# Generalization through Lexical Abstraction in Transformer Models: the Case of Functional Words

**Anonymous ACL submission**

## Abstract

"The researchers wrote the paper" and "They wrote it" share syntactic and semantic information that is easily recognizable for humans. Specifically, the latter is an abstraction of the former. Can language models also recognize the syntactic and semantic parallelism of the two sentences, which relies on lexical abstraction? We present a study that aims to uncover whether a language model encodes words and sentences in a way that reflects this linguistic abstraction.

We compare representations of nouns, on one side, and the pronouns and adverbs (functional words) that can replace these nouns, as well as the corresponding lexicalized and functional sentences, on the other. The shallow analyses show that nouns and functional words inhabit different areas of the embedding space, both when considered in isolation or in the same sentential contexts. Deeper analyses, however, show that the structure shared between the lexicalized sentences and their functional variations is encoded and can be uncovered from their embeddings.

Our results then indicate that, when properly constrained by the structure, the information supporting the generalization through abstraction provided by pronouns and functional words can be revealed.

## 1 Introduction

Large language models (LLMs) are very successful, and much of their success stems from their ability to induce word or token representations that encode the extremely complex language data, with many generative factors (Bengio et al., 2013). It is an ongoing quest to understand these representations, the kind of linguistic information they encode and the way a system is able to successfully manipulate them to solve a wide variety of tasks. It is difficult to attribute their high performance on numerous linguistic and NLP tasks to their understanding of language and its structure (Waldis et al., 2024). One of the criteria for judging the degree of language understanding in LLMs is their capacity to "generalize" well. This question is often approached from a technical, rather than a linguistic, perspective. Generalization is considered a crucial property of a learned model, as it ensures trust in its deployment outside of its training environment – whether this application involves a slightly different task, out-of-distribution data, a different language, or some other level of distinction between the application domain and the one it was trained on (Hupkes et al., 2023). This point of view often involves learning a probe on top of the pretrained model.

There are however other types of generalization, namely linguistic generalizations and abstractions. For example, speakers can easily strip down a sentence to a basic syntactic-semantic structure, such as *Who did what to whom* or *She put that there* or *She does that sometimes.* The use of pronouns or adverbs to reduce a sentence to a "skeleton" does not rely on using out-of-vocabulary items, as pronouns and adverbials such as *somewhere/sometime* are some of the most frequent words in a corpus, and appear in many shared contexts, as their frequent use in coreferring expressions attests. In semantics, pronominal forms are usually treated as variables, placeholders for more structured lexical elements within a sentence and thus highly abstract entities (Büring, 2019).

Is this particular property of functional words – as abstract place-holders for nouns and prepositional phrases – captured in LLMs? To explore this question we start with a shallow exploration of the embeddings of functional words and nouns in isolation and in the same contexts. We then move on to deeper analyses, where we use a system to search for the shared syntactic structure of sentences that differ only in the use of nouns and prepositional phrases compared to using pronouns and adverbs.

The shallower analyses show that functional

words and nouns do not inhabit the same regions of the embedding space, even when they appear in the same context. The in-depth analyses, instead, show that accessing deeper information in sentence embeddings reveals shared features encoding information about syntactic structure, whether these are filled by content or functional words.

## 2 Data

To explore linguistic generalization through abstraction, in the special case of functional words, we use a purposefully generated dataset on verb alternations, with structure and lexical variation at multiple levels. Unlike other linguistic phenomena (e.g. agreement rules), where all necessary syntactic elements for the rule are contained within a single sentence, verb alternations require observing at least two related sentences. They show that the same verb can appear in different sentential contexts, with systematically related syntactic-semantic mappings of their arguments. This allows us to study a variety of sentential contexts and their lexical and functional expressions, within the controlled environment of the same verb meaning.

The dataset is generated from a set of verbs belonging to the change-of-state (COS) and object-drop (OD) classes (Levin, 1993). These classes provide an argument structure minimal pair: they share the same syntactic structure - transitive/intransitive alternation - but differ in their argument structure. The object of the transitive verbs belonging to the COS class bears the same semantic role (Patient) as the subject of the intransitive verb (*The artist opens this door/This door opens*). The transitive form of the verb has a causative meaning. In contrast, for OD verbs the subject bears the same semantic role (Agent) in both the transitive and intransitive forms and the verb does not have a causative meaning (*The artist paints this door/The artist paints*) (Levin, 1993; Merlo and Stevenson, 2001).

Moreover, we divide words into lexical and functional. Lexical elements, or content words, are an open class of words with a meaningful content, corresponding to concepts or entities and events in the world. The role of the closed class of function words, instead, is to express grammatical functions. We focus specifically on pronouns, and a subset of adverbs, those that can express temporal and spatial concepts. These function words can be used as general placeholders for nouns and prepositional phrases: for instance, *The researchers wrote the*

*article last week* can also be expressed more abstractly as *They wrote it then.*

### 2.1 Data templates

The dataset comprises instances that follow the Blackbird Language Matrices framework (Merlo, 2023). Each instance is a multiple-choice puzzle and it consists of (i) a rule-generated *context* sequence of sentences that illustrate the encoded phenomenon. The rules are of two types: rules that described the linguistic property under study (verb alternation) and rules that are not related to it (e.g. presence or absence of a prepositional phrase). One sentence that would make the sequence complete is missing, and must be chosen from (ii) an answer set of minimally differing contrastive sentences – one correct, and each of the others violating a sub-rule.

