# Evaluating language models for the retrieval and categorization of lexical collocations

**Luis Espinosa-Anke[1], Joan Codina-Filbà[2], Leo Wanner[3,2]**

[1]School of Computer Science and Informatics, Cardiff University, UK

[2]TALN Research Group, Pompeu Fabra University, Barcelona, Spain

[3]Catalan Institute for Research and Advanced Studies (ICREA), Barcelona, Spain

`espinosa-ankel@cardiff.ac.uk`, `joan.codina|leo.wanner@upf.edu`

## Abstract

Lexical collocations are idiosyncratic combinations of two syntactically bound lexical items (e.g., *"heavy rain"*, *"take a step"* or *"undergo surgery"*). Understanding their degree of compositionality and idiosyncrasy, as well their underlying semantics, is crucial for language learners, lexicographers and downstream NLP applications alike. In this paper we analyse a suite of language models for collocation understanding. We first construct a dataset of apparitions of lexical collocations in context, categorized into 16 representative semantic categories. Then, we perform two experiments: (1) unsupervised collocate retrieval, and (2) supervised collocation classification in context. We find that most models perform well in distinguishing light verb constructions, especially if the collocation's first argument acts as a subject, but often fail to distinguish, first, different syntactic structures within the same semantic category, and second, finer-grained categories which restrict the set of correct collocates[1].

## 1 Introduction

Language models (LMs) such as BERT (Devlin et al., 2018), and its variants SpanBERT (Joshi et al., 2020), ALBERT (Lan et al., 2019), RoBERTa (Liu et al., 2019), etc. have proven extremely flexible, as they behave as unsupervised multitask learners (Radford et al., 2019), and can be leveraged in a wide array of NLP tasks almost out-of-the-box (see, e.g., the GLUE and SuperGLUE results in Wang et al. (2019b) and Wang et al. (2019a), respectively). They have also been harnessed as supporting resources for knowledge-based NLP (Petroni et al., 2019), as they capture a wealth of linguistic phenomena (Rogers et al., 2020). Recently, a great deal of research analyzed the degree to which they encode, e.g., morphological (Edmiston, 2020), syntactic (Hewitt and Manning, 2019), or lexico-semantic structures (Joshi et al., 2020). However, less work explored so far how LMs interpret phraseological units at various degrees of compositionality. This is crucial for understanding the suitability of different text representations (e.g., static vs contextualized word embeddings) for encoding different types of multiword expressions (Shwartz and Dagan, 2019), which, in turn, can be useful for extracting latent world or commonsense information (Zellers et al., 2018).

One central type of phraselogical units are *lexical collocations*, defined as restricted co-occurrences of two syntactically bound lexical items (Kilgarriff, 2006), such that one of the items (the *base*) conditions the selection of the other item (the *collocate*) to express a specific meaning. For instance, the base *lecture* conditions the collocates *give* or *deliver* to express the meaning 'perform', the base *applause* conditions the selection of the collocate *thunderous* to express the meaning 'intense', and so on. Lexical collocations are of high relevance to lexicography, NLP and second language learning alike, and constitute a challenge for computational models because of their heterogeneity in terms of idiosyncrasy and degree of semantic composition (Mel'čuk, 1995).

In this paper, we analyze a suite of LMs in the context of two tasks that involve lexical collocation modeling. First, **unsupervised collocate retrieval**, where we *mask* a collocation's collocate (e.g., *"heavy"* in *"heavy rain"*), and quantify how well a LM of choice (BERT in particular) predicts, via its masked language modeling (MLM) objective, a valid collocate for that particular base ({*"heavy"*, *"torrential"*, *"violent"*, ...} for the base *"rain"* and the meaning *intense*). Second, **su-**

---

[1]The resources associated with this paper are available at https://github.com/luisespinosaanke/lexicalcollocations.

**pervised in-context collocation categorization**, where we fine-tune LMs on the task of predicting a semantic category of a collocation in terms of its *lexical function* (LF), given its sentential context; cf. Section 3.1. Modeling, recognizing, and classifying collocations in corpora has obvious applications for automatically creating and expanding lexicographic resources, as well as for various downstream NLP applications, among them, e.g., machine translation (Seretan, 2014), word sense disambiguation (Maru et al., 2019), or natural language generation (Wanner and Bateman, 1990). The two main contributions of this paper thus are:

1. A "collocations-in-context" dataset, with instances of collocations of 16 different semantic categories (in terms of LFs) in context, and with a fixed and lexical (i.e., no overlapping) train/dev/test split (Section 3).

2. An evaluation framework for assessing the degree of compositionality of lexical collocations, pivoting around two tasks: unsupervised *collocate* retrieval (Section 4) and in-context *collocation* categorization (Section 5).

Our results suggest that modeling collocations in context is a challenge, even for widely used LMs, and that this is particularly true for less semantic (and thus less compositional and more idiosyncratic) collocations. We also find that jointly recognizing the semantics and the syntactic structure (e.g., whether the collocate acts as subject or object in verbal constructions) of a collocation also constitutes non-trivial challenges for current architectures. Moreover, as a byproduct of our analysis, we also find an interesting behaviour in LMs when modeling antonymy in adjectives, specifically that their representations undergo substantial transformations as they flow through BERT's transformer layers, with many contextualized embeddings clustered together in the tip of a narrow cone that seems to represent adjectives in collocations denoting intensity (*"heavy"* rain) and weakness (*"minor"* issue).

## 2 Related Work

In this section, we discuss related works in two methodological areas that are relevant to this paper, namely conditioning MLMs (Section 2.1) and recognition of multiword epressions (MWEs) (Section 2.2).

### 2.1 Conditioning MLMs

A Masked Language Model (MLM) can be used as a *proxy* for gaining insights into how language is encoded by the weights of the (usually transformer-based) LM architecture. Moreover, simply asking an LM to predict words in context (without task-specific fine-tuning) has proved useful in NLP applications dealing with lexical items (affixes, words or phrases). For example, Wu et al. (2019) use BERT's MLM for augmenting their training data in sentiment analysis tasks; Qiang et al. (2019) use BERT for lexical simplification by *conditioning* the predictions over the [MASK] token by providing the original sentence as context; and Zhou et al. (2019) obtain SotA results in lexical substitution by conditioning BERT via embedding dropout on the target (unmasked) word. Inspired by the findings in these works (especially Qiang et al. (2019)), we will explore the predictions of BERT over masked lexical collocations (with and without conditioning) in Section 4, with the aim to understand whether these predictions can be used to measure the idiosyncrasy of the underlying semantics of a lexical collocation, i.e., whether the restrictions imposed by a collocation's base are due to the frozenness of the phrase itself or, on the contrary, sentential context is neccessary.

