

Word sense generation

Anonymous ACL submission

Abstract

The lexicon makes creative reuse of words to express novel senses. A long-standing effort in natural language processing has been focusing on disambiguating and inducing word senses from context. Little has been explored about how novel word senses may be generated automatically. We consider a paradigm of *word sense generation* (WSG) that enables words to spawn new senses by extending toward novel naturalistic context. We develop a general framework that simulates novel word sense extension by dividing a word into hypothetical child tokens and making inferences about the plausibility of sense extension among the sibling tokens in usage sentences that never appear in training. Our framework combines probabilistic models of chaining with a learning scheme that transforms a language model embedding space to support various types of word sense extensions. We evaluate our framework rigorously against several competitive baselines and show that it is superior in predicting plausible novel senses including metonymic and metaphoric word usages in a large set of 1,500 English verbs. We show that the learned semantic space exhibits systematic patterns of word sense extension while retaining competence in common natural language processing tasks.

1 Introduction

A key property of the lexicon is the creative reuse of words to express novel senses. For example, the English phrase *to arrive at* extended from its original sense “to reach a physical location (e.g., gate)” toward new senses such as “to come to an event (e.g., concert)” and “to reach an abstract, cognitive state (e.g., conclusion)”. The extension of word meaning toward new context may appear to draw on different processes ranging from metonymy to metaphor (see Figure 1 for an illustration), but here we present a general framework that infers how words extend to plausible new senses under novel naturalistic context.

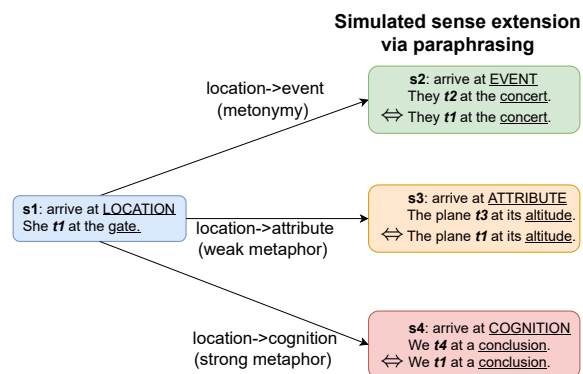


Figure 1: The word sense generation (WSG) framework. A verb (e.g., *arrive*) with a set of sense-labeled usage sentences is partitioned into distinct child tokens (e.g., t_1 - t_4 signifying senses s_1 - s_4 as illustrated). Each token represents a hypothetical word type that replaces its original parent verb in sentences where it expresses a given sense (e.g., *arrive* in sentences where it predicates a location will be substituted by the hypothetical t_1 , and separately substituted by other tokens t_2 - t_4 for sentences expressing the other senses). WSG models infer whether a child token (e.g., t_1) can be extended to express the senses of its siblings via paraphrasing novel sentences that do not appear in training (e.g., t_1 has 3 potential siblings to extend its meaning to: t_2 , t_3 , t_4).

One of the most long-standing efforts in natural language processing (NLP) is to be able to disambiguate word senses from text. This line of work takes a discriminative approach toward the multifaceted aspect of word meaning and has developed models relying on both traditional machine learning techniques (Gale et al., 1992; Kilgarriff and Rosenzweig, 2000; Zhong and Ng, 2010; Iacobacci et al., 2016; Raganato et al., 2017) and modern neural language models (Huang et al., 2019; Wiedemann et al., 2019; Loureiro and Jorge, 2019; Bevilacqua and Navigli, 2020). Related work has developed automated methods for inducing or detecting novel word senses (Lau et al., 2012; Cook et al., 2014; Lau et al., 2014), which can also be considered as unsupervised approaches to sense

059 disambiguation. Here we consider an alternative
060 paradigm that aims to model word senses by taking
061 a generative approach in naturalistic context.

062 Work in computational and cognitive linguistics
063 suggests that word senses often do not extend
064 arbitrarily (Nunberg, 1979; Lehrer, 1990). Lexical
065 semanticists have pointed out that a number of
066 cognitive devices may be applied to generate creative
067 word uses, such as logical metonymy (Cope-
068 stake and Briscoe, 1995; Pustejovsky, 1998) and
069 metaphor (Lakoff and Johnson, 2008; Pustejovsky
070 and Rumshisky, 2010). Cognitive linguists have
071 also suggested that systematic mappings between
072 conceptual domains underlie the metaphorization
073 of word meanings (Brugman and Lakoff, 1988;
074 Lakoff and Johnson, 2008). However, the re-
075 liance on hand-crafted rules of semantic productiv-
076 ity makes it difficult to implement these accounts
077 into computational systems that support flexible
078 and scalable generation of novel word senses.

079 We develop a principled framework termed *word*
080 *sense generation* (WSG). As a starting point, we fo-
081 cus on modelling sense generation of English verbs,
082 which constitute a broad yet notoriously challeng-
083 ing class of productive sense extensions (Puste-
084 jovsky and Rumshisky, 2010). Figure 1 illustrates
085 our WSG framework. Given a verb (e.g., *arrive*),
086 we consider a paradigm that simulates how it can be
087 extended to express plausible novel senses in natu-
088 ralist context. We do so by dividing a word into
089 hypothetical child tokens signifying its different
090 senses. We then infer whether a child token may be
091 extended to express its sibling tokens under novel
092 context (i.e., a simulation of how a word might ex-
093 tend from its existing senses to novel senses). We
094 propose a family of deep probabilistic models for
095 this inference problem that are built on the cog-
096 nitive theory of chaining, which states that word
097 meanings grow by linking novel senses to existing
098 ones that are close in semantic space (Lakoff, 1987;
099 Malt et al., 1999; Ramiro et al., 2018; Habibi et al.,
100 2020; Yu and Xu, 2021). We expect these chaining
101 models to support the incremental extension of a
102 word’s meaning toward a variety of new senses,
103 a process analogous to the gradient cline of verb
104 sense extensions (e.g., weak metaphor) discussed
105 in Pustejovsky and Rumshisky (2010).