**Context set** The syntax-semantics features of the verb alternation, and their combination rules, lead to the construction of the context set. Specifically, (i) the presence of one or two arguments and their attributes (agents, Ag; patients, Pat) ; (ii) the active (Akt) or passive (Pass) voice of the verb. The phenomenon-external factors include an alternation between a NP introduced (i) by any preposition (e.g., *in an instant*, henceforth p-NP) and (ii) by the preposition by (e.g., *by chance*, by-NP), but not agentive (e.g., *by the artist*, by-Ag/by-Pat), which remains a confounding variable. The OD context minimally differs from the COS in the last sentence of the context: the subject of the intransitive is an Agent, and not the Patient.

**Answer Set** All answers have the same structure: (NP V by-NP) consisting of a verb, two nominal constituents (giving rise to a structure of the type NP V NP) and a preposition (by, or the lack of the preposition) between the verb and the second NP. The candidate answers comprise the correct intransitive form of the alternation followed by a by-NP which satisfy the rules of the BLM, and the contrastive incorrect answers obtained by corrupting some properties of the rules (wrong argument, wrong voice of the verb, lack of preposition, wrong constituent of the PP).[1]

The answer set does not change across verb classes, only the label of the correct answer: the cor-

---

[1]Error types: wrong semantic role on the first constituent is a syntax-semantic mapping error (SSM), wrong last constituent introduced by the preposition *by* WRBY, *and the other errors are labelled according to the type of resulting structure – intransitive,* INTR; *transitive,* TRANS; *passive,* PASS.

rect answer for COS is an error for OD, and vice-versa. The BLM-template (context and answers) for COS and OD are presented in Figure 1.

| COS CONTEXT | | | | | COS ANSWERS | | | |
|---|---|---|---|---|---|---|---|---|
| 1 | Ag | Akt | Pat | p-NP | 1 | Pat Akt by-NP | CORRECT |
| 2 | Ag | Akt | Pat | by-NP | 2 | Ag Akt by-NP | SSM-INT |
| 3 | Pat | Pass | by-Ag | p-NP | 3 | Pat Pass by-Ag | PASS |
| 4 | Ag | Pass | by-Ag | by-NP | 4 | Ag Pass by-Pat | SSM-PASS |
| 5 | Pat | Pass | | p-NP | 5 | Pat Akt Ag | TRANS |
| 6 | Pat | Pass | | by-NP | 6 | Ag Akt Pat | SSM-TRANS |
| 7 | Pat | Akt | | p-NP | 7 | Pat Akt by-Ag | WRBY |
| ? | ??? | | | | 8 | Ag Akt by-Pat | SSM-WRBY |

| OD CONTEXT | | | | | OD ANSWERS | | | |
|---|---|---|---|---|---|---|---|---|
| 1 | Ag | Akt | Pat | p-NP | 1 | Pat Akt by-NP | SSM-INT |
| 2 | Ag | Akt | Pat | by-NP | 2 | Ag Akt by-NP | CORRECT |
| 3 | Pat | Pass | by-Ag | p-NP | 3 | Pat Pass by-Ag | SSM-PASS |
| 4 | Ag | Pass | by-Pat | by-NP | 4 | Ag Pass by-Pat | PASS |
| 5 | Pat | Pass | | p-NP | 5 | Pat Akt Ag | SSM-TRANS |
| 6 | Pat | Pass | | by-NP | 6 | Ag Akt Pat | TRANS |
| 7 | Ag | Akt | | p-NP | 7 | Pat Akt by-Ag | SSM-WRBY |
| ? | ??? | | | | 8 | Ag Akt by-Pat | WRBY |

Figure 1: BLM COS and OD contexts and answers.

## 2.2 Levels of lexical abstraction

To explore generalisation through abstraction, we produce two main variants of the data – a lexicalized one (labelled *Lex*), and a functional one, where functional words replace all content words except the main verb (labelled *Fun*). The lexicalised variant comes in different types (type I, II, III), with varying amounts of lexicalisation, for comparison with the small size inventory of the functional words. The groups are exemplified in Figure 2, together with the generation process presented in the next paragraph. Figure 9 and Figure 10 in the appendix examples for type I data for both verb classes.

## 2.3 Main Dataset

The main dataset is built based on thirty (manually chosen) verbs from each of the two classes discussed in Levin (1993). See Table 2 in the appendix for the full list.

The functional lexicon has been manually selected by the authors to maintain the syntactic and semantic acceptability of the sentences[2]. The lexical alternatives were provided by a masked language model (*bert-base-uncased*, (Devlin et al., 2018)). The models received sentences containing only the masked constituent, the verb and the functional elements. For example, to retrieve the arguments for the verb *break*, two masked templates are used: the patient is masked and the agent is in pronominal form (e.g. *she broke (the/a/some/...)*