### 2.2 Distributional Lexical Composition

Building representations that account for non-compositional meanings within the broader spectrum of encoding semantic relations between words is a long-standing problem in computational semantics (Baroni and Zamparelli, 2010; Mitchell and Lapata, 2010; Boleda et al., 2013). Interestingly, there seems to be little agreement on how these representations should be defined, with recent attempts focusing on verbal multiword expressions (see an overview of approaches in Ramisch et al. (2018)), phrases of variable length encoded via LSTMs, based on their definitions (Hill et al., 2016), or arbitrary lexical and commonsense relations between word pairs for downstream NLP. As a testimony of the broad methods explored in the most recent literature, let us refer to, for instance, the combination of word vector averages with conditional autoencoders (Espinosa-Anke and Schockaert, 2018), expectation maximization (Camacho-Collados et al., 2019), LSTMs for predicting word pair contexts (Joshi et al., 2019), and explicit encoding of generalized lexico-syntactic patterns (Washio and Kato,

| LF | BERT input | Masked sentence | Pred. collocates | Orig. collocate |
|---|---|---|---|---|
| Oper1 | Masked sent | iran feared that the u.s and israel may [MASK] an air raid on its controversial nuclear facilities | perform, conduct, mount, | launch |
| | Orig sent [SEP] | iran feared that the u.s and israel may launch an air raid on its controversial nuclear facilities [SEP] | launch, conduct, order | launch |
| | Masked sent | iran feared that the u.s and israel may [MASK] an air raid on its controversial nuclear facilities | | |
| Real1 | Masked sent | if that happens, lindsey will use explosive triggers to [MASK] the landing gear | destroy, stop, remove | lower |
| | Orig sent [SEP] | if that happens , lindsey will use explosive triggers to lower the landing gear [SEP] | destroy, stop, trigger | lower |
| | Masked sent | if that happens , lindsey will use explosive triggers to [MASK] the landing gear | | |
| Magn | Masked sent | the EU is driving a [MASK] bargain over Swiss demands for greater access to EU airspace . | plea, hard, new | hard |
| | Orig sent [SEP] | the EU is driving a hard bargain over Swiss demands for greater access to EU airspace . [SEP] | plea, hard, tough | hard |
| | Masked sent | the EU is driving a [MASK] bargain over Swiss demands for greater access to EU airspace . | | |

Table 1: Illustrative behaviour of BERT when prompted to predict a collocate in the position of a masked token, for three LFs: Oper1 ('launch an air raid'), Real1 ('lower the landing gear') and Magn ('hard bargain'), under two settings, when not conditioned (Masked sent) and when conditioned on the original (unmasked) sentence (Orig sent [SEP] Masked sent).

2018). Parting ways with the above works, in this paper we will follow the experimental setting described in Shwartz and Dagan (2019), based on injecting sentential contexts into multiword expressions (in our case, only lexical collocations) to leverage the contextual nature of current LMs. However, our goal is not to compare different combinations of feature-extraction and training/fine-tuning methods, but rather to understand lexical collocations' learnability, idiosyncrasy and their internal vector-space representations.

## 3 Data and Resources

### 3.1 Lexical Collocations

Let us first introduce the notion of lexical collocation and LF. The term *collocation* has been used in computational linguistics research to denote two different concepts. On the one hand, following Firth (1957), Church and Hanks (1989); Evert (2007); Pecina (2008) and others, a collocation has been assumed to be a combination of words that have the tendency to occur together in discourse. Typical examples are *doctor – hospital*, *mop – bucket*, *real – estate*, *look – for*, etc. On the other hand, for instance, Wanner et al. (2006); Gelbukh and Kolesnikova. (2012); Rodríguez Fernández et al. (2016); Garcia et al. (2017) adopt the definition that is common in lexicography and phraseology (Hausmann, 1985; Cowie, 1994; Mel'čuk, 1995), according to which, a collocation is an idiosyncratic combination of two lexical items, the base and the collocate, as defined above in Section 1. This interpretation states that collocations are phraseological units, although their degree of compositionality can vary. For instance, *win [a] war* is perceived to possess a higher degree of (free) composition than, e.g., *hold [a] meeting*, and *heavy*

*rain* is less compositional than *[a] well-justified argument*. We adopt this definition of the notion of collocation, and in order to avoid any confusion, we refer to it, following Krenn (2000), as *lexical collocation*.

Lexical collocations can be typified with respect to the meaning of the collocate and the syntactic structure formed by the base and the collocate. LFs provide a fine-grained typology of this kind (Mel'čuk, 1996). An LF can be considered a function $f(L)$ that delivers for a base $L$ a set of synonymous collocates that express the meaning of $f$. Where pertinent, $f$ also codifies the subcategorization structure of the base+collocate combination. LFs are assigned Latin acronyms as names; cf., e.g., "Oper1" ('operare'), which means 'perform' and realizes the first argument of the base as subject: Oper1(*lecture*) = {*deliver*, *give*, *hold*}; "Magn" ('magnum'), which stands for 'intense': Magn(*applause*) = {*thunderous*, *loud*, ...}.[2]

The encoding of LFs in NLP research has in recent years revolved around applying word embeddings-based techniques, e.g., in terms of linear projections (Rodríguez Fernández et al., 2016) and semantic generalizations (Espinosa-Anke et al., 2016). Recently, Shwartz and Dagan (2019) analyzed, from the perspective of "static" vs. "contextualized" representations and their applicability to studying compositional phenomena like "meaning shift", one specific type of lexical collocations, namely *light verb constructions* (LVCs), which are well illustrated by the LFs Oper1 and Oper2 (and also, partially, by Real1 and Real2). While we find that the above research directions (i.e., embeddings-based and contextualized representations for modeling MWEs) are complementary, in this work, we

---

[2]For simplicity we will write Oper1(*lecture*) = *deliver*, etc.

| LF | semantic gloss | example |
|---|---|---|
| Oper1 | 'perform'; 1st argument → subject | Oper1(*support*) = *lend* |
| IncepOper1 | 'begin to perform'; 1st argument → subject | IncepOper1(*impression*) = *gain* |
| Oper2 | 'undergo'; 2nd argument → subject | Oper2(*support*) = *find* |
| Real1 | 'realize'; 1st argument → subject | Real1(*accusation*) = *prove* |
| Real2 | 'apply'; 2nd argument → subject | Real2(*support*) = *enjoy* |
| AntiReal2 | 'fail to apply'; 2nd argument → subject | AntiReal2(*war*) = *lose* |
| CausFunc0 | 'cause the existence' | CausFunc0(*hope*) = *raise* |
| Caus1Func0 | 'cause the existence; 1st argument' | Caus1Func0(*hope*) = *gain* |
| LiquFunc0 | 'cause termination of the existence' | LiquFunc0(*hope*) = *destroy* |
| IncepPredPlus | 'increase' | IncepPredPlus(*temperature*) = *rise* |
| Magn | 'intense' | Magn(*smoker*) = *heavy* |
| AntiMagn | 'little', 'weak' | AntiMagn(*smoker*) = *occasional* |
| Ver | 'genuine' | Ver(*demand*) = *legitimate* |
| AntiVer | 'non-genuine' | AntiVer(*demand*) = *illegitimate* |
| Bon | 'positive' | Bon(*performance*) = *good* |
| AntiBon | 'negative' | AntiBon(*performance*) = *poor* |

Table 2: LFs used in this paper. The 'semantic gloss' column provides both a definition and the actantial structure, which is required in cases where one LF may express the same semantics but with a different syntactic structure (e.g., Real1 vs. Real2).

specifically focus on the existing (and learnable) knowledge LMs have concerning lexical collocations, and whether they can be used to recognize and categorize LFs in free text.

For our experiments, we use, as initial lexical collocation source, a **collocations dataset**, LEXFUNC (Espinosa-Anke et al., 2019), which we have extended to cover a wider range of LFs (listed in Table 2). The original LEXFUNC dataset and this extended version are both the result of an initial collection of collocations categorized into LFs made available by Igor Mel'čuk. Each collocation has been manually lemmatized, and bases and collocates have been manually annotated with part-of-speech tags and their syntactic dependency relation.

