informed by an expert, the local clustering step may make errors, especially if the data for a given word are sparse. Allowing the global clustering to merge local clusters enables the correction of local clustering's recall errors using information from the global level.

We also want to highlight that the proposed methodology is generic, in the sense that it is not tied to a specific choice of clustering algorithm.

## 4.2 Local-only and Global-only

Sense-inducing systems (WSI approaches) that create only local clusters of occurrences for each word are said to be Local-only systems. We use them as baseline models that only produce word-level clusters of size 1 and do not reflect synonymy, but still learn polysemy.

On the other hand, consider a system in which each occurrence is mapped to its own local cluster (i.e. no actual local clustering step), and the global step divides occurrences directly into global clusters. We refer to this kind of system as Global-only approaches. They allow to evaluate how useful the local clustering step is in the process: we hypothesize that the local step in Bi-level will reduce potential variance in occurrences by aggregating them, increasing Precision compared to Global-only.

# 5 Experiments

In this section, we evaluate the abilities of the proposed methods to induce concepts and compare the proposed bi-level approach to other methods. We investigate the advantages of the bi-level approach not only for the global viewpoint but also in the local setting.

## 5.1 Settings

**Data.** We choose to use the annotated part of the SemCor 3.0[3] corpus. This dataset contains occurrences for a wide number of words, and morpho-syntactic annotations provide their lemma and their Part-of-Speech tag. Among all lemmas having at least 10 annotated occurrences, we keep only nouns (excluding proper nouns)[4] composed only of alphabetical characters with a minimum length of 3 letters. The resulting lexicon $W$ contains 1,560 different lemmas, for which we gather a corpus $\mathcal{O}$ containing a total of 52,997 occurrences[5]. SemCor is also semantically annotated, with each occurrence of a target lemma assigned to a synset in WordNet, that we consider to be the concept it refers to. We derive a reference partition of the occurrences $\mathcal{C}$ and a reference soft-clustering of the words $\mathcal{C}_W$ from annotations, for a total of 3,855 different concepts (WordNet's synsets) covered in $\mathcal{O}$. This set of concepts is the subset of WordNet corresponding to

---

[3]https://www.eecs.umich.edu/~mihalcea/downloads.html\#semcor

[4]For the sake of simplicity and clarity, this study is focused only on nouns. Indeed, other Parts-of-Speech induce extra difficulties. Verbs for instance required extra preprocessing steps and decisions (e.g. include or exclude gerundive uses, past participle employed like adjectives, etc.). Extension of experiments to other PoS is left to future work.

[5]Sentences in which the lemma appears, paired with its position within them. If the lemma appears multiple times in the same sentence, we create several distinct occurrences, where only the position varies.

the textual data.

**Evaluation of Concept Induction** We compare the learned word clustering $\widehat{\mathcal{C}}_W$ to the reference $\mathcal{C}_W$. We choose to use the BCubed metrics [Bagga and Baldwin1998], obtaining Precision and Recall for the evaluated clustering compared to the reference, as well as an F1 score. To account for overlapping clusters, we use the Extended BCubed metrics proposed by [Amigó et al.2009], which has already been used as evaluation in SemEval2013 WSI task [Jurgens and Klapaftis2013].

Using BCubed metrics, for a given evaluated clustering, low precision would mean that grouped lemmas should not have been clustered together because none of their occurrence annotations map to a shared concept according to annotations. Low recall means that the evaluated system fails to capture clusters of lemmas whose occurrences share a concept according to annotations. The number of common clusters between two words also impacts BCubed metrics: if two lemmas appear together in too many clusters compared to the reference clustering, precision is decreased; if the number of common clusters is too low, recall is decreased.

**Development.** To learn the clustering, candidate systems have access to the full set of occurrences-in-context but not their annotations. To choose the appropriate set of hyperparameters, we create a Dev split of the annotations by randomly sampling 10% of concepts and revealing semantic annotations of the corresponding occurrences. We use them to evaluate Concept Induction for this small set of concepts, and choose the set of hyperparameters that scores best in BCubed F1.