---
[2]Following the discussion in Haspelmath (1997), we add elements like *somebody* as pronominal elements.
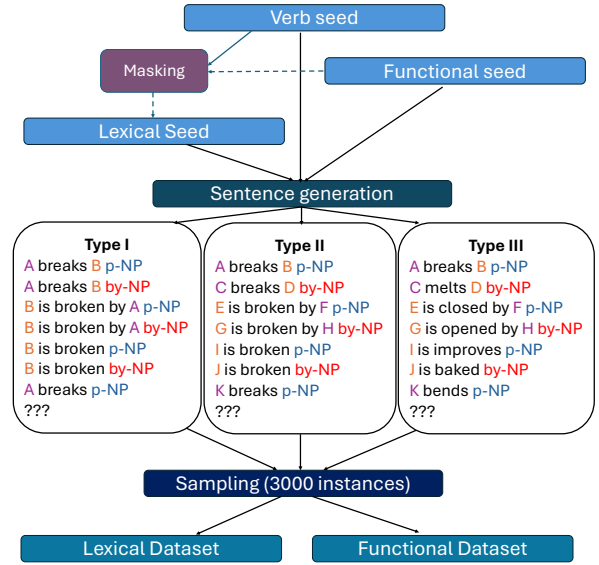


Figure 2: Process of generation of the three levels of lexical variation (type I, II, III), exemplified for COS data. Type I data contains instances with lexically consistent material, with minimal change across the context and the answer set. In type II the verb remains the same while one constituent varies across the context and the answer set. Type III data displays maximal lexical variation in both the context sentences and the answer set.

<MASK>), and the subject of the transitive is masked and the patient is a pronoun *(e.g. (the/a/some/...)* <MASK> *broke it*. Both the lexical seed and the functional seed contain five semantically plausible instances for each constituent class ( Ag, Pat, p-NP and by-NP). We ensured a balanced distribution of tense and number across verbal inflections.

For our experiment, we sampled 3000 instances (out of 38400 combinations of arguments and verbs) for each type, semi-automatically crafted and manually evaluated for plausibility and grammaticality.

## 2.4 Dataset variations

Starting from the main datasets described above, we build several variations that will be used in the different experiments.

**Words** From each sentence in the type I subset of the BLM dataset, we extract the functional words and their corresponding nouns and prepositional phrases. There are 17 functional words and phrases: *he, her, him, it, she, somebody, someone, that, that one, them, these, these ones, they, this, this one, those, those ones* and 204 noun phrases.

**Sentences** We compile parallel versions of the sentences in their lexicalized and functional word

3

forms from the FUN and LEX subsets of the type I BLM dataset. Each sentence has associated its syntactic pattern (the syntactic version of the syntactic-semantic template shown in Figure 1, e.g. *Pron Vpass PP PP*). From these, we sample 4000 sentences, split 80:20 between training and testing, and use 10% of the training data for validation.

**BLM data** Of the thirty verbs, all instances for three of the verbs (3x100) are selected for testing. Of the instances of the other 27 verbs, 2000 are randomly sampled for training. Ten percent of the training data is dynamically selected for validation. The same 27:3 verb split is used for all Fun/Lex and type I/ type II/type III variations. All variations have 2000 instances for training, 300 for testing.

## 3 Analyses and experiments

We aim to understand whether language models encode sentences that we perceive as syntactically and semantically parallel – due to the linguistic abstraction property of pronouns and adverbs relative to nouns and noun phrases – such that this shared information is accessible.

To achieve this, we proceed in several steps. We investigate the relative positions of lexical and functional word embeddings, obtained from isolated words or when presented in similar sentential contexts (Section 3.1). We study the relative positions of the representations of two variations of sentences – with nouns, or with functional words (Section 3.2). We analyse the representation of functional and lexicalized sentences for detecting the shared syntactic structure (Section 3.3). We deploy the BLM linguistic puzzles, whose solution relies on detecting shared structure at the level of input sequence and within each sentence (Section 3.4).

We obtain word and sentence representations (as averaged token embeddings) from an Electra pretrained model (Clark et al., 2020)[3]. We choose Electra because it has been shown to perform better than models from the BERT family on the Holmes benchmark[4], and to also encode information about syntactic and argument structure better (Yi et al., 2022; Nastase and Merlo, 2024).

As a first step of analysis, we use 2D t-SNE projections (Hinton and Roweis, 2002). We project the

---

embeddings of lexical and functional words, when considered in isolation, or within parallel sentential contexts. t-SNE is designed to project high-dimensional data into a lower dimensional space while preserving neighbourhood information. Considering that the embedding space was built based on the notion of similarity and similarity metrics, this type of visualization provides a valid first level of analysis of the properties of the lexicalized and functional word and sentence embeddings.

### 3.1 Word embeddings

**Stand-alone embeddings** Figure 3 shows the t-SNE projection of the word embeddings (as averages over the respective token embeddings) for the functional words and noun phrases in our data, obtained in isolation (when presented to the pretrained model alone). Functional words appear isolated in this space, which indicates that the shared information between the functional elements and the nouns they can replace, should there by any, is not to be found at a shallow level.
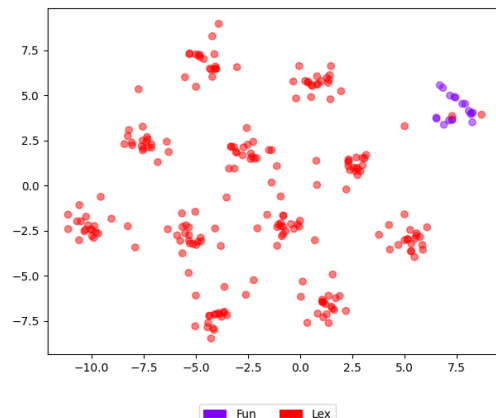


Figure 3: t-SNE projection of the embeddings of functional words and nouns, without a sentential context.