With the lexical collocations of the extended LEXFUNC dataset at hand, we first compile from the English Gigaword[3] a **collocations corpus**, which contains the occurrences of these lexical collocations. In principle, the identification of a given collocation in corpora is a straightforward procedure, as we know its elements (base and collocate) and the syntactic dependency relation between them. However, automatic dependency parsing is far from perfect, which complicates the task. Therefore, and in order not to lose any relevant collocation occurrence in the GigaWord corpus, we apply a cascaded procedure for their identification on the lemmatized and POS- and head-modifier relation tagged Giga-Word.[4] In the first stage, we identify sentences in which between the collocation elements in question one of the relevant syntactic dependency relations has been identified. In the second (more relaxed) stage, we match adjacent lemmatized collocation elements and their PoS tags. In the third stage, finally, we match lemmatized collocation elements and their PoS tags within a distance of up to 5 tokens. While this procedure inevitably introduces some noise (we might retrieve sentences where base and collocate co-occur, but not as a collocation), we performed a manual inspection on a random sample, and calculated precision of our collocation retrieval strategy, which resulted in >0.95. This confirms the quality of our retrieval strategy, and hence, our resource.

In terms of corpus statistics, Table 3 indicates the number of sentences for each LF distributed across training (70% of the sentences), development (15%) and test (15%) sets. The split was done maintaining this proportion across all LF. Note that these splits are constructed such that there are no overlapping collocations, in an effort to avoid the well-known phenomenon of *lexical memorization* (Levy et al., 2015), which may artificially inflate the results on the test set. The number of different collocations per split, globally and for each LF, also maintains the same proportions (70/15/15 ±1%), such that, e.g., AntiReal2 has 55 different collocations in the 942 sentences of the training set and 11 different collocations in the development and test sets, distributed across 205 and 200 sen-

| Label | Train | Dev | Test | Total |
|---|---|---|---|---|
| Magn | 28,748 | 6,113 | 6,164 | 41,025 |
| Oper1 | 11,746 | 2,517 | 2,493 | 16,756 |
| Real1 | 3,481 | 746 | 743 | 4,970 |
| AntiMagn | 2,959 | 638 | 649 | 4,246 |
| IncepOper1 | 2,489 | 541 | 539 | 3,569 |
| Oper2 | 2,408 | 515 | 520 | 3,443 |
| AntiVer | 1,874 | 397 | 405 | 2,676 |
| AntiBon | 1,815 | 385 | 393 | 2,593 |
| CausFunc0 | 1,714 | 370 | 367 | 2,451 |
| Real2 | 1,570 | 336 | 337 | 2,243 |
| Bon | 1,471 | 298 | 315 | 2,084 |
| LiquFunc0 | 1,398 | 301 | 297 | 1,996 |
| AntiReal2 | 942 | 203 | 200 | 1,345 |
| Ver | 926 | 198 | 196 | 1,320 |
| Caus1Func0 | 686 | 151 | 149 | 986 |
| IncepPredPlus | 624 | 147 | 138 | 909 |
| Total | 64,851 | 13,856 | 13,905 | |

Table 3: Statistics (in number of sentences) of our *collocations in context* dataset (ordered by frequency).

tences respectively. In the overall corpus, there is an average of 18 samples per collocation (*work hard* being the most frequent one with 102 samples). *Hope*, *attack*, *criticism*, *fire* and *thread* are bases that each co-occur with more than 30 different collocates, across most LF. These bases are also among the ones with more samples in the corpus. On the other side, half of the bases are combined with one single collocate only. Overall, the statistical properties of our dataset arguably make it a faithful replica of the distribution of collocations in, at least, newswire corpora. At the same time, it is a challenging dataset, as the results we report in this paper suggest.

## 4 Experiment 1: Collocate Retrieval

### 4.1 Setup

In the first experiment, we aim to analyze how well an MLM retrieves valid collocates for a given base when being provided with the original (sentence-level) context. We use BERT (`bert-base`) (Devlin et al., 2018), as it is the *de-facto* model on top of most specialized and distilled/quantized language models. Its behaviour should thus be a good proxy for the general distributional behaviour of lexical collocations. This experiment serves, first, as an opportunity to understand *how much* semantics that is underlying LFs can be encoded via a MLM pretraining objective, and second, as a

testbed for exploring conditioning strategies often used in tasks involving data augmentation and lexical substitution and simplification (cf. Section 2.1). Since this is an "in-context collocate retrieval" task, we consider it a ranking problem. Intuitively, if BERT is able to retrieve a base's valid collocates (e.g., {*heavy*, *torrential*, *violent*, . . . } for *rain* as base for Magn) in the position of a masked token, this could mean that: (1) the sentence is giving enough context for the model to "know" the lexical restrictions involved in that collocation, and/or (2) the LF is sufficiently frozen, and therefore the base alone may restrict which collocates are acceptable. For the first point, and continuing with the *heavy rain* example, consider the following sentence.

(1) *Hurricane Katrina brought* [MASK] *rain to Louisiana.*

Intuitively, we would expect the sentential context to be informative enough for the model to select *heavy* or any other collocate denoting the notion of intensity, and restricted by the presence of the base *rain*. In fact, here, BERT predicts *heavy* with 79.5% probability. However, in example (2)

(2) *Policeman earns applause for staying on duty in* [MASK] *rain.*

there is much lesser evidence for the *rain* to be 'intense', and in fact BERT predicts here *'the'* with 85.1% probability. This disparity lets us investigate ways to prompt BERT to select *heavy* or any other valid collocate for example (2). Thus, in addition to simply passing one masked sentence, we explore an approach based on passing the masked sentence concatenated with the original sentence, which is a natural way to encode not only the context surrounding the word, but also the meaning of the target word itself. This strategy was successfully used for the task of unsupervised lexical simplification (Qiang et al., 2019). For the second point above, the works of Espinosa-Anke et al. (2019); Shwartz and Dagan (2019) already point to the fact that light verb constructions (LVCs) are easy to recognize in text. Further evidence is provided in Table 1, which shows BERT's top predictions for three sentences containing Oper1, Real1 and Magn collocations. It is immediately obvious that the nature of the LF itself, as well as the amount of the information provided by the sentential context, are crucial. Note that, in the case of Oper1, by providing the original sentence as context, BERT's grasp of the LF improves, as it tends to predict the correct collocate, whereas for Real1 or Magn, this

improvement results in assigning higher probability to tokens which are more similar to the LF's abstract meaning (e.g., *trigger the landing gear* for Real1 or *tough bargain* for Magn).

With the above considerations in mind, we run BERT's MLM head over our GigaWord-based test set, and compute *Mean Reciprocal Rank* (MRR) and *Mean Average Precision* (MAP) over each LF. Recall that we consider as valid *hits* all the collocates for a given base and its corresponding LF. In practice, this means that for bases that have just one valid collocate, both metrics yield the same score. We lemmatize BERT's predictions using SpaCy's lemmatizer.[5]

## 4.2 Results and Discussion

Our results (Table 4) show, first, that conditioning BERT's MLM by passing the original sentence as additional context for the [MASK] token is useful for predicting an embedding whose semantics is more related to the original collocate. The improvements are particularly relevant for LVCs (Oper1, and, to a certain extent, Real1 and Real2), suggesting that these LFs, while perhaps easy to distinguish from others (cf. Section 5.1), they do benefit from additional contexts to be well represented. Interestingly, the Magn LF has small gains in both MRR and MAP, clearly showing that additional context helps little, and thus highlighting a strong semantic dependency between sentence meaning and the collocation's base.