106 We make three contributions: 1) we formulate
107 word sense generation as a novel probabilistic in-
108 ference task whereby a language model, after learn-
109 ing a set of partitioned tokens signifying different

senses of a polysemous word, automatically infers
whether sibling tokens can be used interchangeably
under novel context absent in training; 2) we de-
velop a family of WSG models motivated by the
cognitive theories and models of semantic chaining,
and a new learning scheme to capture regular pat-
terns of word sense extension; 3) we collect a new
dataset of word sense generation examples which
includes natural usages for approximately 22,000
senses of over 1,500 common English verbs.¹

2 Related work

2.1 Theories of word meaning extension

110 Researchers in lexical semantics and cognitive lin-
111 guistics have both proposed theories to account
112 for the malleable nature of lexical meaning. The
113 Generative Lexicon theory by Pustejovsky (1998)
114 argues that a fixed set of generative devices, such
115 as type-coercion and co-composition, can operate
116 on the lexical structure a word to produce various
117 related meaning interpretations. Copestake and
118 Briscoe (1995) also illustrates how formal lexical
119 rules such as grinding and portioning can be ap-
120 plied to produce novel word usages such as logical
121 metonymy. In cognitive linguistics, Lakoff and
122 Johnson (2008) argues that systematic mappings
123 across conceptual domains result in the abundance
124 of metaphorical word senses in natural language,
125 while these mappings can be motivated by cog-
126 nitive processes such as chaining (Lakoff, 1987)
127 and image schema transformation (Brugman and
128 Lakoff, 1988; Dewell, 1994; Gibbs Jr and Colston,
129 2008). Our work connects the cognitive and formal
130 approaches to word meaning extension by show-
131 ing that a cognitively inspired chaining-based word
132 sense generation framework can learn systematic
133 patterns of meaning extension discussed in the tra-
134 dition of generative lexical semantics.

2.2 Non-literal language generation

147 Our framework also relates to research on auto-
148 mated generation of non-literal word uses such as
149 metaphor and metonymy. Recent work has ex-
150 plored using contextualized language models to
151 generate metaphorical paraphrases for literal usage
152 sentences of a word (Tong et al., 2021; Stowe et al.,
153 2021; Chakrabarty et al., 2021). It has also been
154 shown that chaining mechanisms can be incorpo-
155 rated into contextualized language models to pre-

¹We release the code and data for our work here:
PLACEHOLDER link.

dict unconventional word usages such as slang (Sun et al., 2021). On the other hand, generation of more conventionalized senses, such as metonymy, have remained underexplored (Rambelli et al., 2020), because most contextualized language models have already been exposed to usages of the most common senses of a word type during pretraining, while evaluation on sense generation requires models without such prior knowledge. Our framework circumvents this circularity problem by creating novel tokens that reflect partial usages of a polysemous word and utilizing language models to learn them from scratch, thereby allowing us to model the generative processes of word senses in a zero-shot setting.

3 Framework of word sense generation

Our framework of word sense generation involves three interrelated components: 1) A procedure for partitioning polysemous words in the lexicon into child tokens signifying their different senses; 2) a probabilistic formulation for inferring a child token to paraphrase one of its siblings under a novel linguistic context; and 3) a representational learning algorithm for a transformed semantic space to learn flexible extensions of word senses.

3.1 Sense-based word type partitioning

Let $V = \{w_1, \dots, w_{|V|}\}$ be our vocabulary of polysemous English verbs, where each verb w_i has a sense inventory $S(w_i) = \{s_i^{(1)}, \dots, s_i^{(n_i)}\}$. Assume that for each sense s_i^j of w_i , there is a collection of its representative usage sentences $U(s_i^j)$ in which w_i exhibits sense s_i^j . We wish to investigate whether a language model, which has knowledge only about a partial set of senses of w_i , is able to generate a usage that reflects a novel sense of w_i . In particular, we define a *partition* of a word type w as a grouping of its sense inventory $S(w_i)$ into a collection of K_i distinct sense subsets $\{S_1^i, \dots, S_{K_i}^i\}$ – for example, as illustrated in Figure 1, the sense inventory of the verb *arrive* $S = \{s_1, s_2, s_3, s_4\}$ can be partitioned into four singleton sets (represented by the four colored rectangles). For each sense subset S_k^i with an associated usage sentence set $U(S_k^i) = \bigcup_{s \in S_k^i} U(s)$, we replace all mentions of w_i in each $u \in U(S_k^i)$ with a novel child token t_k^i (e.g., the sentence “We *arrive* at a conclusion” in which *arrive* expresses the abstract state achievement sense s_4 will be converted into “We t_4 at a conclusion”). We then use a contextualized lan-

guage model to learn semantic representations for each child token from the replaced usage sentences $U(S_k^i)$ via the task of masked language modeling (MLM). To prevent information smuggling, the language model is initialized from scratch and therefore does not have any *a priori* knowledge about either the partitioned tokens or their parent verb types. Next, we explain how the task of WSG can be formulated as a paraphrasing problem of inferring a partitioned child token to substitute one of its siblings under a given context.

3.2 Probabilistic formulation of WSG

3.2.1 WSG as partitioned token paraphrasing

Let S_k, S_l be two partitioned sense subsets of the sense inventory $S(w)$ of w , and t_k, t_l be their corresponding child tokens. We say that a language model generates a novel sense for token t_k (called the *source* token) if infers that t_k can serve as a good paraphrase token to substitute its sibling t_l (called the *target* token) in a usage sentence $u^* \in U(S_l)$ containing t_l that does not appear in the MLM training set. For instance, if the LM initially learns two child tokens spawned from the the verb *arrive* that reflect its two distinct senses $s_1 =$ “to come to a physical location” and $s_2 =$ “to achieve a goal” respectively, we would expect the LM to predict that the source token t_1 denoting the concrete sense s_1 can be used to paraphrase its target sibling t_2 denoting the abstract sense s_2 in usages such as “They t_2 at a conclusion after a debate”. We cast WSG as inference of the following word choice probability:

$$P(t_k \rightarrow S_l; u^*) = P(t_k \Leftrightarrow t_l | u^*) \quad (1)$$

Here $t_k \rightarrow S_l$ means that t_k can be extended to express novel senses drawn from S_l , and $t_k \Leftrightarrow t_l$ means that t_k, t_l can be used interchangeably under context u^* . Next, we introduce several models that infer the paraphrase probability in Eq. 1.