**Contextual word embeddings** We use the parallel versions of the sentences – with content words or functional words – to build contextualized word embeddings, and verify whether the added constraints of belonging in the same sentential contexts brings the word embeddings closer together. Each point in the plot in Figure 4 corresponds to the contextual embedding of a functional word or noun in each of the input sentences. Figure 4 shows that even when embedded in the same context, the embeddings of the functional words remain apart from the embeddings of the nouns.
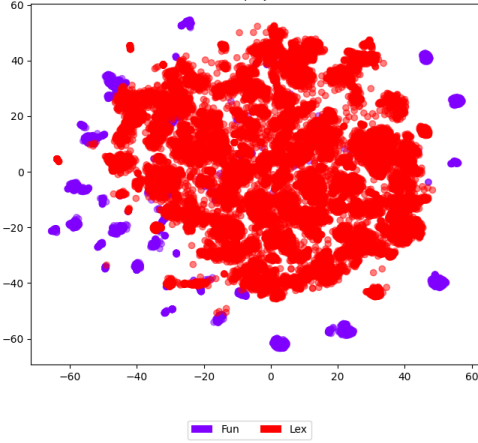
4

Figure 4: tSNE projection of the embeddings of functional words and nouns obtained from parallel contexts. Each point is a contextual embedding.
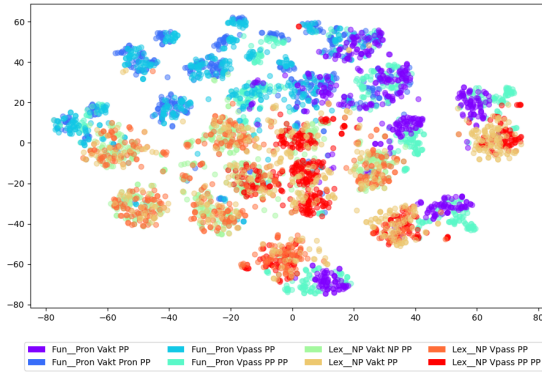


Figure 5: t-SNE projection of sentence representations (averaged token embeddings) coloured by their syntactic pattern and the use of lexicalized or functional words.

### 3.2 Sentence embeddings

Figure 5 shows the t-SNE projection of the representations of the two variations of each sentence. They also occupy different regions of the embedding space, just as the contextualised or out-of-context word embeddings.

### 3.3 Shared structure

The projections of the word and sentence embeddings show that the functional words and the nouns inhabit different regions of the embedding space. The distinctions we observe in these analyses, however, may be only superficial. According to the principle of superposition (Bengio et al., 2013; Elhage et al., 2022), each dimension can contribute to several features, and a feature may be encoded by a combination of dimensions. It is however difficult to define what features are, and how they are encoded in a deep learning model.

We mine for information about the structure of the sentences: these are our shared "features". To

| test on / train on | Fun | Lex |
|---|---|---|
| Fun | 1.000 | 0.441 |
| Lex | 0.493 | 0.990 |
| Mixed | 0.995 | 0.990 |

Table 1: F1 scores on predicting the sentence with the same structure as the input, through a variational encoder-decoder system. For all eperiments the system uses 2000 training instances, 10% of which are dynamically selected in each experiment for validation.

reflect this notion of features, we use sentences that are parallel in grammatical structure and semantic roles. We use the sentences extracted from the BLM data, as described in Section 2.4, and form instances by pairing an input sentence $s_i$ with structure $str_i$ with a sentence $s_j \neq s_i$ that has the same structure ($str_j = str_i$), and with several (7) negative examples $s_k$ that have different structures ($str_k \neq str_i$). The structure information is only used to build the dataset and obtain a deeper evaluation of the results, but will not be provided to the system. We built separate datasets for Fun and Lex.

To mine for the structure of the sentences we follow the approach described in Nastase and Merlo (2024), which uses a variational encoder-decoder to compress sentences into representations that capture syntactic and semantic information. To encourage the desired information – in this case syntactic-semantic structure – to be encoded on the latent layer, input sentences are paired to correct outputs that have the same internal structure, and use additional contrastive negative candidates that have different structure than the input. There is no overt signal about a sentence's structure.

This approach enables a two-fold evaluation: (i) in terms of performance in detecting the correct structure, by choosing the candidate answer that has the same syntactic-semantic information as the input; (ii) in terms of the compressed representation on the latent layer, which captures these syntactic and semantic properties.