A potential limitation of this setup, however, is that we cannot possibly include all possible collocates for all the bases in our resource. An estimate of the quality of BERT's predictions can be obtained by measuring the semantic similarity (for instance, by cosine distance) between the original masked collocate and the predicted collocates. In the example we already referred to above, *heavy rain*, the similarity between 'the' and 'heavy' is low, whereas, if the model predicts *hard* or even any other adjective, it should be considered less wrong. We obtain a broad picture of the quality of BERT's predictions by plotting a histogram (Figure 1) of the similarities obtained by comparing the original collocate's and BERT's predicted GloVe embeddings (Pennington et al., 2014) under both settings (MASKED and CONDITIONED) for the same three LFs as in Table 1), namely Magn, Oper1 and Real1. The conditioning strategy is helpful; it con-
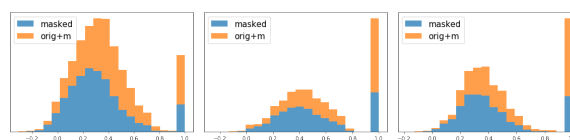
Figure 1: Histograms showing the distribution of similarities between gold and predicted (lemmatized) collocates for the three LFs Magn, Oper1 and Real1 (from left to right).

| | MASKED | | CONDITIONED | |
| --- | --- | --- | --- | --- |
| | MRR | MAP | MRR | MAP |
| AntiBon | 13.78 | 13.57 | 61.64 | 57.37 |
| AntiMagn | 30.16 | 27.11 | 82.18 | 66.89 |
| AntiReal2 | 39.14 | 36.32 | 70.58 | 62.13 |
| AntiVer | 11.38 | 10.68 | 40.5 | 37.12 |
| Bon | 25.13 | 24.87 | 62.49 | 59.34 |
| Caus1Func0 | 52.77 | 45.93 | 95.94 | 80.36 |
| CausFunc0 | 61.14 | 53.53 | 86.2 | 73.52 |
| IncepOper1 | 54.38 | 48.36 | 88.52 | 74.61 |
| IncepPredPlus | 7.27 | 6.375 | 12.27 | 10.67 |
| LiquFunc0 | 45.51 | 43.12 | 71.14 | 63.60 |
| Magn | 33.43 | 31.67 | 74.72 | 66.92 |
| Oper1 | 73.12 | 62.58 | 95.22 | 81.15 |
| Oper2 | 63.84 | 53.86 | 93.63 | 76.86 |
| Real1 | 59.87 | 53.31 | 90.96 | 75.02 |
| Real2 | 55.36 | 46.87 | 76.60 | 64.64 |
| Ver | 33.06 | 29.37 | 72.02 | 63.18 |

Table 4: Results of the collocate retrieval experiment when passing to the MLM the masked sentence (MASKED) or with the original sentence as context (CONDITIONED).

tributes not only to retrieving the original collocate (which would be trivial if we do not mask it), but also candidates with clearly similar meanings. We see, for instance, more cases for Oper1 and Real1, where the correct verb is predicted, whereas for Magn we see a more sustained improvement across all similarities, but not necessarily for retrieving the original collocate.

## 5 Experiment 2: Collocation categorization

In the second experiment, we test the performance of a number of well-known LMs for the task of LF categorization using the train/test splits we sampled and annotated from GigaWord (Section 3). This experiment serves two purposes. First, we expect to learn about the *predictability* of LFs in context,

which is a long-standing problem in computational lexicography and the cornerstone of automatic construction of collocation resources. Second, previous work has shown that some LFs are quite easy to distinguish, without (Espinosa-Anke et al., 2019) and with sentential context (Shwartz and Dagan, 2019). However, it is still unclear whether, by focusing exclusively on the phenomenon of collocations, and excluding, e.g., idiomatic expressions or non-compositional phrasal verbs (which are not only semantically but, more importantly, syntactically different from collocations), an LM can indeed be used to construct a resource for second language learners, or whether (and to what extent) an LM can be trained to select appropriate collocates. Our setting is essentially a sentence-pair classification problem, where the second sentence is the lexical collocation itself. Specifically, a training instance is a tuple <sentence, collocation, label>, as in Example (3) (where we use ; as a wildcard concatenation token, as these are different across LMs):

(3) *The US military **launched an airstrike** on what it described as a safehouse in the Iraqi town of Fallujah ; launched an airstrike → **Oper1***

We use as labels the LFs listed in Table 2, with their respective training/test splits, and train all LMs with the same hyperparameters.[6] The considered LMs are BERT (base and large, uncased) (Devlin et al., 2018), RoBERTa (base and large) (Liu et al., 2019), DistilBERT (Sanh et al., 2019), ALBERT (Sanh et al., 2020) and XLNet (base and large) (Yang et al., 2019). We use the implementation in the Transformers Python library (Wolf et al., 2020).[7]

## 5.1 Results and discussion

The results of this experiment (cf. Table 5) clearly highlight what has already pointed out in Shwartz and Dagan (2019): the prototypical LVCs (as modeled by Oper1) can be identified with a rather high quality. Interesting enough, this is not true for LVCs captured by Oper2, whose only difference to Oper1 is the subcategorization frame: while in Oper1, it is the 1st argument of the base that is realized as the grammatical subject, in Oper2, it is the 2nd argument. Frequency cannot explain this discrepancy since, e.g., IncepOper1, which appears

in our corpus in nearly the same number of sentences as Oper2, is categorized with a significantly higher quality.

Results are also lower for some other verbal LFs with more semantic load, among them, e.g., Real1/2 and Caus1Func0, suggesting that the semantics expressed in the notions of 'realize' and 'cause', especially when the 2nd argument of the collocation functions as a subject, are more challenging. Again, these results cannot be fully explained by the amount of training data and neither by the semantic load. Thus, the categorization of IncepPredPlus achieves the highest score (the best model on IncepPredPlus obtains an average F1 of 95.21 with little variability across runs), and it is clearly an LF with a semantic load, namely 'increase'. Interestingly, Ver ('genuine') and Bon ('positive') are the worst categorized LFs in our sample, while their antonyms AntiVer and Anti-Bon are categorized considerably better.

As for the considered LMs, the best overall performing model is the RoBERTa family, with an overall F1 score of **71.19%** for RoBERTa-base and 70.6% for RoBERTa-large, and both models accounting for the best results on 7 of the 16 target LFs. The second best results are achieved by XLNet (base and large), with XLNet-large being the best model on both Magn and AntiMagn, two LFs which have been traditionally challenging to tell apart due to the fact that the representations of antonyms are clustered together in distributional spaces. We also note that, interestingly, DistilBERT is the best at categorizing Oper1 and Real2, which may suggest that small models may be sufficient to obtain good performnace on categorizing LVCs.

In order to gain further insights on *why* a LM may err in the task of in-context collocation categorization, we display a confusion matrix obtained from random runs for the two LMs with the highest avg score for Oper1 (Distilbert) and Oper2 (XLNet-base) (Figure 2) – the two LVC LFs that differ only in terms of their subcategorization patterns (cf. above). We may hypothesize that the categorization of a collocation based mainly on its actantial structure is challenging, and indeed, we observe that for these two models,[8] syntax-based categorization over the same semantics proves hard. Specifically, XLNet-base has as the greatest source for confusion regarding Oper1, precisely, Oper2;

---

[6]All models are trained for 1 epoch, with a learning rate of $4e-05$, the Adam with weight decay optimizer and a warmup ratio of 0.06.

[7]https://huggingface.co/transformers/.

[8]As a matter of fact, it can be assumed that this applies to all the LMs evaluated in this paper.

and this also occurs with Real1 vs. Real2 (which also differ only with respect to their subcategorization pattern). The results for DistilBERT show a greater spread among the misclassifications of Oper1, namely across Oper2, Caus1Func0 and IncepOper1, and for Oper2 across Oper1 and CausFunc0. Caus1Func0 and IncepOper1 have the same subcategorization pattern as Oper1 (but different semantics). In the case of CausFunc0 ('cause the existence', e.g., CausFunc0(*hope*) = *raise*), the subcategorization pattern is very similar to Caus1Func0, only that the grammatical subject of the corresponding syntactic construction is not an argument of the base. As we can observe, CausFunc0 is easily miscategorized as a full LVC. Finally, let us highlight the fact that while Magn is generally well categorized, the few misclassifications come, as would be expected, from collocations which convey a similar notion of amplification (e.g., Bon), but interestingly, also collocations that convey opposite semantics, such as AntiMagn or AntiBon.