3.2.2 Baseline models of WSG

We first consider two simple baseline models: the masked language modeling (MLM) baseline ignores information about t_l and predicts $P(t_k \Leftrightarrow t_l | u^*)$ simply as the infilling probability of t_k under a masked sequence of u^* :

$$P(t_k \Leftrightarrow t_l | u^*) = P_{MLM}(t_k | u_{\setminus t_l}^*) \quad (2)$$

Here $P_{MLM}(t_k|u^*_{t_l})$ denotes the probability of choosing t_k to infill a masked sequence of u^* with t_l replaced by a placeholder, as determined by the contextualized language model. The semantic textual similarity (STS) baseline instead predicts the paraphrase probability as proportional to the cosine similarity between the contextualized representations $h(t_l|u^*), h(t_k|u^*)$ of t_l and t_k under the context u^* :

$$P(t_k \Leftrightarrow t_l|u^*) \propto \text{cosine-sim}(h(t_l|u^*), h(t_k|u^*)) \quad (3)$$

3.2.3 Chaining-based models of WSG

A common issue of the two baseline models described is that they do not model the relations or semantic similarities between the novel use of t_k in u^* and its existing usages $U(S_k)$. We therefore propose a family of WSG models that draws inspirations from the cognitive theory of chaining and memory-augmented deep learning. These models predict that a token t_k is a good paraphrase for its sibling t_l under u^* if the collection of conventional usages $U(S_k)$ of t_k bears a close overall proximity to u^* in the contextualized embedding space:

$$P(t_k \Leftrightarrow t_l|u^*) \propto \text{sim}(t_k, t_l) \quad (4)$$

$$= \text{sim}(H(t_k), h(t_l; u^*)) \quad (5)$$

Here $H(t_k) = \{h(t_k; u)\}_{u \in U(S_k)}$ is the collection of contextualized embeddings of t_k in its conventional usages. We next describe two commonly used chaining models that specify the similarity function $\text{sim}(H(t_k), h(t_l; u^*))$.

WSG-Prototype model. The prototype model draws inspirations from prototypical network for few-shot learning (Snell et al., 2017) and follows the prototype theory of categorization (Rosch, 1975) in cognitive psychology. It assumes that the meaning of a child token t_k can be summarized by a global mean of its contextualized embeddings taken from all of its conventional usages $U(S_k^i)$, so that the probability of t_k being a good paraphrase for t_l under u^* is proportional to the semantic similarity between the contextualized embedding $h(t_l|u^*)$ of t_l and the summary prototype of its sibling:

$$\text{sim}(t_k, t_l) = \exp(-\|h(t_l|u^*) - z(t_k)\|^2) \quad (6)$$

$$z(t_k) = \frac{1}{|U(S_k^i)|} \sum_{u \in U(S_k^i)} h(t_k|u) \quad (7)$$

Here $z(t_k)$ is the mean contextualized embedding of t_l over all of its existing usages, and we have defined semantic similarity as the negative exponential Euclidean distance between two embeddings.

WSG-Exemplar model. The exemplar model resembles the memory-augmented matching network in deep few-shot learning (Vinyals et al., 2016), and formalizes the exemplar theory of categorization (Nosofsky, 1986). This model postulates that the meaning of t_k is represented by the entire collection of its usages $u \in U(S_k^i)$. The probability that t_k paraphrases t_l under context u^* is then proportional to the mean negative exponential Euclidean distance between $h(t_l|u^*)$ and each contextualized embedding of t_k :

$$\text{sim}(t_k, t_l) = \frac{1}{|U(S_k^i)|} \sum_{u \in U(S_k^i)} \exp(-\|h(t_l|u^*) - h(t_k|u^*)\|^2) \quad (8)$$

3.3 Learning sense-extensional semantic space

Chaining relies on identifying close semantic relations between senses, and we therefore develop a learning scheme that transforms a standard semantic space to one that is sensitive to regular relations attested in sense extension. For instance, if a WSG model has observed how verb *grasp* relates its literal sense (e.g., to *grasp* an item) to the extended metaphorical sense (e.g., to *grasp* an idea), under the transformed semantic space the model should also predict similar but novel non-literal sense extensions for other verbs that involve such metaphorical mappings (e.g., to *get* someone’s idea, which also reflects the conceptual metaphor IDEAS ARE OBJECTS).

We follow work in deep few-shot learning and propose an episodic learning algorithm for a sense-extensional semantic space: at each episode, we sample a mini-batch of N source-target token pairs $\{(t_k^{(n)}, t_l^{(n)})\}_{n=1}^N$ partitioned from N different parent word types, and sample a usage sentence $u^{*(n)}$ for each target token $t_l^{(n)}$. The model then learns to perform in-batch WSG by choosing a paraphrase token for every target $t_l^{(n)}$ in $u^{*(n)}$ among the set of N candidate source tokens $\{(t_k^{(1)}, \dots, t_k^{(N)})\}$, with $t_k^{(n)}$ being the ground-truth paraphrase. For the two chaining-based models, learning can be performed by minimizing the in-batch classification loss:

$$\mathcal{J} = \sum_{n=1}^N -\log \frac{\text{sim}(H(t_k^{(n)}), h(t_l^{(n)}; u^{*(n)}))}{\sum_{n'} \text{sim}(H(t_k^{(n')}), h(t_l^{(n)}; u^{*(n)}))} \quad (9)$$

Sense type	No. of verb types	No. of spawned tokens (Full partition)	No. of spawned tokens (Leave-one-out partition)	No. of usage sentences
Domain-based	1,199	14,735	5,860	481,654
Synset-based	1,468	6,463	2,936	371,920

Table 1: Summary statistics of the two collected WSG datasets.

Here $\text{sim}(\cdot, \cdot)$ can either be a prototype-based similarity function in Eq.6, or be its exemplar-based counterpart specified in Eq.8. The two baseline models can also be trained directly on the same set of examples by maximizing their word choice probabilities (i.e., Eq.2 and Eq.3).

4 Data

We collect usage sentences for English verbs from the Wikitext-103 linguistic corpus (Merity et al., 2016) that is commonly used as a language modeling benchmark dataset. It is well-known that word senses are highly fuzzy linguistic categories, and there does not exist a set of word senses that is suitable for every NLP task (Kilgarriff, 1997; Rumshisky and Batiukova, 2008). We therefore consider two different definitions of verb senses in our datasets, as well as two different ways of partitioning the sense inventory of a verb.

4.1 Domain-based vs. synset-based verb sense

Our first sense definition of an English verb is based on the semantic domain of its syntactic object nouns. In particular, for each usage of a verb with a noun object, we label the noun with its supersense category defined in the WordNet lexical database (Miller, 1995), and define a verb sense as the collection of usages whose object nouns share the same supersense label. This sense definition helps us investigate the regularity underlying sense extension cases of different words. For instance, a common type of logical polysemy involves coercing an event nominal to denote the physical location where it takes place (e.g. to arrive at the *theatre* \rightarrow to arrive at the *concert*), which can be captured by extending the domain-based sense “to arrive at a LOCATION” to “to arrive at an EVENT”². We also consider a second, more established type of word sense as recorded in lexicographic resources. In particular, for each word, we apply a state-of-the-art word sense disambiguation algorithm (Bevilacqua and Navigli, 2020) on each of its usage sentence to

²See Appendix B for a full list of supersense categories and their descriptions.

identify the its evoked WordNet sense, as indicated by one of its associated synset ID.