**Evaluation splits** In the final evaluation phase, we compute scores on all concepts/all occurrences, including the Dev split, as concepts in it are part of the whole subset of WordNet described by SemCor's annotations. In the full data, we found that 88% of the concepts were instantiated using only a single lemma. To better evaluate cases of synonymy, we also evaluate systems on a subset of the corpus, denoted "Synon", that contains only occurrences of concepts showing synonymy (the remaining 12% of concepts, instantiated through at least 2 distinct lemmas). Statistics are provided in Table 5 in Appendix B. Note that it only changes the set of concepts/lemmas for which the system is being evaluated, not the clustering's training data.

## 5.2   Systems and baselines

**Clustering Algorithms.** We try two different clustering algorithms relying on different paradigms: Kmeans (used in [Chronis and Erk2020]), a centroid-based algorithm with a fixed number of clusters, and Agglom-

erative clustering (used in [Saidi and Jarray2023, Velasco et al.2023]; dubbed "Agglo" for short), a deterministic hierarchical approach using a distance threshold to create a dynamic number of clusters instead of using a fixed one. Another difference between Kmeans and Agglo is that the former assumes that expected clusters are of nearly-spherical shape and balanced in number of points, while the latter does not make assumptions on the shape of data. Details of tested hyperparameter values are provided in Appendix C.

**Representations.** Following [Chronis and Erk2020] and [Eyal et al.2022], we use BERT Large [Devlin et al.2019], a masked language model with 24 layers and 345M parameters. This allows for direct comparisons with these approaches. Also, BERT Large was found by [Haber and Poesio2021] to allow for better grouping of sense interpretations than other LLMs.[6] We average subwords' embeddings if needed. It is a common practice in previous work on semantic-related tasks to use the average of the last 4 layers to get embeddings; we decided to adopt the same "4 layers average pooling" strategy, but trying with different possible sets of layers (see Appendix C). Therefore, for a set of four layers, we average hidden states across the selected layers to get a single 1024-dimensional vector. We found that layers 14 to 17 obtained the best results on Dev for all methods (global/local-only and bi-level).

**Sense-inducing systems.** Comparison to Local-only systems will give a (strong) baseline just by inducing senses without aiming at concepts. We used the same clustering algorithms. We also implement the WSI method proposed by [Eyal et al.2022]. It relies on a different paradigm, using the Language Model for substitution instead of word embeddings. From lists of substitutes, they build a graph of substitutes in which they find communities and then assign each occurrence to a community of substitutes to find the word senses. Because Local-only methods only induce senses, their hyperparameters are chosen to maximize a WSI objective on polysemous words of the dev split.

**Baselines** We construct a candidate clustering $\widehat{\mathcal{C}}_W$ where each lemma has its own cluster. This baseline model is referred to as the "Lemmas" baseline. This is to evaluate the extent to which the information contained by the lemma alone can be used to induce concepts without any knowledge on word senses neither on context. As a second baseline, we create for each lemma as many singletons as the number of different concepts its occurrences are annotated with. All created clusters are of size 1: we account perfectly for polysemy but not at all for synonymy. This second baseline is dubbed "Oracle WSI".

---

[6]We leave to further work the use of autoregressive and/or newer Language Models.

### 5.3   Concept Induction in SemCor

In Table 1 we display the Concept Induction scores (F1) of proposed baselines and systems on the full SemCor data and on the Synon. split. On the full data, both the Lemmas and Oracle WSI baselines achieve very good performance because they have, by design, a perfect precision (they do not cluster lemmas at all and do not overestimate the number of clusters) and because 88% of concepts are instantiated with only a single lemma (thus their recall is still good). However, they are very limited on the Synon. split of the data, where concepts are instantiated with multiple lemmas.

The proposed Concept Induction systems reach scores ranging from .56 to .66 on the full data, half of them outperforming the Lemmas baseline, and from .59 to .62 on the Synon. split, outperforming all other systems. While still challenging, it exhibits that it is indeed possible to induce WordNet-based concepts in a corpus using LMs hidden layers vectors.

We also see that Kmeans-based approaches are consistently outperformed by Agglomerative methods. This indicates that the representational spaces in LM hidden layers are not organized in a nearly-spherical fashion as Kmeans algorithm assumes, but rather are populated less uniformly. This is reflected in precision and recall: Agglomerative systems reach a higher precision than Kmeans with similar recall.