Table 1 shows the averaged F1 scores over three experiments. The results on test data of the same type as the training are very different from those on the test of the other type. This indicates that for each of the Fun and Lex data variations, the system discovers different clues to match two sentences with the same structure. The high results when training on the sentences with functional
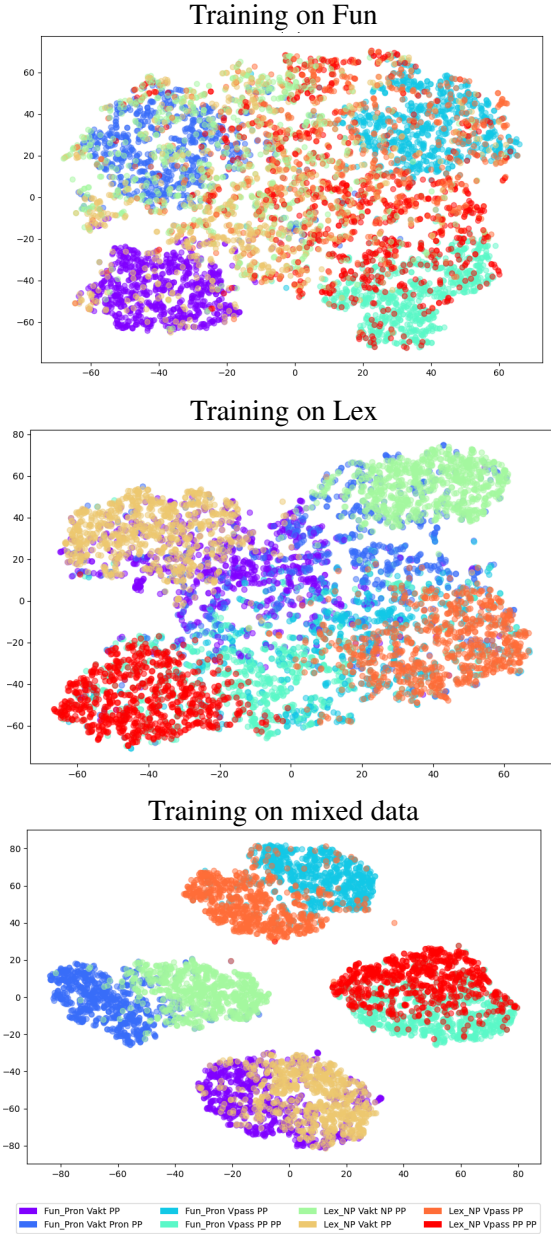
5

Figure 6: Latent representation analysis: t-SNE projection of vectors on the latent layer for the sentences in the training instances.

words may also indicate overfitting. Additional information comes from the analysis of the compressed representations on the latent layer, which are expected to capture the sentence structure that is shared by the functional and lexicalized data. The top two plots of Figure 6 show the projection on the latent layer of the sentence representations with functional and content words, when trained on the sentences with functional words (top) or on the sentences with content words (middle). The plots, matching the F1 scores, show clear clusters for the data that matches the training type, but only slight separation for the data points from the other type.

To test whether there is a shared level of information between sentences with functional or content words, despite what the shallow analyses in Sections 3.1 and 3.2 indicate, we train the system with a dataset containing a mixture of instances. Evidence for shared information will come from two directions: high results on both test sets when training with the mixed training data, and overlapping clusters for the compressed representations on the latent layer. If there is no shared information, the results may be high on each test set (because separately they have been very well modelled), but the clusters of the compressed representations would be separate.

The results in Table 1 shows very high results for both datasets for the mixed data training. The analysis of the representations on the latent layer, at the bottom of Figure 6, shows that the system has discovered a shared space between the sentences with functional and those with content words. What these sentences have in common is the syntactic and semantic structure, and the overlapping clusters of the compressed representations on the latent layer confirms that the system has uncovered this shared structure.

### 3.4 Task solving

We add another step to the investigation into how the shared structure that supports abstraction is encoded in sentence embeddings. Instead of presenting the system with isolated sentences, we present it with change-of-state (COS) and object-drop (OD) verb paradigms, as described in section 2. To choose the correct answer, the relevant linguistic objects (verbs and noun phrases) and their properties (grammatical and semantic roles in the given contexts) must be identified. This dataset also allows us to test generalization at several levels, because of the several levels of lexical variation.

We use the system described by Nastase and Merlo (2024), that solves the BLM problem in two steps: compresses the sentence into a representation that encodes the structure relevant to the BLM puzzle, and use these compressed representations to solve the multiple-choice puzzle. The system construct the representation of an answer, then chooses the closest one from the given options. The two steps are encoded through interconnected variational encoder-decoders, as illustrated in Figure 7, which are trained together. The learning objective is to maximize the score of the correct answer from the candidate answer set, and minimize
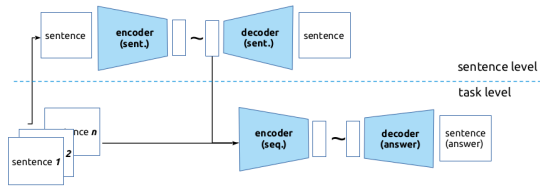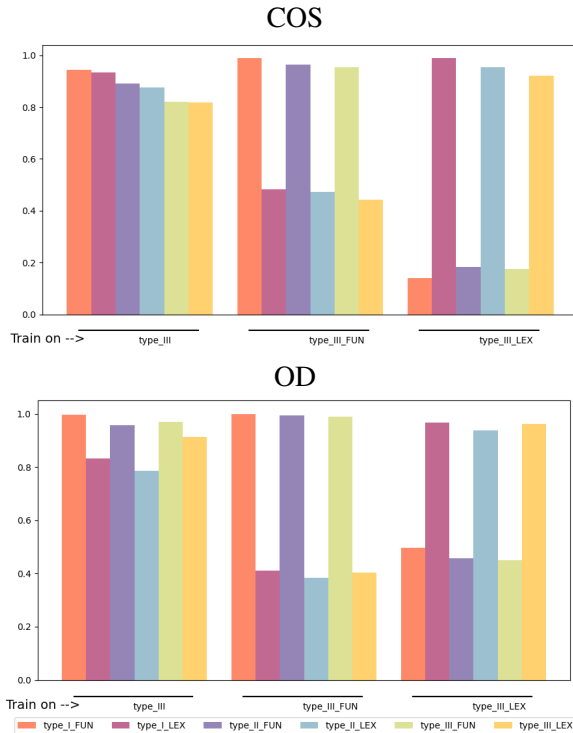
6

Figure 7: Two-step VAE BLM solver



Figure 8: Results in terms of average F1 over three runs for solving the type III (maximal lexical variation) COS and OD BLM tasks for three models. Joint training vs. separate training.

that of the incorrect ones.