4.2 Full vs. leave-one-out word type partition

For each WSG dataset, we also consider two types of sense inventory partitioning: the *full partition* creates a child token for each sense of a parent verb (e.g., the partition in Fig.1), and a WSG model is evaluated on predicting extensions between all possible child token pairs spawned from the same word type. The *leave-one-out partition* instead randomly samples one sense for each verb to create the target token, and group all remaining senses into a single source token. Summary statistics of the two resulting datasets are shown in Table 1.

5 Evaluation and results

5.1 Model implementation

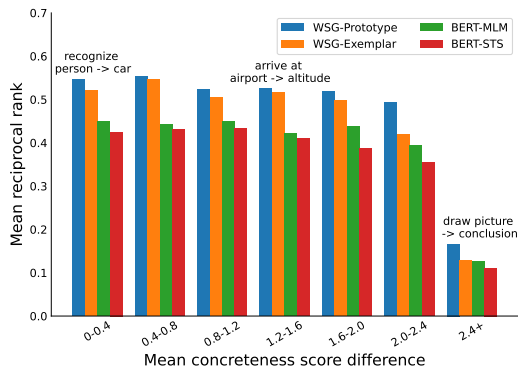
We use a BERT language model (Devlin et al., 2019) to build both the contextualized baselines and the chaining-based WSG models. All BERT models are implemented from scratch and are not pretrained on any NLP tasks. In the masked language modeling step, we increase the vocabulary size of each model by replacing all parent verbs in our datasets with their spawned child tokens, and increase the size of the model’s embedding layer and final classification layer accordingly. During sense-extensional semantic space learning, we randomly choose 70% of the parent verb types in each dataset, and take usage sentences containing their spawned tokens as training set. Sentences containing partitioned tokens spawned by the remaining 30% verb types will be taken as the test set, so that there is no overlap in the vocabulary of partitioned tokens or parent verb types between training and testing.³

5.2 Evaluation on WSG

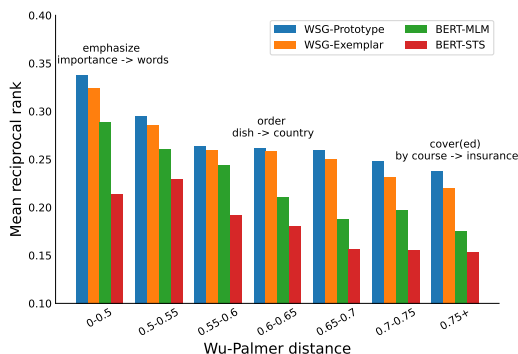
We first evaluated our models on the task of predicting paraphrase partitioned tokens formulated in Eq.1: given a target token and its sample usage

³Implementation details are described in Appendix A.

417 sentence, how likely is the model to predict one
 418 of its sibling source tokens as a good paraphrase?
 419 At each trial, we pick one ground-truth source to-
 420 ken and 99 negative tokens with different parent
 421 verb types, and ask the model to rank the 100 can-
 422 didates based on their infilling likelihoods. We
 423 then compute the mean reciprocal ranks for each
 424 ground-truth source token among the 100 candi-
 425 dates (MRR-100) to assess model performance.



(a) Model performance (MRR) vs. Concreteness score difference between supersenses.



(b) Model performance (MRR) vs. Wu-Palmer distance between WordNet synsets.

Figure 2: Model performance vs. sense relatedness between fully-partitioned tokens.

426 Table 2 summarizes the averaged results over
 427 five independently sampled candidate sets. The un-
 428 supervised/supervised columns correspond to mod-
 429 els with/without learning a sense-extensional se-
 430 mantic space. We observed that 1) all BERT-based
 431 models benefit from learning a sense-extensional
 432 semantic space, suggesting the presence of regu-
 433 larity shared among examples of sense extension
 434 across verb types; 2) both the prototype and ex-
 435 emplar WSG models consistently supersede other
 436 baseline models in both unsupervised and super-
 437 vised setups, indicating that chaining mechanisms
 438 are useful inductive biases for modeling the gener-

439 ative processes of word meaning extension.⁴

440 Table 3 shows example predictions on 6 poly-
 441 semous verbs made by a supervised prototype
 442 WSG model (the full model) and a supervised
 443 BERT-MLM baseline trained on the same datasets.
 444 The full model successfully predicts many types
 445 of sense extension, including weak metaphor (the
 446 *pass* example, where the location-type argument of
 447 a predicate is weakened to its super-type of scalar
 448 attribute), strong metaphor (the *throw* and *stretch*
 449 examples, where verbs with concrete noun argu-
 450 ments are extended to predicate abstract terms)
 451 and logical metonymy (the *appear* (in) example
 452 of event-for-place type coercion). In contrast, the
 453 MLM baseline exhibits a greater tendency to pre-
 454 dict a “literal” paraphrase for a partitioned token
 455 (e.g., all of its top-3 predicted paraphrase tokens for
 456 “throw COGNITION” have senses from the same
 457 abstract domain). In addition, both chaining-based
 458 and baseline models still struggle in predicting
 459 some usages that involve strong non-literal sense
 460 extension (e.g., the *grasp* example).

461 6 Model interpretation and analysis

462 6.1 Sense relatedness and model predictability

463 Prior work in psycholinguistics suggests that both
 464 adults and children often find it easier to infer a
 465 new intended meaning of a word if they can access
 466 a highly related conventional sense of that word
 467 to constrain their interpretation (Clark and Ger-
 468 rig, 1983; Klepousniotou et al., 2008; Rodd et al.,
 469 2012). Here we investigate whether our WSG mod-
 470 els exhibit human-like sensitivity to the conceptual
 471 relatedness of source-target sense pairs in the full
 472 partition setup. We quantify the degree of con-
 473 ceptual relatedness for domain-based partitioned
 474 tokens by computing the difference in mean con-
 475 creteness score (Brysbaert et al., 2014) of their
 476 attested object nouns, and for tokens representing
 477 synset-based senses, we take the Wu-Palmer se-
 478 mantic distance (Wu and Palmer, 1994) between
 479 their associated WordNet synsets as the similarity
 480 measurement. Figure 2 shows the performance of 4
 481 WSG models by binning sense pairs based on their
 482 degrees of conceptual similarity. We observe that
 483 all models yield worse prediction for sense pairs
 484 that are conceptually more distant (e.g., metaphors),
 485 while generally performing better on pairs that are

⁴We also summarize results on a subset of the domain-based dataset where the target sense is more abstract than the source – see Appendix C.