Overall, results are in favor of Bi-level approaches over Global-only systems, with substantial improvements in F1 on the full data while obtaining (nearly) identical performance on concepts of multiple lemmas, and large increases in precision while the loss in recall is minimal. This demonstrates that considering the local (lemma-centric) perspective is beneficial to a global (cross-lexicon) view when inducing concepts. The local clustering, with the subsequent representation averaging, helps reducing variance in occurrences and therefore allow to reach higher levels of precision in the global clustering compared to Global-only. We would also like to emphasize that, while Global-only systems are more simple in design, their computational cost is usually higher than Bi-level ones, especially when the clustering algorithm's time complexity is quadratic with respect to the number of occurrences.

### 5.4   Qualitative Analysis of Concepts Clusters

We manually annotate word clusters (obtained from our best-performing approach, the Agglo Bi-level system) containing at least 2 lemmas according to the semantic similarity between lemmas. Distribution of cluster sizes (in number of lemmas) can be found in Appendix D. We distinguish four categories: *synonyms*

13

Table 1: Concept Induction BCubed Precision (P), Recall (R) and F1 on the SemCor data averaged over 5 runs.

| | Full data | | | Synon. | | |
|---|---|---|---|---|---|---|
| | P | R | F1 | P | R | F1 |
| **Baselines** | | | | | | |
| Lemmas | 1.0 | .43 | .61 | .31 | .70 | .44 |
| Oracle WSI | 1.0 | .15 | .26 | .35 | .50 | .41 |
| **Local-only Systems** | | | | | | |
| Kmeans Local | .92 | .35 | .50 | .37 | .39 | .38 |
| Agglo Local | .73 | .92 | .81 | .48 | .65 | .56 |
| Eyal et al. (2022) | .39 | .50 | .44 | .53 | .56 | .54 |
| **CI Systems** | | | | | | |
| Kmeans Global | .70 | .59 | .64 | .61 | .60 | .60 |
| Kmeans Bi-level | .75 | .60 | .66 | .68 | .54 | .60 |
| Agglo Global | .82 | .47 | .59 | .82 | .50 | .62 |
| Agglo Bi-level | .86 | .49 | .62 | .86 | .49 | .62 |

when lemmas are cognitive synonyms (e.g. "necessity" and "need"), *near-synonyms* for lemmas close to be synonyms but showing slight difference in meaning (e.g. "duty" and "task", the former being stronger than the latter)[7], *related* when lemmas show a topical (e.g. "dirt", "sand" and "mud") or lexical relations (e.g., antonyms like "man" and "woman") and *invalid* clusters when lemmas show no semantic relation (e.g. "child" and "idea").

Proportions of these annotations are displayed in Table 2 with respect to the cluster size, the number of lemmas in the cluster. For a given cluster size, if the number of clusters exceeds 50, we randomly sample 50 clusters to be annotated. Overall, the proportion of synonyms and near-synonyms is generally above 50% and less than 10% of clusters are invalid, indicating that most learned concepts are reliable and meaningful. We argue that the remaining related term clusters, while not synonyms, may still be interesting in less fine-grained studies. The portion of related clusters is in line with findings from previous work showing that BERT was also reflective of other lexical relations, such as hypernymy [Hanna and Marecek2021].

## 5.5 Benefits at the Local Level

We now turn back to the local level and assess whether the information brought at the global level helps distinguishing senses of individual words. Here we do not evaluate the word-level soft clustering, but

---

[7]Notions of cognitive synonymy and near-synonymy are discussed by [Stanojević2009].

Table 2: Qualitative manual evaluation of obtained word clusters of size $\geq 2$.

| Cluster size | 2 | 3 | 4 | >4 |
|---|---|---|---|---|
| Nb. of annotated clusters | 50 | 50 | 23 | 23 |
| Category (% of annotated clusters) | | | | |
| Synonyms | 42 | 38 | 17 | 24 |
| Near-synonyms | 24 | 24 | 35 | 26 |
| Related | 36 | 48 | 48 | 42 |
| Invalid | 8 | 2 | 0 | 8 |

the occurrence-level division of SemCor's data, considering each word independently. In other words, we evaluate WSI in SemCor using annotations as the reference sense clustering.

**Evaluation of induced senses** For each word $w \in W$, we compare how its set of occurrences $\mathcal{O}_w$ is divided in $\widehat{\mathcal{C}}$ to how it is divided in the reference $\mathcal{C}$ provided by annotations using BCubed metrics, and we average scores obtained across $W$. We display the WSI BCubed F1, as in previous WSI tasks like [Jurgens and Klapaftis2013]. Following [Amrami and Goldberg2019], we report $\rho$ the Spearman correlation coefficient between the number of clusters a lemma is assigned to and its number of senses according to annotations, to ensure that the number of created senses actually scales with the actual degree of polysemy.