Figure 2: Confusion matrix for the two best performing models on average on Oper2 (Xlnet-base, left) and Oper1 (DistilBERT, right).

## 6 Subspace Analysis

In this section, we further explore the semantics of some selected LFs. We generate visualizations of PCA-projected BERT vectors for all collocation mentions of Magn, AntiMagn, Oper1 and Oper2. These four LFs are sufficiently frequent, and they encode different morphosyntactic structures.[9] We can see that antonymy (Ono et al., 2015; Schwartz et al., 2015; Nguyen et al., 2016) is relatively well captured in contextualized models, although the subspaces are clearly different between the embedding and the last transformer layer. More specifically, as the representations of collocates for Magn

[9]<Magn,AntiMagn> are most frequently expressed by *adj+noun* combinations, whereas <Oper1,Oper2> are always realized by a *verb+noun* pattern.



Figure 3: Oper1 (red) and Oper2 (blue) collocate embeddings for BERT's embedding layer (top row, left), and for the 1st (top row, right), and 5th and 12th transformer layers (second row, left and right, respectively). The bottom quadrant corresponds to Magn (blue) vs AntiMagn (red), with the same arrangements (embedding, 1st, 5th and 12th layer).

and AntiMagn undergo the self-attention-based transformations through BERT's layers, many of these contextualized embeddings tend to group in a narrow cone, with many antonymic collocates indistinguishably overlapping with each other. Similarly, we also observe a tendency of representation overlap in the Oper1 vs Oper2 case, with the embeddings in the last transformer layer showing a cluttered distribution, suggesting that there is little inherent knowledge in BERT to categorize a collocation into the syntactic typification of a LF.

| | BERT-base | | | BERT-large | | | Albert | | | DistilBERT | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Precision | Recall | F1 | Precision | Recall | F1 | Precision | Recall | F1 | Precision | Recall | F1 |
| AntiBon | $67.92_{\pm5.99}$ | $72.01_{\pm4.03}$ | $69.85_{\pm4.67}$ | $71.73_{\pm11.50}$ | $72.43_{\pm3.03}$ | $71.81_{\pm6.29}$ | $77.00_{\pm3.90}$ | $75.31_{\pm1.66}$ | $76.11_{\pm2.21}$ | $69.89_{\pm5.13}$ | $70.90_{\pm4.55}$ | $70.35_{\pm4.42}$ |
| AntiMagn | $85.59_{\pm3.30}$ | $78.22_{\pm4.83}$ | $81.60_{\pm1.14}$ | $85.05_{\pm2.93}$ | $74.93_{\pm3.54}$ | $79.66_{\pm3.26}$ | $84.55_{\pm1.90}$ | $74.83_{\pm0.69}$ | $79.38_{\pm0.58}$ | $85.29_{\pm3.56}$ | $75.55_{\pm2.09}$ | $80.10_{\pm2.38}$ |
| AntiReal2 | $81.67_{\pm1.48}$ | $78.50_{\pm2.29}$ | $80.02_{\pm0.56}$ | $84.74_{\pm3.38}$ | $78.50_{\pm3.50}$ | $81.40_{\pm0.77}$ | $81.39_{\pm12.70}$ | $82.50_{\pm4.33}$ | $81.58_{\pm6.90}$ | $71.23_{\pm2.33}$ | $79.16_{\pm1.89}$ | $74.98_{\pm2.02}$ |
| AntiVer | $61.32_{\pm4.03}$ | $72.51_{\pm4.82}$ | $66.28_{\pm1.91}$ | $61.21_{\pm2.44}$ | $71.76_{\pm6.26}$ | $66.01_{\pm3.67}$ | $62.35_{\pm3.00}$ | $79.25_{\pm1.28}$ | $69.78_{\pm2.37}$ | $63.33_{\pm2.20}$ | $73.33_{\pm10.40}$ | $67.63_{\pm3.18}$ |
| Bon | $50.95_{\pm4.34}$ | $43.06_{\pm3.96}$ | $46.54_{\pm2.64}$ | $53.52_{\pm9.98}$ | $56.71_{\pm4.09}$ | $54.91_{\pm6.98}$ | $55.23_{\pm7.97}$ | $49.52_{\pm8.27}$ | $52.19_{\pm8.06}$ | $54.52_{\pm13.50}$ | $51.64_{\pm10.80}$ | **$53.02_{\pm12.10}$** |
| Caus1Func0 | $67.18_{\pm5.53}$ | $46.97_{\pm15.10}$ | **$53.77_{\pm10.10}$** | $51.49_{\pm7.12}$ | $41.83_{\pm9.77}$ | $45.47_{\pm5.95}$ | $59.19_{\pm3.67}$ | $43.62_{\pm7.56}$ | $50.13_{\pm6.25}$ | $64.21_{\pm10.60}$ | $42.05_{\pm8.06}$ | $50.51_{\pm8.25}$ |
| CausFunc0 | $72.37_{\pm10.20}$ | $59.30_{\pm7.90}$ | $64.41_{\pm2.97}$ | $71.10_{\pm7.40}$ | $58.21_{\pm6.02}$ | $63.75_{\pm4.08}$ | $78.63_{\pm7.92}$ | $60.39_{\pm5.61}$ | $68.07_{\pm3.99}$ | $74.12_{\pm2.79}$ | $65.57_{\pm1.59}$ | $69.56_{\pm1.52}$ |
| IncepOper1 | $79.30_{\pm3.50}$ | $76.43_{\pm1.95}$ | $77.82_{\pm2.28}$ | $78.44_{\pm7.49}$ | $75.13_{\pm5.46}$ | $76.40_{\pm1.19}$ | $81.60_{\pm4.03}$ | $76.87_{\pm1.31}$ | $79.14_{\pm2.53}$ | $79.61_{\pm1.22}$ | $77.55_{\pm4.01}$ | $78.53_{\pm2.39}$ |
| IncepPredPlus | $97.21_{\pm3.03}$ | $91.06_{\pm0.41}$ | $94.02_{\pm1.63}$ | $93.85_{\pm6.16}$ | $90.57_{\pm1.91}$ | $92.08_{\pm2.32}$ | $97.73_{\pm2.21}$ | $85.99_{\pm6.69}$ | $91.43_{\pm4.68}$ | $98.46_{\pm2.01}$ | $90.57_{\pm0.72}$ | $94.34_{\pm0.68}$ |
| LiquFunc0 | $88.02_{\pm6.34}$ | $91.02_{\pm2.36}$ | $89.43_{\pm3.96}$ | $88.35_{\pm3.32}$ | $87.54_{\pm2.75}$ | $87.94_{\pm2.94}$ | $88.70_{\pm2.86}$ | $91.80_{\pm5.57}$ | $90.10_{\pm1.22}$ | $85.72_{\pm2.55}$ | $87.76_{\pm1.40}$ | $86.70_{\pm0.93}$ |
| Magn | $92.47_{\pm0.75}$ | $93.91_{\pm0.77}$ | $93.18_{\pm0.32}$ | $92.87_{\pm0.80}$ | $93.67_{\pm2.23}$ | $93.26_{\pm1.50}$ | $92.74_{\pm0.52}$ | $93.78_{\pm0.93}$ | $93.26_{\pm0.62}$ | $92.80_{\pm1.16}$ | $94.58_{\pm0.32}$ | $93.68_{\pm0.74}$ |
| Oper1 | $78.89_{\pm1.08}$ | $91.69_{\pm0.92}$ | $84.80_{\pm0.27}$ | $78.01_{\pm1.94}$ | $90.26_{\pm1.54}$ | $83.69_{\pm1.76}$ | $78.69_{\pm0.36}$ | $92.77_{\pm2.82}$ | $85.14_{\pm1.18}$ | $79.58_{\pm1.50}$ | $91.89_{\pm1.94}$ | **$85.28_{\pm1.27}$** |
| Oper2 | $70.16_{\pm2.87}$ | $48.07_{\pm6.46}$ | $56.79_{\pm3.97}$ | $66.19_{\pm1.82}$ | $52.17_{\pm5.14}$ | $58.24_{\pm3.04}$ | $71.03_{\pm9.65}$ | $51.98_{\pm2.25}$ | $59.85_{\pm3.83}$ | $67.82_{\pm4.49}$ | $48.01_{\pm2.22}$ | $56.17_{\pm2.34}$ |
| Real1 | $70.03_{\pm2.29}$ | $59.71_{\pm4.38}$ | $64.44_{\pm3.48}$ | $69.91_{\pm2.31}$ | $59.21_{\pm2.24}$ | $64.12_{\pm2.24}$ | $70.52_{\pm4.10}$ | $62.18_{\pm1.63}$ | $66.02_{\pm1.60}$ | $70.81_{\pm8.46}$ | $62.44_{\pm4.33}$ | $66.25_{\pm5.57}$ |
| Real2 | $58.36_{\pm2.86}$ | $53.41_{\pm7.12}$ | $55.48_{\pm3.12}$ | $64.99_{\pm5.87}$ | $54.79_{\pm5.44}$ | $59.13_{\pm1.98}$ | $63.52_{\pm2.53}$ | $47.47_{\pm8.75}$ | $54.14_{\pm6.64}$ | $69.69_{\pm0.79}$ | $54.30_{\pm12.80}$ | $60.51_{\pm8.24}$ |
| Ver | $40.15_{\pm7.96}$ | $24.31_{\pm3.82}$ | $30.23_{\pm4.92}$ | $37.59_{\pm12.30}$ | $22.27_{\pm2.90}$ | $27.38_{\pm4.40}$ | $31.54_{\pm16.10}$ | $22.61_{\pm11.60}$ | $25.95_{\pm13.40}$ | $43.15_{\pm25.20}$ | $20.91_{\pm9.93}$ | $28.07_{\pm14.50}$ |
| Average | $72.60_{\pm4.10}$ | $67.51_{\pm4.45}$ | $69.29_{\pm3.00}$ | $71.82_{\pm5.42}$ | $67.50_{\pm4.11}$ | $69.08_{\pm3.22}$ | $73.40_{\pm5.21}$ | $68.18_{\pm4.43}$ | $70.14_{\pm4.13}$ | $73.14_{\pm5.47}$ | $67.89_{\pm4.82}$ | $69.73_{\pm4.41}$ |