Sense type	Model	Mean reciprocal rank (MRR-100)			
		Full partition		Leave-one-out partition	
		Unsupervised	Supervised	Unsupervised	Supervised
Domain-based	Random Baseline	5.2±0.0	–	5.2±0.0	–
	BERT-STS	13.7±1.2	33.1±2.4	38.5±1.9	62.5±1.6
	BERT-MLM	15.2±1.5	35.4±1.9	44.4±2.0	64.8±1.3
	WSG-Prototype	19.8±1.3	52.3±0.3	55.0±1.4	86.9±1.6
	WSG-Exemplar	19.5±1.1	47.9±1.5	78.2±2.7	86.6±1.1
Synset-based	Random Baseline	5.2±0.0	–	5.2±0.0	–
	BERT-STS	14.2±2.0	16.9±1.2	26.2±2.0	55.0±1.8
	BERT-MLM	17.3±0.8	19.1±0.7	25.7±0.8	59.9±0.5
	WSG-Prototype	20.7±0.5	30.8±0.6	36.0±1.1	73.4±1.0
	WSG-Exemplar	19.6±0.8	29.5±1.2	68.0±2.9	72.9±0.6

Table 2: Summary of model MRR-100 scores (%) for word sense generation in the two datasets. Numbers after \pm are standard deviations over 5 sets of independently sampled negative candidate tokens.

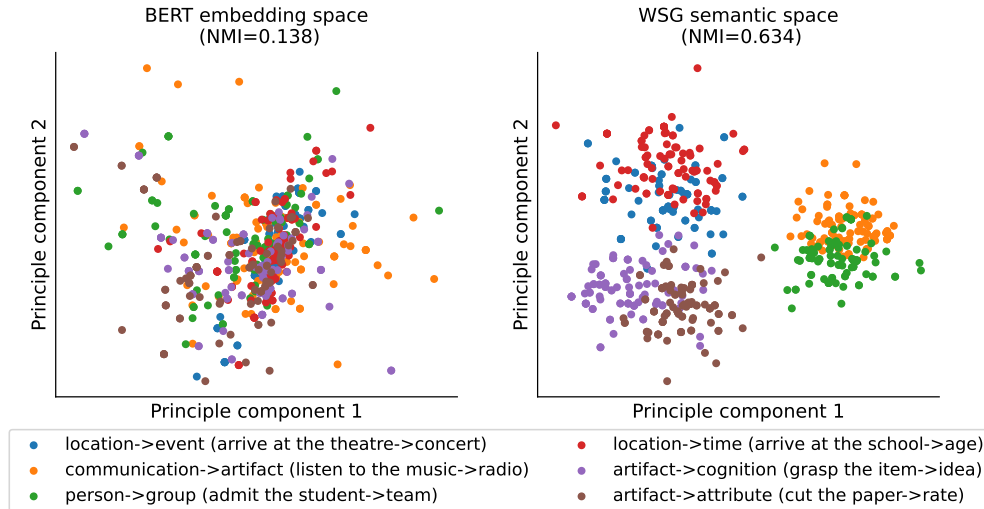


Figure 3: Principal component visualization of prototype difference embeddings of partitioned source-target token pairs, under BERT-MLM baseline (left) and prototype WSG model in sense-extensional semantic space (right).

conceptually more related (e.g., metonymy).

6.2 Interpreting the learned semantic space

To better understand the effect of the sense-extensional semantic space, we compute the mean Euclidean distance between the prototypes (i.e., global mean) of embedded usage sentences for all fully-partitioned source-target token pairs yielded by two types of chaining-based models. As shown in Table 4, both models benefit from learning the extensional semantic space by bringing closer novel and existing senses of the same word. Moreover, pushing partitioned tokens closer for verbs in the training set also results in a more compact embedded sense inventory for unseen verbs in the evaluation set, suggesting that the WSG models have captured some regularity shared across the meaning

transformations of various word types.

To further interpret the information captured in the WSG models, we also performed unsupervised K-means clustering on prototype difference vectors of all fully partitioned domain-based source-target token pairs taken from 1) the embedding space of the BERT-MLM baseline and a sense-extensional semantic space learned by a prototype WSG model. We then compute the normalized mutual information (NMI) between the cluster labels yielded by the K-means algorithm and the ground-truth pairings of source-target domains for each example. Figure 3 shows the clustering results for two semantic spaces, together with their NMI scores against the ground-truth sense extension type labels. We observe that the learned space captures systematic sense extensional patterns by

Model	Top-3 source tokens predicted by model	Predicted rank
1.1 Verb: <i>throw</i> ; target sense: to throw a COGNITION ; true source sense: to throw an ARTIFACT Usage context: Constructors ... which returned null upon failure were changed to <i>throw</i> an exception instead.		
BERT-MLM	create COGNITION, develop COGNITION, provide COGNITION	55/100
WSG-Prototype	create COGNITION, throw ARTIFACT, build ARTIFACT	2/100
1.2 Verb: <i>pass</i> ; target sense: to pass an ATTRIBUTE ; true source sense: to pass an ARTIFACT Usage context: After 22 minutes of flight, the aircraft <i>passed</i> its assigned altitude.		
BERT-MLM	retain ATTRIBUTE, lose ATTRIBUTE, regain ATTRIBUTE	41/100
WSG-Prototype	pass ARTIFACT, enter LOCATION, return to ARTIFACT	1/100
1.3 Verb: <i>appear (in)</i> ; target sense: to appear in an EVENT ; true source sense: to appear in a LOCATION Usage context: He ... <i>appeared in</i> 24 league matches as well as United's FA Cup defeat to Burnley.		
BERT-MLM	host EVENT, achieve EVENT, lose EVENT	62/100
WSG-Prototype	come to GROUP, star in ACT, compete in ACT	6/100
2.1 Verb: <i>cover</i> ; target sense: be responsible for reporting (news) ; true source sense: to form a cover over Usage context: Generally, only reporters who <i>cover</i> breaking news are eligible.		
BERT-MLM	work (operate in a place), take (be a student), write (communicate by writing)	78/100
WSG-Prototype	practice (carry out as job), sponsor (be a client of), monitor (supervise someone)	10/100
2.2 Verb: <i>stretch</i> ; target sense: to extend the scope or meaning of ; true source sense: to make longer by pulling Usage context: ... the usage of a new "entwinement" standard ... <i>stretched</i> the doctrine beyond its permissible limits.		
BERT-MLM	spoil (alter), meet (satisfy), understand (comprehend)	59/100
WSG-Prototype	stem (remove the stem), stretch (make longer by pulling), extend (stretch to a greater length)	2/100
2.3 Verb: <i>grasp</i> ; target sense: to get the meaning of ; true source sense: to hold firmly Usage context: Madonna later acknowledged that she had not <i>grasped</i> the concept of her mother dying.		
BERT-MLM	appreciate (realize fully), understand (comprehend), enjoy (get pleasure from)	86/100
WSG-Prototype	understand (comprehend), resolve (understand the meaning of), read (interpret)	38/100

Table 3: Example novel sense predictions from the full prototype model and the BERT-MLM baseline (both trained on WSG) on domain-based dataset (first 3 examples) and synset-based dataset (last 3 examples).