Note that, for CI systems, we evaluate the division of occurrences provided by the final clustering $\widehat{\mathcal{C}}$ (i.e. how occurrences are clustered after the global step and its potential merge operations). The quality of sense clusters induced by the local-step only is actually evaluated with Local-only systems.

**Local results.** Results of this local evaluation are displayed in Table 3. Let us recall that Local-only systems' hyperparameters are chosen to maximize the WSI F1 on the dev split, while those of CI systems maximize the Concept Induction F1. Nonetheless, one can observe that all CI systems outperform their Local-only counterparts, achieving higher WSI F1 and $\rho$ even though their hyperparameters are not chosen to match the WSI itself. This indicates that the information brought at the global level by considering cross-lexicon relations may indeed help improving WSI, and benefits between local and global levels go both ways.

We explain the relatively poor performance of State-of-the-art WSI system by the fact that we are in a particular setting, where the number of occurrences per lemma is relatively low in SemCor (30 per lemma on average) and so is the average number of occurrences per concept. Data sparsity is a favorable ground for word senses to be misrepresented. As such, methods meant to be applied on larger datasets like the one of

Table 3: WSI BCubed F1 and sense number correlation coefficient $\rho$ on SemCor full data. Not computed for Kmeans because the number of cluster is constant.

|  | WSI F1 | $\rho$ |
|---|---|---|
| **Local-only Systems** |  |  |
| Kmeans Local | .61 | NA |
| Agglo Local | .76 | .78 |
| Eyal et al. (2022) | .77 | .46 |
| **CI Systems** |  |  |
| Kmeans Global | .42 | NA |
| Kmeans Bi-level | .38 | NA |
| Agglo Global | .17 | .04 |
| Agglo Bi-level | .24 | .51 |

[Eyal et al.2022] may not work as well as expected. Our results show the limitations of these systems when the amount of training data is low and the interest of aiming at concept induction to get senses. This scenario is motivated in areas where data are not available in large quantities and still require to induce senses. In the case of the study of Lexical Semantic Change (the evolution of word meanings over time), recent works perform WSI in diachronic corpora that are often unbalanced and small [Tahmasebi et al.2021].

# 6  Extrinsic Evaluation with Concept-aware Embeddings

In their work, [Eyal et al.2022] derive sense-aware static embeddings from their WSI method, training them on the Wikipedia dataset and used them for the word-in-context (WiC) task. They win nearly-SotA results on the dataset proposed by [Pilehvar and Camacho-Collados2019], and report to be outperformed only by methods using external lexical knowledge and resources. We proceed to the same extrinsic evaluation of our work, constructing concept-aware embeddings using concept clusters of concept induction systems (Global-only and Bi-level Agglo). To obtain such embeddings, we average all vectors representing occurrences in SemCor contained in each global cluster to get one vector per concept cluster.

The WiC task consists of determining whether two occurrences of a target lemma $w$ correspond to the same sense. The WiC dataset's target words are nouns and verbs, but like in the rest of this paper, we restrict our scope to nouns.

To solve the task, we use BERT Large to create representations of the two target occurrences. Each of them is assigned to a concept by finding the closest concept-aware embedding using cosine distance. The

Table 4: Accuracy scores on the nouns of the WiC test dataset [Pilehvar and Camacho-Collados2019].

| Model | Acc. |
|---|---|
| Eyal et al. (2022) (CBOW) | 59.3 |
| Eyal et al. (2022) (Skip-Grams) | 61.9 |
| Ours (Agglo global) | 58.8 |
| Ours (Agglo bi-level) | 60.1 |

decision depends on whether the two occurrences are mapped to the same concept (true) or to distinct ones (false). Results are displayed in Table 4.

Our concept-aware embeddings obtain very similar results to those of their sense-aware embeddings, with ours derived from our bi-level approach even outperforming their CBOW method. Interestingly, our embeddings were trained with far fewer resources than theirs, as we used 52,997 occurrences from the SemCor dataset while they used a dump of Wikipedia, gathering millions of occurrences. This emphasizes the value of concept-aware embeddings: the use of cross-lexicon information allows competitive results with fewer resources.