Figure 8 shows the F1 results (as averages over three runs) of joint vs. separate training for the two BLM tasks: change of state (COS) and object drop (OD). The results are for type III data, with maximum lexical variation. The complete results are in Tables 3 and 4 in the appendix.

Processing separately datasets of sentences with and without functional words leads to high results within each task, but leads to low results when testing across tasks. This shows, as in the case of the mining for the shared sentence structure, that for each of the Fun and Lex subsets, the systems discovers and exploits different regularities in the training data. Using a mixed training dataset, instead, encourages all systems to find a shared feature space.[5]

---

[5]Other architectures – a feed-forward neural network, and

## 4 Discussion

The primary goal of this paper is to investigate if sentence representations produced by LLMs encode an abstract notion of nominal and prepositional phrase and, as such, if LLMs can generalise through abstraction. Specifically, we investigate whether the contextualised word embeddings and sentence embeddings of structurally identical sentences are similar, whether they contain noun phrases and prepositional phrases or their homologous pronouns and functional place-holders.

**Embeddings of words, and variation of sentences with content or functional words occupy different regions of the embedding space.** This result aligns with observations that LLMs generalize based on idiosyncratic lexical similarity, not on structure (Baroni, 2019; Nikolaev and Padó, 2023). It also indicates that pronouns are not represented as place-holders of lexical nominal expressions. It is interesting to remark, in this respect, that the semantic literature also contains proposals suggesting that pronominal forms are not place-holders, but are better considered as equivalent to noun phrase (NP) descriptions, where they refer to a less abstract, fuller expression in context, in relevant environments (Elbourne, 2002; Lewis, 2022). However, the fact that functional words are represented separately does not immediately imply they cannot be used as place-holders by a process of mapping onto the homologous nominal expressions in a more structured environment. The result of separation of spaces, though, stems from a shallow analysis, and may hide similarities at a deeper level.

**We can detect information about the shared syntactic structure in the embeddings of the functional and lexical variations of the same sentences, in the right environment.** Our follow-up experiments uncover information about shared syntactic structure in Fun and Lex variation of sentences, and of a larger linguistic puzzle. The results show, though, that to find this information we must use both types of data, to direct the system to the right abstraction. It is likely that in absence of this constraint a system may exploit other regularities in the data. It is well-known that this is one of the weaknesses of deep learning systems, stemming from their main strength of discovering and exploiting patterns in data. Contrary to out conclusion that the system has discovered a shared space based on the abstraction of nouns, one might argue that the

---

a variational encoder-decoder – show the same result pattern.

shared space we find is due only to the shared verb. But, had that been the case, the cross-testing results, when training on separate data types, would have been closer to the results on mixed data, given that the verb is not replaced by a functional category and it remains the same across all types of data and sentences. This argument is especially true for the type III subset of the BLM task, which has maximal lexical variation.

We think instead that the results indicate that the model trained on the functional data, which has a very small and consistent vocabulary, relies on shallower features, while the model learned on the lexicalized data is more robust, but not sufficiently abstract. Training the system with mixed data leads not only to a model that performs very well on both data variations, but all sentences are projected into the same compressed embedding space, establishing the necessary links between nominal expressions and thir functional equivalents that support abstraction and generalisation.

## 5   Related work

A generalization taxonomy based on an extensive analysis of publications in NLP that deal with the topic of generalization is proposed in Hupkes et al. (2023). They distinguish five main dimensions for generalization analysis: motivation (concerning the higher-level aims of the model), generalization type (the properties of language or domain or model the model is intended to capture), shift type (the kind of differences between training and testing data distributions), shift source (the source of the difference in data distributions) and shift locus (where in the pipeline does the shift in data distributions occurs). This analysis reflects the focus in the NLP community on the model, and its properties from a machine learning point of view.

Language has its own generalization and abstraction dimensions, which could be at the lexical level (Regneri et al., 2024; Sukumaran et al., 2024), concern verb frames (Wilson et al., 2023; Yi et al., 2022), grammar (Kim and Smolensky, 2021) or a combination of these (Wang et al., 2024). The results of such investigations do not reveal a clear picture. While Kim and Smolensky (2021) observe a limited degree of generalization based on grammatical categories, they note that the results may not have been driven by abstraction. Yi et al. (2022) show that both verb and sentence representations encode information about a verb's alternation class, but the linguistic generalization within the verb argument structure is limited, as models fail on unseen contexts. In experiments on an entailment graph that contains abstract concepts entailed by components of events (nouns, verbs, the event as a whole), Wang et al. (2024) show that the LLMs have difficulty understanding abstract knowledge, but they can be improved with fine-tuning.

Structural priming is used in Michaelov et al. (2023) to investigate the degree of grammatical abstraction in LLMs for three verb alternations: active/passive, dative alternation and two forms of possessive. In monolingual and cross-lingual settings, they find evidence for abstract grammatical representations of these phenomena.

Close to the topic of this paper, Regneri et al. (2024) investigate whether hyponymy is encoded in the transformer by analysing the attention matrices when presented with hyponymous noun pairs. In our work, instead, we have analysed the output of a pretrained language model, and whether the word and sentence embeddings it produces encode particular linguistic information that would allow us to establish a parallel between lexicalized and abstract expressions of a sentence.