Model	Mean Euclidean distance	
	Training	Testing
BERT	11.85±2.93	11.96±2.96
+ WSG-Prototype	1.14±0.38	1.67±0.43
+ WSG-Exemplar	1.21±0.34	1.82±0.50

Table 4: Mean Euclidean distance between the prototypes of source and target tokens.

Model	MLM perplexity	STS correlation
BERT	0.337	0.665
+WSG-Prototype	0.339	0.632
+WSG-Exemplar	0.352	0.634

Table 5: Performance on common NLP tasks for the vanilla BERT model and WSG models.

forming well-separated clusters for tokens pairs from the same extensional type.

6.3 Evaluation on common NLP tasks

We finally examined how learning WSG might affect a language model on common NLP tasks. Table 5 shows performance on two general NLP tasks for the BERT encoders with and without supervised learning on WSG: 1) masked language modeling (MLM) on Wikitext-103 dataset; and 2) seman-

tic textual similarity (STS) prediction on the STS benchmark dataset (Cer et al., 2017). We found that contextualized language models trained on WSG largely preserved their linguistic competence on both tasks, suggesting that our WSG learning framework is not simply a fine-tuning technique designed for a specific task, but rather a more fundamental exercise of learning flexible embedding spaces that may improve the semantic generalizability of language models.

7 Conclusion

We have presented a framework of word sense generation that supports lexical items to spawn new senses in scenarios that involve novel context. Our results show that chaining provides a general mechanism for extending to novel senses including metaphor and metonymy, and learning a transformed sense-extensional space enables systematic generalization in word sense extension. In contrast with the established traditions in word sense disambiguation and induction, our work emphasizes a generative approach to model word senses that is scalable and applicable to natural sentences and a broad set of words in the lexicon. Future work may extend our framework to incorporate temporal dimensions, other word classes and languages.

8 Ethical considerations

In this section, we discuss the limitations and potential risks of our work.

8.1 Limitations

Our work has some limitations. First, the current study considers sense generation for English verbs, and it therefore does not account for all common types of sense extension in English (e.g., English noun-to-verb conversion, as discussed in Clark and Clark 1979; Yu and Xu 2022) and other languages. Future work can consider extending our WSG framework more broadly to other word classes and languages.

Second, our framework does not explicitly consider the temporal order via which word senses have emerged. In particular, in the data collection step, we choose source-target sense pairs based on random partitioning of word sense inventories, whereas an alternative approach would be to sort all senses of a word chronologically by their times of emergence in history, and use the model to incrementally predict each sense of a word based on usages of its older senses. However, we found that it is infeasible to find accurate timestamps of senses in natural corpora at a comprehensive scale. Another approach is to have human annotators evaluate the plausibility of each possible source-target sense pairs against sampled alternatives, which is a potential area for future extension.

8.2 Potential risks

All scientific artifacts in this study have been made publicly available and are consistent with their intended use and access conditions. We acknowledge that our focus on English might introduce linguistically or culturally specific biases in model-generated outputs. For instance, we observe that the WSG models trained on English sentences learn to generate a metaphorical expression “to spend some time” for the English verb *spend*, which is common in English but differ in other languages (e.g., Hungarian speakers instead tend to say “to fill some time” as in Kövecses et al. 2010). We believe that by training WSG models cross-linguistically to cover various innovative lexical uses should help alleviate this issue.

References

- Michele Bevilacqua and Roberto Navigli. 2020. Breaking through the 80% glass ceiling: Raising the state of the art in word sense disambiguation by incorporating knowledge graph information. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2854–2864.
- Claudia Brugman and George Lakoff. 1988. Cognitive topology and lexical networks. In *Lexical ambiguity resolution*, pages 477–508. Elsevier.
- Marc Brysbaert, Amy Beth Warriner, and Victor Kuperman. 2014. Concreteness ratings for 40 thousand generally known english word lemmas. *Behavior research methods*, 46(3):904–911.
- Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. Semeval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14.
- Tuhin Chakrabarty, Xurui Zhang, Smaranda Muresan, and Nanyun Peng. 2021. Mermaid: Metaphor generation with symbolism and discriminative decoding. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4250–4261.
- Eve V Clark and Herbert H Clark. 1979. When nouns surface as verbs. *Language*, pages 767–811.
- Herbert H Clark and Richard J Gerrig. 1983. Understanding old words with new meanings. *Journal of verbal learning and verbal behavior*, 22(5):591–608.
- Paul Cook, Jey Han Lau, Diana McCarthy, and Timothy Baldwin. 2014. Novel word-sense identification. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 1624–1635.
- Ann Copestake and Ted Briscoe. 1995. Semi-productive polysemy and sense extension. *Journal of semantics*, 12(1):15–67.
- 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.
- Robert B Dewell. 1994. Overagain: Image-schema transformations in semantic analysis. *Cognitive Linguistics*, 5(4).