# 7 Limitations

The formal framework we defined uses terminology and notions from rather structuralist/relational assumptions of the language's lexical system (e.g. senses, discrete concepts, etc.). We made this choice based on how lexical databases like WordNet (and its derivatives), or others like the Historical Thesaurus of English for instance, are designed using the "word/sense/concept" structure. From a purely practical point of view, this choice makes sense as these resources would be the primary source for task data's annotations. Conceptually, senses are also a notion widely used in computational linguistics and we wanted to propose Concept Induction as a step "beyond" this conventional aspect and its related tasks. Future research may explore definitions/extensions of Concept Induction outside of this structuralist/relational framework, towards cognitive semantics for instance [Geeraerts2010].

Evaluating Concept Induction is mainly limited by the low amount of suitable annotated corpora. Not only the data need to be annotated in concepts, but these annotations must cover a wide variety of lemmas for synonymy to be sufficiently represented in the corpus. Future work may find or create datasets meeting these requirements to evaluate Concept Induction outside of SemCor.

For now, the study is limited to nouns. Performances of benchmarked algorithms and systems may change with other Part-of-speech tags.

Our Bi-level method allows the global clustering to merge local clusters, leveraging lexicon-level information to be used to correct Word Sense Induction errors at the lemma-level. By its sequential nature, our method does not allow to split local clusters using global-level information, which could lead to better results. Further research directions include creating an iterative version of our methodology (alternating local and global clustering), or attempting to tackle both clustering objectives simultaneously with bi-level constrained clustering.

Our results about sense-induction at the local level showed that usual WSI methods may not be robust in our setting where there are few occurrences for some lemmas. We demonstrated that, in this setting, concept-inducing methods provided a better division in word senses. In many fields of linguistics, corpora are not very large and do not contain hundreds of occurrences for each word. Nonetheless, it is still uncertain if this observed advantage of CI systems would still hold on bigger datasets with many occurrences per lemma, a setting better-suited for usual WSI methods.

In this paper, we limited our study to Nouns, the morpho-syntactic class exhibiting the most prominent semantic features. We leave to further research the study of Concept Induction for Verbs, Adjectives, or the heterogeneous family of Adverbs.

## 8 Ethical Considerations

Our methodology uses pretrained Contextualized Language Models, which are known to encode and replicate social biases contained in their training data and sometimes amplify them. While we do not observe surface-level biases arising when manually annotating concept clusters, it is still an open question of how these social biases may influence or even change results when inducing concepts in SemCor.

## Acknowledgements

# References

[Amigó et al.2009] Enrique Amigó, Julio Gonzalo, Javier Artiles, and Felisa Verdejo. 2009. A comparison of extrinsic clustering evaluation metrics based on formal constraints. *Information retrieval*, 12:461–486.

[Amrami and Goldberg2019] Asaf Amrami and Yoav Goldberg. 2019. Towards better substitution-based word sense induction.

[Bagga and Baldwin1998] Amit Bagga and Breck Baldwin. 1998. Entity-based cross-document coreferencing using the vector space model. In *36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, Volume 1*, pages 79–85, Montreal, Quebec, Canada. Association for Computational Linguistics.

[Bizzoni et al.2014] Yuri Bizzoni, Federico Boschetti, Harry Diakoff, Riccardo Del Gratta, Monica Monachini, and Gregory Crane. 2014. The making of Ancient Greek WordNet. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC '14)*, pages 1140–1147, Reykjavik, Iceland. European Language Resources Association (ELRA).

[Chronis and Erk2020] Gabriella Chronis and Katrin Erk. 2020. When is a bishop not like a rook? when it's like a rabbi! multi-prototype BERT embeddings for estimating semantic relationships. In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 227–244, Online. Association for Computational Linguistics.

[Devlin et al.2019] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

[Ethayarajh2019] Kawin Ethayarajh. 2019. How contextual are contextualized word representations? Comparing the geometry of BERT, ELMo, and GPT-2 embeddings. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference*

*on Natural Language Processing (EMNLP-IJCNLP)*, pages 55–65, Hong Kong, China. Association for Computational Linguistics.

[Eyal et al.2022] Matan Eyal, Shoval Sadde, Hillel Taub-Tabib, and Yoav Goldberg. 2022. Large scale substitution-based word sense induction. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4738–4752, Dublin, Ireland. Association for Computational Linguistics.