All this work shows an unclear picture of sentence embeddings, and the information – and its degree of abstractness – it encodes. Our work provides further linguistically-oriented evidence to clarify the relation between embeddings, abstraction and generalisation.

## 6   Conclusions

Our study contributes to the discussion of generalization in language models, and in particular studies linguistic generalization, rather than task or model generalization. It starts from the assumption that generalisation must proceed by a process of abstraction, which is encoded in the word and sentence embeddings. While the initial shallow analysis of isolated and in-context word embeddings, and the embeddings of the parallel (lexicalized and functional) sentences indicate little superficial shared information, a deeper analysis, searching for sentence structure, has shown that structural information is shared between the representation of lexicalized and functional sentence variations. These conclusions are further reinforced by the results on a problem solving task task, the BLM task, whose solution relies on the proper detection of linguistic objects and their relations.

## 7 Limitations

We use a synthetic dataset, for controlled experimentation, which primarily consists of simple sentence structures. The dataset, then, may not fully capture the complexity of language. Future extensions will include many more structures and variations. Another limitation is the all-or-nothing pronominalisation of sentences, where each sentence is either fully categorized into a predefined functional element or not. Future work will have to modulate the amount of pronominalisation and study different patterns of interactions between nominal expressions and their pronominal equivalent. Moreover, at the moment, we do not have comparable results with a human experiment, which could shed light on more human-like abstraction processes. Finally, this study relies exclusively on English data. While many pronominal systems are structured like the one of English, many other pronominal systems exist. Future studies should add a cross-linguistic dimension.

## References

Marco Baroni. 2019. Linguistic generalization and compositionality in modern artificial neural networks. *Philosophical Transactions of the Royal Society B*, 375.

Yoshua Bengio, Aaron Courville, and Pascal Vincent. 2013. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828.

Daniel Büring. 2019. *1. Pronouns*, pages 1–32. De Gruyter Mouton, Berlin, Boston.

Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. ELECTRA: Pretraining text encoders as discriminators rather than generators. In *ICLR*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805.

Paul Elbourne. 2002. *Situations and individuals*. Ph.D. thesis, Massachusetts Institute of Technology.

Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. 2022. Toy models of superposition. *Preprint*, arXiv:2209.10652.

Martin Haspelmath. 1997. Indefinite pronouns. *Oxford Studies in Typology and Linguistic Theory)/Clarendon Press*.

Geoffrey E Hinton and Sam Roweis. 2002. Stochastic neighbor embedding. In *Advances in Neural Information Processing Systems*, volume 15. MIT Press.

Dieuwke Hupkes, Mario Giulianelli, Verna Dankers, Mikel Artetxe, Yanai Elazar, Tiago Pimentel, Christos Christodoulopoulos, Karim Lasri, Naomi Saphra, Arabella Sinclair, Dennis Ulmer, Florian Schottmann, Khuyagbaatar Batsuren, Kaiser Sun, Koustuv Sinha, Leila Khalatbari, Maria Ryskina, Rita Frieske, Ryan Cotterell, and Zhijing Jin. 2023. A taxonomy and review of generalization research in nlp. *Nature Machine Intelligence*, 5(10):1161–1174.

Najoung Kim and Paul Smolensky. 2021. Testing for grammatical category abstraction in neural language models. In *Proceedings of the Society for Computation in Linguistics 2021*, pages 467–470, Online. Association for Computational Linguistics.

Beth Levin. 1993. *English Verb Classes and Alternations A Preliminary Investigation*. University of Chicago Press, Chicago and London.

Karen S Lewis. 2022. Descriptions, pronouns, and uniqueness. *Linguistics and Philosophy*, 45(3):559–617.

Paola Merlo. 2023. Blackbird language matrices (BLM), a new task for rule-like generalization in neural networks: Motivations and formal specifications. *ArXiv*, cs.CL 2306.11444.

Paola Merlo and Suzanne Stevenson. 2001. Automatic verb classification based on statistical distributions of argument structure. *Computational Linguistics*, 27(3):373–408.

James Michaelov, Catherine Arnett, Tyler Chang, and Ben Bergen. 2023. Structural priming demonstrates abstract grammatical representations in multilingual language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3703–3720, Singapore. Association for Computational Linguistics.

Vivi Nastase and Paola Merlo. 2024. Are there identifiable structural parts in the sentence embedding whole? In *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 23–42, Miami, Florida, US. Association for Computational Linguistics.

Dmitry Nikolaev and Sebastian Padó. 2023. Representation biases in sentence transformers. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 3701–3716, Dubrovnik, Croatia. Association for Computational Linguistics.

Michaela Regneri, Alhassan Abdelhalim, and Soeren Laue. 2024. Detecting conceptual abstraction in LLMs. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 4697–4704, Torino, Italia. ELRA and ICCL.

Priyanka Sukumaran, Conor Houghton, and Nina Kazanina. 2024. Investigating grammatical abstraction in language models using few-shot learning of novel noun gender. In *Findings of the Association for Computational Linguistics: EACL 2024*, pages 747–765, St. Julian's, Malta. Association for Computational Linguistics.

Andreas Waldis, Yotam Perlitz, Leshem Choshen, Yufang Hou, and Iryna Gurevych. 2024. Holmes a benchmark to assess the linguistic competence of language models. *Transactions of the Association for Computational Linguistics*, 12:1616–1647.