651	William A Gale, Kenneth Church, and David Yarowsky.	Jey Han Lau, Paul Cook, Diana McCarthy, David New-	702
652	1992. Estimating upper and lower bounds on the per-	man, and Timothy Baldwin. 2012. Word sense in-	703
653	formance of word-sense disambiguation programs.	duction for novel sense detection. In <i>Proceedings</i>	704
654	In <i>30th Annual Meeting of the Association for Com-</i>	<i>of the 13th Conference of the European Chapter of</i>	705
655	<i>putational Linguistics</i> , pages 249–256.	<i>the Association for Computational Linguistics</i> , pages	706
		591–601.	707
656	Raymond W Gibbs Jr and Herbert L Colston. 2008.	Adrienne Lehrer. 1990. Polysemy, conventionality, and	708
657	Image schema. In <i>Cognitive Linguistics: Basic Read-</i>	the structure of the lexicon. <i>Cognitive Linguistics</i> ,	709
658	<i>ings</i> , pages 239–268. De Gruyter Mouton.	1(2).	710
659	Amir Ahmad Habibi, Charles Kemp, and Yang Xu.	Daniel Loureiro and Alipio Jorge. 2019. Language	711
660	2020. Chaining and the growth of linguistic cate-	modelling makes sense: Propagating representations	712
661	gories. <i>Cognition</i> , 202:104323.	through wordnet for full-coverage word sense disam-	713
662	Luyao Huang, Chi Sun, Xipeng Qiu, and Xuan-Jing	biguation. In <i>Proceedings of the 57th Annual Meet-</i>	714
663	Huang. 2019. Glossbert: BERT for word sense dis-	<i>ing of the Association for Computational Linguistics</i> ,	715
664	ambiguation with gloss knowledge. In <i>Proceedings</i>	pages 5682–5691.	716
665	<i>of the 2019 Conference on Empirical Methods in Nat-</i>	Barbara C Malt, Steven A Sloman, Silvia Gennari,	717
666	<i>ural Language Processing and the 9th International</i>	Meiyi Shi, and Yuan Wang. 1999. Knowing ver-	718
667	<i>Joint Conference on Natural Language Processing</i>	us naming: Similarity and the linguistic categoriza-	719
668	<i>(EMNLP-IJCNLP)</i> , pages 3509–3514.	tion of artifacts. <i>Journal of Memory and Language</i> ,	720
		40(2):230–262.	721
669	Ignacio Iacobacci, Mohammad Taher Pilehvar, and	Stephen Merity, Caiming Xiong, James Bradbury, and	722
670	Roberto Navigli. 2016. Embeddings for word sense	Richard Socher. 2016. Pointer sentinel mixture mod-	723
671	disambiguation: An evaluation study. In <i>Proceed-</i>	els. <i>arXiv preprint arXiv:1609.07843</i> .	724
672	<i>ings of the 54th Annual Meeting of the Association for</i>	George A Miller. 1995. Wordnet: a lexical database for	725
673	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	english. <i>Communications of the ACM</i> , 38(11):39–41.	726
674	pages 897–907.	Gregory L Murphy. 1997. Reasons to doubt the present	727
675	Adam Kilgarriff. 1997. I don’t believe in word senses.	evidence for metaphoric representation. <i>Cognition</i> ,	728
676	<i>Computers and the Humanities</i> , 31(2):91–113.	62(1):99–108.	729
677	Adam Kilgarriff and Joseph Rosenzweig. 2000. Frame-	Robert M Nosofsky. 1986. Attention, similarity, and the	730
678	work and results for english senseval. <i>Computers</i>	identification–categorization relationship. <i>Journal of</i>	731
679	<i>and the Humanities</i> , 34(1):15–48.	<i>Experimental Psychology: General</i> , 115(1):39.	732
680	Diederik P Kingma and Jimmy Ba. 2015. Adam: A	Geoffrey Nunberg. 1979. The non-uniqueness of seman-	733
681	method for stochastic optimization. In <i>ICLR (Poster)</i> .	tic solutions: Polysemy. <i>Linguistics and philosophy</i> ,	734
682	Ekaterini Klepousniotou, Debra Titone, and Carolina	pages 143–184.	735
683	Romero. 2008. Making sense of word senses: the	James Pustejovsky. 1998. <i>The generative lexicon</i> . MIT	736
684	comprehension of polysemy depends on sense over-	press.	737
685	lap. <i>Journal of Experimental Psychology: Learning,</i>	James Pustejovsky and Anna Rumshisky. 2010. Mecha-	738
686	<i>Memory, and Cognition</i> , 34(6):1534.	nisms of sense extension in verbs.	739
687	Zoltán Kövecses et al. 2010. Metaphor and culture.	Alessandro Raganato, Jose Camacho-Collados, and	740
688	<i>Acta Universitatis Sapientiae, Philologica</i> , 2(2):197–	Roberto Navigli. 2017. Word sense disambiguation:	741
689	220.	A unified evaluation framework and empirical com-	742
690	George Lakoff. 1987. <i>Women, fire, and dangerous</i>	parison. In <i>Proceedings of the 15th Conference of</i>	743
691	<i>things: What categories reveal about the mind</i> . Uni-	<i>the European Chapter of the Association for Compu-</i>	744
692	versity of Chicago press.	<i>tational Linguistics: Volume 1, Long Papers</i> , pages	745
693	George Lakoff and Mark Johnson. 2008. <i>Metaphors we</i>	99–110.	746
694	<i>live by</i> . University of Chicago press.	Giulia Rambelli, Emmanuele Chersoni, Alessandro	747
695	Jey Han Lau, Paul Cook, Diana McCarthy, Spandana	Lenci, Philippe Blache, and Chu-Ren Huang.	748
696	Gella, and Timothy Baldwin. 2014. Learning word	2020. Comparing probabilistic, distributional and	749
697	sense distributions, detecting unattested senses and	transformer-based models on logical metonymy inter-	750
698	identifying novel senses using topic models. In <i>Pro-</i>	pretation. In <i>Proceedings of the 1st Conference of the</i>	751
699	<i>ceedings of the 52nd Annual Meeting of the Associa-</i>	<i>Asia-Pacific Chapter of the Association for Compu-</i>	752
700	<i>tion for Computational Linguistics (Volume 1: Long</i>	<i>tational Linguistics and the 10th International Joint</i>	753
701	<i>Papers)</i> , pages 259–270.	<i>Conference on Natural Language Processing (ACL-</i>	754
		<i>IJCNLP)</i> .	755