| | XLNet-base | | | XLNet-large | | | RoBERTa-base | | | RoBERTa-large | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Precision | Recall | F1 | Precision | Recall | F1 | Precision | Recall | F1 | Precision | Recall | F1 |
| AntiBon | $67.74_{\pm9.42}$ | $74.72_{\pm2.49}$ | $70.88_{\pm6.14}$ | $70.30_{\pm6.48}$ | $76.75_{\pm3.99}$ | $73.25_{\pm4.11}$ | $71.57_{\pm11.90}$ | $75.82_{\pm9.03}$ | $73.51_{\pm10.20}$ | $77.44_{\pm10.10}$ | $76.42_{\pm10.10}$ | **$76.82_{\pm10.00}$** |
| AntiMagn | $85.26_{\pm5.43}$ | $76.52_{\pm2.05}$ | $80.62_{\pm3.36}$ | $86.18_{\pm5.10}$ | $77.86_{\pm4.67}$ | **$81.78_{\pm4.62}$** | $90.02_{\pm4.44}$ | $77.19_{\pm3.90}$ | $83.10_{\pm3.90}$ | $84.52_{\pm9.63}$ | $75.86_{\pm9.63}$ | $79.84_{\pm5.16}$ |
| AntiReal2 | $81.96_{\pm3.72}$ | $77.83_{\pm2.92}$ | $79.83_{\pm3.16}$ | $85.06_{\pm10.40}$ | $78.83_{\pm1.60}$ | **$81.65_{\pm5.21}$** | $81.72_{\pm12.30}$ | $80.00_{\pm2.78}$ | $80.61_{\pm6.89}$ | $80.64_{\pm3.95}$ | $82.66_{\pm3.95}$ | $81.64_{\pm4.08}$ |
| AntiVer | $65.29_{\pm2.14}$ | $67.57_{\pm1.99}$ | $66.37_{\pm0.36}$ | $67.60_{\pm3.54}$ | $70.20_{\pm6.34}$ | $68.86_{\pm4.91}$ | $69.37_{\pm1.08}$ | $75.39_{\pm9.19}$ | **$72.11_{\pm4.89}$** | $69.93_{\pm9.95}$ | $72.92_{\pm9.95}$ | $71.18_{\pm5.85}$ |
| Bon | $51.70_{\pm6.46}$ | $47.40_{\pm2.16}$ | $49.41_{\pm4.10}$ | $52.08_{\pm15.90}$ | $47.51_{\pm3.42}$ | $49.24_{\pm8.88}$ | $53.95_{\pm7.90}$ | $46.87_{\pm2.06}$ | $50.01_{\pm4.13}$ | $50.84_{\pm7.58}$ | $54.70_{\pm7.58}$ | $52.70_{\pm7.85}$ |
| Caus1Func0 | $64.71_{\pm9.52}$ | $44.29_{\pm15.40}$ | $50.52_{\pm9.85}$ | $64.46_{\pm24.90}$ | $40.93_{\pm11.70}$ | $47.48_{\pm5.79}$ | $68.68_{\pm12.20}$ | $47.87_{\pm22.40}$ | $52.47_{\pm15.40}$ | $59.92_{\pm11.70}$ | $42.05_{\pm11.70}$ | $49.41_{\pm10.20}$ |
| CausFunc0 | $75.65_{\pm9.76}$ | $59.58_{\pm7.49}$ | $66.35_{\pm6.80}$ | $79.15_{\pm14.70}$ | $60.39_{\pm2.38}$ | $68.17_{\pm5.83}$ | $78.06_{\pm6.82}$ | $67.21_{\pm2.45}$ | **$72.15_{\pm3.85}$** | $75.56_{\pm8.38}$ | $63.30_{\pm8.38}$ | $68.81_{\pm6.47}$ |
| IncepOper1 | $73.37_{\pm3.87}$ | $78.16_{\pm1.08}$ | $75.65_{\pm1.94}$ | $79.35_{\pm4.77}$ | $75.20_{\pm2.35}$ | $77.13_{\pm1.86}$ | $79.38_{\pm4.43}$ | $79.03_{\pm1.76}$ | $79.16_{\pm2.54}$ | $83.29_{\pm8.40}$ | $76.49_{\pm8.40}$ | **$79.61_{\pm4.24}$** |
| IncepPredPlus | $99.21_{\pm0.01}$ | $91.54_{\pm2.32}$ | **$95.21_{\pm1.26}$** | $96.31_{\pm4.49}$ | $90.33_{\pm0.83}$ | $93.18_{\pm1.71}$ | $96.31_{\pm1.49}$ | $91.54_{\pm1.20}$ | $93.43_{\pm1.49}$ | $92.79_{\pm12.40}$ | $91.06_{\pm12.40}$ | $91.58_{\pm6.04}$ |
| LiquFunc0 | $87.38_{\pm2.63}$ | $92.48_{\pm4.18}$ | $89.83_{\pm2.63}$ | $87.95_{\pm4.60}$ | $91.91_{\pm6.18}$ | $89.68_{\pm1.09}$ | $88.75_{\pm0.35}$ | $93.04_{\pm4.01}$ | $90.82_{\pm1.89}$ | $90.15_{\pm2.41}$ | $91.58_{\pm2.41}$ | **$90.86_{\pm2.81}$** |
| Magn | $91.93_{\pm0.85}$ | $93.91_{\pm0.48}$ | $92.91_{\pm0.54}$ | $93.41_{\pm1.66}$ | $94.24_{\pm1.66}$ | **$93.82_{\pm1.62}$** | $92.57_{\pm1.31}$ | $94.95_{\pm1.60}$ | $93.74_{\pm1.43}$ | $92.90_{\pm1.29}$ | $94.73_{\pm1.29}$ | $93.81_{\pm1.68}$ |
| Oper1 | $77.80_{\pm0.97}$ | $91.09_{\pm0.20}$ | $83.92_{\pm0.48}$ | $78.90_{\pm0.94}$ | $91.80_{\pm0.78}$ | $84.86_{\pm0.22}$ | $78.33_{\pm2.36}$ | $91.18_{\pm1.78}$ | $84.25_{\pm1.46}$ | $78.41_{\pm1.85}$ | $92.40_{\pm1.85}$ | $84.81_{\pm1.79}$ |
| Oper2 | $72.13_{\pm8.13}$ | $52.69_{\pm3.66}$ | **$60.85_{\pm5.12}$** | $74.05_{\pm4.18}$ | $50.70_{\pm7.88}$ | $59.86_{\pm4.83}$ | $72.05_{\pm1.20}$ | $50.70_{\pm8.57}$ | $59.18_{\pm5.79}$ | $70.19_{\pm5.83}$ | $52.56_{\pm5.83}$ | $59.97_{\pm1.06}$ |
| Real1 | $72.00_{\pm1.43}$ | $59.84_{\pm0.99}$ | $65.36_{\pm0.96}$ | $66.09_{\pm3.48}$ | $64.60_{\pm4.19}$ | $65.33_{\pm3.82}$ | $71.11_{\pm4.52}$ | $61.05_{\pm4.51}$ | $65.66_{\pm4.10}$ | $72.59_{\pm4.37}$ | $62.49_{\pm4.37}$ | **$66.90_{\pm5.70}$** |
| Real2 | $66.37_{\pm1.17}$ | $51.23_{\pm4.45}$ | $57.77_{\pm3.13}$ | $64.01_{\pm12.30}$ | $54.20_{\pm7.40}$ | $58.67_{\pm9.52}$ | $66.12_{\pm3.52}$ | $52.91_{\pm8.66}$ | $58.67_{\pm6.65}$ | $64.58_{\pm7.64}$ | $51.03_{\pm7.64}$ | $56.92_{\pm9.77}$ |
| Ver | $32.58_{\pm12.90}$ | $15.98_{\pm10.50}$ | $21.37_{\pm12.90}$ | $32.94_{\pm19.20}$ | $26.70_{\pm22.00}$ | $29.21_{\pm21.10}$ | $48.10_{\pm29.10}$ | $22.78_{\pm5.97}$ | **$30.38_{\pm11.30}$** | $41.46_{\pm4.40}$ | $18.02_{\pm4.40}$ | $24.77_{\pm3.23}$ |
| Average | $72.82_{\pm4.90}$ | $67.18_{\pm3.90}$ | $69.18_{\pm3.92}$ | $73.71_{\pm8.38}$ | $68.20_{\pm5.55}$ | $70.15_{\pm5.33}$ | $75.38_{\pm6.75}$ | $69.15_{\pm5.59}$ | **$71.19_{\pm5.38}$** | $74.08_{\pm6.87}$ | $68.64_{\pm6.87}$ | $70.60_{\pm5.37}$ |