[Fellbaum1998] Christiane Fellbaum. 1998. *WordNet: An electronic lexical database*. MIT press.

[François2022] Alexandre François. 2022. Lexical tectonics: Mapping structural change in patterns of lexification. *Zeitschrift für Sprachwissenschaft*, 41(1):89–123.

[Geeraerts2010] Dirk Geeraerts. 2010. *Theories of Lexical Semantics*. Oxford University Press.

[Ghanem et al.2023] Sana Ghanem, Mustafa Jarrar, Radi Jarrar, and Ibrahim Bounhas. 2023. A benchmark and scoring algorithm for enriching Arabic synonyms. In *Proceedings of the 12th Global Wordnet Conference*, pages 274–283, University of the Basque Country, Donostia-San Sebastian, Basque Country. Global Wordnet Association.

[Haber and Poesio2021] Janosch Haber and Massimo Poesio. 2021. Patterns of polysemy and homonymy in contextualised language models. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2663–2676, Punta Cana, Dominican Republic. Association for Computational Linguistics.

[Haber and Poesio2024] Janosch Haber and Massimo Poesio. 2024. Polysemy—Evidence from linguistics, behavioral science, and contextualized language models. *Computational Linguistics*, 50(1):351–417.

[Hanna and Marecek2021] Michael Hanna and David Marecek. 2021. Analyzing BERT's knowledge of hypernymy via prompting. In *Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 275–282, Punta Cana, Dominican Republic. Association for Computational Linguistics.

[Haspelmath2023] Martin Haspelmath. 2023. Coexpression and synexpression patterns across languages: comparative concepts and possible explanations. *Frontiers in Psychology*, 14.

[Jurgens and Klapaftis2013] David Jurgens and Ioannis Klapaftis. 2013. SemEval-2013 task 13: Word sense induction for graded and non-graded senses. In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*, pages 290–299, Atlanta, Georgia, USA. Association for Computational Linguistics.

[Khan et al.2022] Fahad Khan, Francisco J. Minaya Gómez, Rafael Cruz González, Harry Diakoff, Javier E. Díaz Vera, John P. McCrae, Ciara O'Loughlin, William Michael Short, and Sander Stolk. 2022. Towards the construction of a WordNet for Old English. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 3934–3941, Marseille, France. European Language Resources Association.

[Kutuzov and Giulianelli2020] Andrey Kutuzov and Mario Giulianelli. 2020. UiO-UvA at SemEval-2020 task 1: Contextualised embeddings for lexical semantic change detection. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 126–134, Barcelona (online). International Committee for Computational Linguistics.

[Manandhar et al.2010] Suresh Manandhar, Ioannis Klapaftis, Dmitriy Dligach, and Sameer Pradhan. 2010. SemEval-2010 task 14: Word sense induction & disambiguation. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 63–68, Uppsala, Sweden. Association for Computational Linguistics.

[Martinc et al.2020] Matej Martinc, Syrielle Montariol, Elaine Zosa, and Lidia Pivovarova. 2020. Capturing evolution in word usage: Just add more clusters? In *Companion Proceedings of the Web Conference 2020, WWW'20*, page 343–349, New York, NY, USA. Association for Computing Machinery.

[McDaid et al.2011] Aaron F McDaid, Derek Greene, and Neil Hurley. 2011. Normalized mutual information to evaluate overlapping community finding algorithms. *arXiv preprint arXiv:1110.2515*.

[Miller1995] George A Miller. 1995. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41.

[Nair et al.2020] Sathvik Nair, Mahesh Srinivasan, and Stephan Meylan. 2020. Contextualized word embeddings encode aspects of human-like word sense knowledge.

[Pilehvar and Camacho-Collados2019] Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. WiC: the word-in-context dataset for evaluating context-sensitive meaning representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1267–1273, Minneapolis, Minnesota. Association for Computational Linguistics.

[Raganato et al.2017] Alessandro Raganato, Jose Camacho-Collados, and Roberto Navigli. 2017. Word sense disambiguation: A unified evaluation framework and empirical comparison. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 99–110, Valencia, Spain. Association for Computational Linguistics.

[Saidi and Jarray2023] Rakia Saidi and Fethi Jarray. 2023. Sentence transformers and distilbert for arabic word sense induction.