756	Christian Ramiro, Mahesh Srinivasan, Barbara C Malt, and Yang Xu. 2018. Algorithms in the historical emergence of word senses. <i>Proceedings of the National Academy of Sciences</i> , 115(10):2323–2328.	810
757		811
758		812
759		813
760	Jennifer M Rodd, Richard Berriman, Matt Landau, Theresa Lee, Carol Ho, M Gareth Gaskell, and Matthew H Davis. 2012. Learning new meanings for old words: Effects of semantic relatedness. <i>Memory & Cognition</i> , 40(7):1095–1108.	814
761		815
762		816
763		817
764		818
765	Eleanor Rosch. 1975. Cognitive representations of semantic categories. <i>Journal of Experimental Psychology: General</i> , 104(3):192.	819
766		820
767		821
768	Anna Rumshisky and Olga Batiukova. 2008. Polysemy in verbs: Systematic relations between senses and their effect on annotation. In <i>Coling 2008: Proceedings of the workshop on Human Judgements in Computational Linguistics</i> , pages 33–41, Manchester, UK. Coling 2008 Organizing Committee.	822
769		823
770		824
771		825
772		826
773		827
774	Jake Snell, Kevin Swersky, and Richard Zemel. 2017. Prototypical networks for few-shot learning. In <i>Advances in Neural Information Processing Systems</i> , pages 4077–4087.	828
775		829
776		830
777		831
778	Mahesh Srinivasan and Susan Carey. 2010. The long and the short of it: On the nature and origin of functional overlap between representations of space and time. <i>Cognition</i> , 116(2):217–241.	832
779		833
780		834
781		835
782	Kevin Stowe, Tuhin Chakrabarty, Nanyun Peng, Smaranda Muresan, and Iryna Gurevych. 2021. Metaphor generation with conceptual mappings. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6724–6736.	836
783		837
784		
785		
786		
787		
788		
789		
790	Zhewei Sun, Richard Zemel, and Yang Xu. 2021. A computational framework for slang generation. <i>Transactions of the Association for Computational Linguistics</i> , 9:462–478.	
791		
792		
793		
794	Paul Thibodeau and Frank H Durgin. 2008. Productive figurative communication: Conventional metaphors facilitate the comprehension of related novel metaphors. <i>Journal of Memory and Language</i> , 58(2):521–540.	
795		
796		
797		
798		
799	Xiaoyu Tong, Ekaterina Shutova, and Martha Lewis. 2021. Recent advances in neural metaphor processing: A linguistic, cognitive and social perspective. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4673–4686.	
800		
801		
802		
803		
804		
805		
806	Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. 2016. Matching networks for one shot learning. <i>Advances in Neural Information Processing Systems</i> , 29:3630–3638.	
807		
808		
809		
	Gregor Wiedemann, Steffen Remus, Avi Chawla, and Chris Biemann. 2019. Does bert make any sense? interpretable word sense disambiguation with contextualized embeddings. <i>arXiv preprint arXiv:1909.10430</i> .	
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2020. Transformers: State-of-the-art natural language processing. In <i>Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations</i> , pages 38–45.	
	Zhibiao Wu and Martha Palmer. 1994. Verbs semantics and lexical selection. In <i>Proceedings of the 32nd annual meeting on Association for Computational Linguistics</i> , pages 133–138.	
	Lei Yu and Yang Xu. 2021. Predicting emergent linguistic compositions through time: Syntactic frame extension via multimodal chaining. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 920–931.	
	Lei Yu and Yang Xu. 2022. Noun2Verb: Probabilistic Frame Semantics for Word Class Conversion. <i>Computational Linguistics</i> , pages 1–36.	
	Zhi Zhong and Hwee Tou Ng. 2010. It makes sense: A wide-coverage word sense disambiguation system for free text. In <i>Proceedings of the ACL 2010 system demonstrations</i> , pages 78–83.	

A Implementations of WSG models

We use the bert-base-uncased configuration provided by Hugging Face (Wolf et al., 2020) to initialize all BERT-based WSG models (two baselines and two chaining-based models). During MLM pretraining of BERT models to learn novel partitioned tokens, we randomly mask 15% of tokens in each sentence, and train each model on predicting the masked tokens. Learning is performed using the Adam optimizer (Kingma and Ba, 2015), with a learning rate of 5e-5 and a batch size of 128, for 8 epochs (after which all models achieved highest evaluation accuracy). During sense-extensional semantic space learning, all chaining-based models are trained on the objective function in Eq.3.3 using Adam, with a mini-batch size of 16 and a learning rate of 2e-5, for 8 epochs (after which all models achieved highest evaluation accuracy). Two baseline models are trained on the same set of usage data using Eq.2 and Eq.3 respectively, with an Adam learning rate of 2e-5 and a batch size of 32, for 4 epochs (after which all models achieved highest evaluation accuracy). All experiments are run on machines with 4 NVIDIA Tesla V100 GPUs, with an average training time of 25 minutes per epoch for MLM pretraining, and 10 minutes per epoch for sense-extensional semantic space learning.

B Description of supersense-based domains

WordNet organize noun synsets into 26 supersense domains based on their semantic coherence. During data collection, we first run the WSD model of (Bevilacqua and Navigli, 2020) on usage sentences to identify the WordNet synset label for each noun object, and then tag it with the associated supersense label. We then take the top-11 most frequent supersenses with at least 10,000 identified noun objects to construct our domain-based WSG dataset. Table 6 shows detailed information about the 11 selected supersense categories.

C Results on WSG examples with decreasing sense concreteness

Research in cognitive science suggests that concrete and embodied word senses tend to be extended to more abstract senses to achieve better communication efficiency and learnability (Murphy, 1997; Srinivasan and Carey, 2010; Thibodeau

and Durgin, 2008). We therefore also compute performance scores of our WSG models on subsets of the domain-based usage dataset where the meaning of a verb is extended from concrete senses to a more abstract one (e.g. the metaphorical extensions of *arrive* as shown in Figure 1). In particular, we define a source-target sense pair to be an example of concrete-to-abstract sense extension if the average object noun concreteness score of the target token is lower than that of its source sibling. Table C summarizes the results, from which we observe that almost all models perform slightly worse on concrete-to-abstract extensions compared to the general setup, while still making significantly better predictions than the random baseline.

Supersense	Definition	Sample nouns
noun.act	nouns denoting acts or actions	attempt, performance, exchange
noun.artifact	nouns denoting man-made objects	aircraft, phone, guitar
noun.attribute	nouns denoting attributes of people and objects	popularity, style, power
noun.cognition	nouns denoting cognitive processes and contents	viewpoint, imagination, scheme
noun.communication	nouns denoting communicative processes and contents	proposal, screenplay, film
noun.event	nouns denoting natural events	tournament, final, competition
noun.group	nouns denoting groupings of people or objects	family, league, party
noun.location	nouns denoting spatial position	province, area, neighborhood
noun.person	nouns denoting people	officer, visitor, mother
noun.state	nouns denoting stable states of affairs	existence, friendship, injury
noun.time	nouns denoting time and temporal relations	era, day, season

Table 6: Definitions and sample members of the WordNet noun supersenses used to create the domain-based WSG dataset.

Model	Mean reciprocal rank (MRR-100)			
	Full partition		Leave-one-out partition	
	Unsupervised	Supervised	Unsupervised	Supervised
BERT-STS	11.5±1.2	34.6±2.0	37.7±1.8	61.6±1.6
BERT-MLM	13.6±1.4	35.8±1.7	43.0±2.2	64.1±0.9
BERT-Prototype	17.5±1.1	49.9±0.5	51.6±1.5	81.4±1.9
BERT-Exemplar	17.8±1.1	45.8±1.4	75.3±2.9	79.7±1.0

Table 7: Summary of model MRR-100 scores (%) on the subset of domain-based word sense generation dataset where the source sense has higher concreteness score than its target sibling token (i.e. extensions from concrete to abstract senses). Numbers after \pm are standard deviations over 5 sets of independently sampled negative candidate tokens.