Table 5: Average Precision, Recall and F1 results for the collocation classification experiment, computed by averaging the results of three independent runs. We also report standard deviation figures. Results are provided per LF as well as the average over each metric (Average).

## 7 Conclusions

We have analyzed LMs in tasks revolving around modeling, recognizing and categorizing lexical collocations. We conclude that some prominet types of LVCs require little context to be well encoded, as opposed to other LFs involving, e.g., nouns and adjectives, and that predictability of LFs is challenging, not a function of training data, and that syntax plays a major role.

## 8 Future Work

In the future, we will make this work multilingual using linguistic equivalences as anchors, in the spirit of cross-lingual embedding research, in order to align collocations of the same LF across languages (e.g., in English and Norwegian we *take a nap*, in German, we 'make' it, in Portuguese we 'pull' it, in Spanish, we 'throw' it, etc.). We would also like to explore the idea of "semantic masking" for collocate discovery, where we would train models for dynamically *masking* (or removing) idiosyncratic information such that only the semantics of the collocate remain, thus largely corresponding to a latent abstraction over the LF. This approach has been applied recently in the lexical substitution task, with the limitation, however, that the dropout rate was tuned in a validation set, whereas a promising avenue to explore would be to automatically learn the embedding dropout in a fully supervised setting. Finally, motivated by the observed large gap in performance between the categorization of, e.g., Oper1 and Oper2, Bon and AntiBon, Ver and AntiVer, we plan to investigate in more depth the codification of collocational information in pretrained LMs.

## Acknowledgements

# References

Marco Baroni and Roberto Zamparelli. 2010. Nouns are vectors, adjectives are matrices: Representing adjective-noun constructions in semantic space. In *Proceedings of the 2010 conference on empirical methods in natural language processing*, pages 1183–1193.

Gemma Boleda, Marco Baroni, Louise McNally, et al. 2013. Intensionality was only alleged: On adjective-noun composition in distributional semantics. In *Proceedings of the 10th International Conference on Computational Semantics (IWCS 2013): long papers; 2013 Mar 20-22; Postdam, Germany. Stroudsburg (USA): Association for Computational Linguistics (ACL); 2013. p. 35-46.* ACL (Association for Computational Linguistics).

Jose Camacho-Collados, Luis Espinosa-Anke, Shoaib Jameel, and Steven Schockaert. 2019. A latent variable model for learning distributional relation vectors. In *International Joint Conferences on Artificial Intelligence*.

Kenneth W. Church and Patrick Hanks. 1989. Word Association Norms, Mutual Information, and Lexicography. In *Proceedings of the 27th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 76–83, Vancouver, Canada.

Anthony P. Cowie. 1994. Phraseology. In R.E. Asher and J.M.Y. Simpson, editors, *The Encyclopedia of Language and Linguistics, Vol. 6*, pages 3168–3171. Pergamon, Oxford.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Daniel Edmiston. 2020. A systematic analysis of morphological content in bert models for multiple languages. *arXiv preprint arXiv:2004.03032v1*.

Luis Espinosa-Anke, Jose Camacho-Collados, Sara Rodríguez Fernández, Horacio Saggion, and Leo Wanner. 2016. Extending wordnet with fine-grained collocational information via supervised distributional learning. In *Proceedings of COLING 2016: Technical Papers. The 26th International Conference on Computational Linguistics; 2016 Dec. 11-16; Osaka (Japan): COLING; 2016. p. 900-10.* COLING.

Luis Espinosa-Anke and Steven Schockaert. 2018. Seven: Augmenting word embeddings with unsupervised relation vectors. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2653–2665.