[Scarlini et al.2020] Bianca Scarlini, Tommaso Pasini, and Roberto Navigli. 2020. Sensembert: Context-enhanced sense embeddings for multilingual word sense disambiguation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8758–8765.

[Stanojević2009] Marija Stanojević. 2009. Cognitive synonymy: A general overview. *Facta Universitatis Series: Linguistics and Literature*, 07:193–200.

[Tahmasebi et al.2021] Nina Tahmasebi, Lars Borin, and Adam Jatowt. 2021. Survey of computational approaches to lexical semantic change detection. *Computational approaches to semantic change*, 6(1).

[Velasco et al.2023] Dan John Velasco, Axel Alba, Trisha Gail Pelagio, Bryce Anthony Ramirez, Jan Christian Blaise Cruz, Unisse Chua, Briane Paul Samson, and Charibeth Cheng. 2023. Towards automatic construction of Filipino WordNet: Word sense induction and synset induction using sentence embeddings. In *Proceedings of the First Workshop in South East Asian Language Processing*, pages 1–12, Nusa Dua, Bali, Indonesia. Association for Computational Linguistics.

[Zhang et al.2021] Jingqing Zhang, Luis Bolanos Trujillo, Tong Li, Ashwani Tanwar, Guilherme Freire, Xian Yang, Julia Ive, Vibhor Gupta, and Yike Guo. 2021. Self-supervised detection of contextual synonyms in a multi-class setting: Phenotype annotation use case. In *Proceedings of the 2021 Con-*

*ference on Empirical Methods in Natural Language Processing*, pages 8754–8769, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

# A  Extended BCubed to Evaluate CI and WSI

The extension of BCubed for overlapping clusters rely on two quantities, Multiplicity Precision (MP) and Multiplicity Recall (MR). In the case of Concept Induction, MP and MR between two lemmas are defined as follows:

$$MP(w_1, w_2) = \frac{\min[|g(w_1) \cap g(w_2)|, |f(w_1) \cap f(w_2)|]}{|f(w_1) \cap f(w_2)|}$$

$$MR(w_1, w_2) = \frac{\min[|g(w_1) \cap g(w_2)|, |f(w_1) \cap f(w_2)|]}{|g(w_1) \cap g(w_2)|}$$

with $w_1$ and $w_2$ two lemmas, and $g$ a reference clustering function and $f$ the clustering function we want to evaluate. MP (resp. MR) can be computed for every lemma $w_1$ with every other lemma $w_2$ sharing at least one cluster with $w_1$ in $f$ (resp. in $g$). We denote $MP(w_1, \cdot)$ and $MR(w_1, \cdot)$ the obtained averages. In the case of non-overlapping clusters, this formulation gives the same result as the original (non-extended) BCubed. To evaluate WSI, the formulation is the same but we do not evaluate at the word-level but at the occurrence-level.

Precision, Recall and F-score are obtained as follows:

$$\text{Precision} = \frac{1}{|W|} \sum_{w \in W} MP(w, \cdot)$$

$$\text{Recall} = \frac{1}{|W|} \sum_{w \in W} MR(w, \cdot)$$

$$F_\beta = \frac{(1 + \beta^2) \times \text{Recall} \times \text{Precision}}{\beta^2 \times \text{Precision} + \text{Recall}}$$

By default we fix $\beta = 1$, as we compare the learned clustering and the reference clustering as equals and therefore do not find that Precision and Recall should be weighted differently.

[Amigó et al.2009] showed that the benefits of BCubed over other clustering scores. For instance,

Rand Index does not handle well the case of many small clusters, which is likely to be the case for Concept Induction. We also prefer Extended BCubed over Overlapping Normalized Mutual Information [McDaid et al.2011] as the latter is matching-based. That is, the repetition (or non-repetition) of identical clusters will have no impact on the measure. However, we can easily imagine identical clusters of words to be repeated as they may correspond to distinct concepts. In Extended BCubed, repeated clusters are taken in account as we measure the number of times two lemmas are clustered together. The denominator of MP ensures that over-estimating the number of common clusters is also penalized, and those of MR ensures that under-estimating is penalized. Min operators are there to prevent both quantities to grow over 1.

# B    Splits and dataset statistics

In Table 5 we display statistics over the different splits we used. Dev is a subset containing a sample of 10% of concepts and their occurrences. Synon. is a subset containing only concepts instantiated with 2 lemmas or more, and their occurrences.