Luis Espinosa-Anke, Steven Schockaert, and Leo Wanner. 2019. Collocation classification with unsupervised relation vectors. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5765–5772.

Stefan Evert. 2007. Corpora and Collocations. In A. Lüdeling and M. Kytö, editors, *Corpus Linguistics. An International Handbook.* Mouton de Gruyter, Berlin.

John R. Firth. 1957. Modes of Meaning. In J.R. Firth, editor, *Papers in Linguistics, 1934-1951*, pages 190–215. Oxford University Press, Oxford.

Marcos Garcia, Marcos García Salido, and Margarita Alonso Ramos. 2017. Using bilingual word-embeddings for multilingual collocation extraction. In *Proceedings of the 13th Workshop on Multiword Expressions (MWE 2017)*, pages 21–30.

A. Gelbukh and O. Kolesnikova. 2012. *Semantic Analysis of Verbal Collocations with Lexical Functions*. Springer, Heidelberg.

Franz Josef Hausmann. 1985. Kollokationen im Deutschen Woerterbuch: ein Beitrag zur Theorie des lexicographischen Biespiels. *Lexikographie und Grammatik*.

John Hewitt and Christopher D Manning. 2019. A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4129–4138.

Felix Hill, Kyunghyun Cho, Anna Korhonen, and Yoshua Bengio. 2016. Learning to understand phrases by embedding the dictionary. *Transactions of the Association for Computational Linguistics*, 4:17–30.

Mandar Joshi, Eunsol Choi, Omer Levy, Daniel S Weld, and Luke Zettlemoyer. 2019. pair2vec: Compositional word-pair embeddings for cross-sentence inference. In *NAACL-HLT (1)*.

Mandar Joshi, Kenton Lee, Yi Luan, and Kristina Toutanova. 2020. Contextualized representations using textual encyclopedic knowledge. *arXiv preprint arXiv:2004.12006*.

Adam Kilgarriff. 2006. Collocationality (And How to Measure it). In *Proceedings of the 12th Euralex International Congress on Lexicography (EURALEX)*, pages 997–1004, Turin, Italy. Springer-Verlag.

Brigitte Krenn. 2000. *The Usual Suspects: Data-Oriented Models for the Identification and Representation of Lexical Collocations*, volume 7 of *Saarbrücken Dissertations in Computational Linguistics and Language Technology*. DFKI and Universität des Saarlandes.

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations.

Omer Levy, Steffen Remus, Chris Biemann, Ido Dagan, and Israel Ramat-Gan. 2015. Do supervised distributional methods really learn lexical inference relations? In *NAACL 2015*, Denver, Colorado, USA. Association for Computational Linguistics.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Marco Maru, Federico Scozzafava, Federico Martelli, and Roberto Navigli. 2019. Syntagnet: challenging supervised word sense disambiguation with lexical-semantic combinations. 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 3525–3531.

I.A. Mel'čuk. 1996. Lexical functions: A tool for the description of lexical relations in the lexicon. In Leo Wanner, editor, *Lexical Functions in Lexicography and Natural Language Processing*, pages 37–102. Benjamins Academic Publishers, Amsterdam.

Igor Mel'čuk. 1995. Phrasemes in Language and Phraseology in Linguistics. In M. Everaert, E.-J. van der Linden, A. Schenk, and R. Schreuder, editors, *Idioms: Structural and Psychological Perspectives*, pages 167–232. Lawrence Erlbaum Associates, Hillsdale.

Jeff Mitchell and Mirella Lapata. 2010. Composition in distributional models of semantics. *Cognitive science*, 34(8):1388–1429.

Kim Anh Nguyen, Sabine Schulte im Walde, and Ngoc Thang Vu. 2016. Integrating distributional lexical contrast into word embeddings for antonym-synonym distinction. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 454–459.

Masataka Ono, Makoto Miwa, and Yutaka Sasaki. 2015. Word embedding-based antonym detection using thesauri and distributional information. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 984–989.

Pavel Pecina. 2008. A Machine Learning Approach to Multiword Expression Extraction. In *Proceedings of the LREC 2008 Workshop Towards a Shared Task for Multiword Expressions)*, pages 54–57, Marrakech, Morocco.

Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, 12.

Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? 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 2463–2473.

Jipeng Qiang, Yun Li, Yi Zhu, and Yunhao Yuan. 2019. A simple bert-based approach for lexical simplification. *arXiv preprint arXiv:1907.06226*.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9.

Carlos Ramisch, Silvio Cordeiro, Agata Savary, Veronika Vincze, Verginica Mititelu, Archna Bhatia, Maja Buljan, Marie Candito, Polona Gantar, Voula Giouli, et al. 2018. Edition 1.1 of the parseme shared task on automatic identification of verbal multiword expressions.

Sara Rodríguez Fernández, Luis Espinosa-Anke, Roberto Carlini, and Leo Wanner. 2016. Semantics-driven recognition of collocations using word embeddings. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics; 2016 Aug. 7-12; Berlin (Germany).[place unknown]: ACL; 2016. Vol. 2, Short Papers; p. 499-505.* ACL (Association for Computational Linguistics).

Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in bertology: What we know about how bert works. *arXiv preprint arXiv:2002.12327*.

V. Sanh, L. Debut, J. Chaumond, and Th. Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. In *Proceedings of the 5th Workshop on Energy Efficient Machine Learning and Cognitive Computing*, Vancouver, BC, Canada.

V. Sanh, L. Debut, J. Chaumond, and Th. Wolf. 2020. ALBERT: A lite BERT for self-supervised learning of language representations. In *Proceedings of ICLR*.

Roy Schwartz, Roi Reichart, and Ari Rappoport. 2015. Symmetric pattern based word embeddings for improved word similarity prediction. In *Proceedings of the nineteenth conference on computational natural language learning*, pages 258–267.

V. Seretan. 2014. On collocations and their interaction with parsing and translation. *Informatics*, 1(1):11–31.

Vered Shwartz and Ido Dagan. 2019. Still a pain in the neck: Evaluating text representations on lexical composition. *Transactions of the Association for Computational Linguistics*, 7:403–419.

Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019a. SuperGLUE: A stickier benchmark for general-purpose language understanding systems. In *Proceedings of NeurIPS*.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019b. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of ICLR*.

Leo Wanner and John A. Bateman. 1990. A collocational based approach to salience sensitive lexical selection. In *Proceedings of the 5th International Workshop on Natural Language Generation*, Dawson, PA.

Leo Wanner, Bernd Bohnet, and Mark Giereth. 2006. Making sense of collocations. *Computer Speech and Language*, 20(4):609–624.

Koki Washio and Tsuneaki Kato. 2018. Neural latent relational analysis to capture lexical semantic relations in a vector space. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 594–600.

Thomas Wolf, Julien Chaumond, Lysandre Debut, Victor Sanh, Clement Delangue, Anthony Moi, Pierric Cistac, Morgan Funtowicz, Joe Davison, Sam Shleifer, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45.

Xing Wu, Shangwen Lv, Liangjun Zang, Jizhong Han, and Songlin Hu. 2019. Conditional bert contextual augmentation. In *International Conference on Computational Science*, pages 84–95. Springer.

Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. Salakhutdinov, and Q.V. Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *Proceedings of the 33rd Conference on Neural Information Processing Systems (NeurIPS 2019)*, Vancouver, BC, Canada.

Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. Swag: A large-scale adversarial dataset for grounded commonsense inference. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 93–104.

Wangchunshu Zhou, Tao Ge, Ke Xu, Furu Wei, and Ming Zhou. 2019. Bert-based lexical substitution. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3368–3373.