Table 5: Statistics on the different data splits in annotated SemCor. The split "synon" only contains occurrences of concepts instantiated with multiple lemmas (cases of synonymy). $d_{\mathrm{lemmas}}$ is the average number of unique lemmas per concept, $d_{\mathrm{polysemy}}$ is the average number of distinct concepts per lemma.

| | #Occs | #Lemmas | #Concepts | #Occs/Concept | #Occs/Lemma | $d_{\mathrm{lemmas}}$ |
|---|---|---|---|---|---|---|
| $d_{\mathrm{polysemy}}$ | | | | | | |
| Full data 2.47 | 52,997 | 1,560 | 3,855 | 13.75 | 33.97 | 1.14 |
| Dev 2.24 | 4,795 | 389 | 386 | 12.42 | 12.33 | 1.14 |
| Synon. 1.59 | 13,158 | 630 | 447 | 29.44 | 20.89 | 2.83 |

# C    Used hyperparameters and layers

## C.1    CLM layers

Prior work like [Ethayarajh2019] showed that later layers usually correlates with deeper levels of contextualization and more semantic information, [Chronis and Erk2020] showed that moderately late were preferred for lexical similarity while very last layers were preferred for semantic relatedness. To get embeddings, we

try 4 sets of layers corresponding to different depths: first layers (1 to 4), moderately early layers (8 to 11), moderately late (14 to 17), and last layers (21 to 24). To get the representation of a word's occurrence, we simply average its embeddings from the four chosen layers into one single 1024-dimensional embedding. For Concept Induction, we find that best results were obtained using layers 14 to 17, that are the reported results.

## C.2  Hyperparameters

For [Eyal et al.2022], we tried different resolution, varying it from $10^{-3}$ to 10, for the Louvain clustering but found very little to no effect.

For Kmeans at the local level, we varied the number of clusters $k$ between 2 and 10. For Agglomerative clustering at both levels, we tried single, average and complete linkage.

The distance threshold in Agglo $\tau$ was indexed on the distribution of distances. We fixed an hyperparameter $\gamma$ and derived $\tau = \text{avg}(d) + \gamma \cdot \text{std}(d)$ with $d$ the distribution of distances between clustered instances. We made $\gamma$ vary between $-4$ and $+8$. For global Kmeans, the number of clusters was indexed using a proportion $\pi$ on the number of lemmas (e.g. $120\% \times |W|$), $\pi$ varying from 40% to 400%. This may help transferring hyperparameters to other dataset in future research.

Best hyperparameters choices are in Table 6.

Table 6: Best hyperparameters on the Dev split

| Systems | Best hyperparameters |
|---|---|
| Local-only Kmeans | $k = 3$ |
| Local-only Agglo | linkage=average, $\gamma = 1.0$ |
| Global-only Kmeans | $k = 8, \pi = 120\%$ |
| Global-only Agglo | linkage=average, $\gamma = 3.5$ |
| Bi-level Kmeans | $k = 10$, linkage=average, $\gamma_{\text{local}} = 0.0, \pi = 120\%$ |
| Bi-level Agglo | linkage=average (both), $\gamma_{\text{local}} = 0.0, \gamma_{\text{global}} = 4.0$ |
| Bi-level Kmeans (local Agglo) | $k = 10$, linkage=average, $\gamma_{\text{local}} = 4.0$ |
| Bi-level Agglo (local Kmeans) | linkage=average, $k = 3$ |

# D  Concept Clusters Size Distribution

The distribution of the concept cluster size (in number of lemmas) obtained with Bi-level Agglo system can be found in Figure 2.
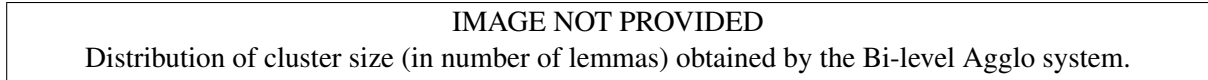
> IMAGE NOT PROVIDED
> Distribution of cluster size (in number of lemmas) obtained by the Bi-level Agglo system.

Figure 2: Distribution of cluster size (in number of lemmas) obtained by the Bi-level Agglo system.

# E  Scientific Artifacts

We used WordNet and SemCor, both properties of Princeton University. Licence can be found at `https://wordnet.princeton.edu/license-and-commercial-use`.