Zhaowei Wang, Haochen Shi, Weiqi Wang, Tianqing Fang, Hongming Zhang, Sehyun Choi, Xin Liu, and Yangqiu Song. 2024. AbsPyramid: Benchmarking the abstraction ability of language models with a unified entailment graph. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3991–4010, Mexico City, Mexico. Association for Computational Linguistics.

Michael Wilson, Jackson Petty, and Robert Frank. 2023. How abstract is linguistic generalization in large language models? experiments with argument structure. *Transactions of the Association for Computational Linguistics*, 11:1377–1395.

David Yi, James Bruno, Jiayu Han, Peter Zukerman, and Shane Steinert-Threlkeld. 2022. Probing for understanding of English verb classes and alternations in large pre-trained language models. In *Proceedings of the Fifth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 142–152, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

# A Data

| Class | Verb |
|-------|------|
| COS | *bake, bend, blacken, break, brighten, caramelize, chip, close, corrode, crinkle, defrost, empty, expand, fry, harden, harmonize, heat, improve, increase, intensify, melt, open, propagate, purify, sharpen, shrink, sweeten, tear, whiten, widen.* |
| OD | *clean, cook, draw, drink, eat, fish, hum, iron, knead, knit, mend, milk, nurse, paint, play, plow, polish, read, recite, sculpt, sew, sing, sow, study, sweep, teach, wash, weave, whittle, write.* |

Table 2: Verbs categorized by class

| | COSFUN - CONTEXT | | COSFUN - ANSWERS |
|---|---|---|---|
| 1 | She broke it with this | 1 | **It broke by those there** |
| 2 | She broke it by those there | 2 | She broke by those there |
| 3 | It was broken by her with this | 3 | It was broken by her |
| 4 | It was broken by her by those there | 4 | She was broken by it |
| 5 | It was broken with this | 5 | It broke her |
| 6 | It was broken by those there | 6 | She broke it |
| 7 | It broke with this | 7 | It broke by her |
| ? | ??? | 8 | She broke by it |
| | COSLEX - CONTEXT | | COSLEX - ANSWERS |
| 1 | The archaeologist broke a vase in the lab | 1 | **The vase broke by mistake** |
| 2 | The archaeologist broke a vase by mistake | 2 | The archaeologist broke by mistake |
| 3 | The vase was broken by the archaeologist in the lab | 3 | The vase was broken by the archaeologist |
| 4 | The vase was broken by the archaeologist by mistake | 4 | The archaeologist was broken by the vase |
| 5 | The vase was broken in the lab | 5 | The vase broke the archaeologist |
| 6 | The vase was broken by mistake | 6 | The archaeologist broke the vase |
| 7 | The vase broke in the lab | 7 | The vase broke by the archaeologist |
| ? | ??? | 8 | The archaeologist broke by the vase |

Figure 9: Examples of FUN and LEX for the English verb *break*, one of the verbs belonging to COS class.

| | ODLEX - CONTEXT | | ODFUN - ANSWERS |
|---|---|---|---|
| 1 | They paint it with this | 1 | It painted by that |
| 2 | They paint it by that | 2 | **They painted by that** |
| 3 | It was painted by them with this | 3 | It was painted by them |
| 4 | It was painted by them by that | 4 | They were painted by it |
| 5 | It was painted with this | 5 | It painted them |
| 6 | It was painted by that | 6 | They painted it |
| 7 | They painted with this | 7 | It painted by them |
| ? | ??? | 8 | They painted by it |
| | COSLEX - CONTEXT | | COSLEX - ANSWERS |
| 1 | These artists paint a portrait with a brush | 1 | A portrait painted by the lake |
| 2 | These artists paint a portrait by the lake | 2 | **These artists painted by the lake** |
| 3 | A portrait was painted by these artists with a brush | 3 | A portrait was painted by the artists |
| 4 | A portrait was painted by these artists by the lake | 4 | These artists were painted by a portrait |
| 5 | A portrait was painted with a brush | 5 | A portrait painted these artists |
| 6 | A portrait was painted by the lake | 6 | These artists painted a portrait |
| 7 | These artists painted with a brush | 7 | A portrait painted by these artists |
| ? | ??? | 8 | These artists painted by a portrait |

Figure 10: Examples of Type_I FUN and LEX data for the English verb *paint*, one of the verbs belonging to OD class

# B  BLM task results

The experiments were run on an HP PAIR Workstation Z4 G4 MT, with an Intel Xeon W-2255 processor, 64G RAM, and a MSI GeForce RTX 3090 VENTUS 3X OC 24G GDDR6X GPU.

| test on | train on | | |
|---|---|---|---|
| | Joint training | | |
| | type_I | type_II | type_III |
| type_I_Fun | 0.983 | 0.987 | 0.997 |
| type_I_Lex | 0.763 | 0.723 | 0.833 |
| type_II_Fun | 0.857 | 0.897 | 0.957 |
| type_II_Lex | 0.690 | 0.680 | 0.787 |
| type_III_Fun | 0.920 | 0.967 | 0.970 |
| type_III_Lex | 0.837 | 0.887 | 0.913 |
| | Training on Fun | | |
| | type_I_Fun | type_II_Fun | type_III_Fun |
| type_I_Fun | 1.000 | 1.000 | 1.000 |
| type_I_Lex | 0.510 | 0.553 | 0.410 |
| type_II_Fun | 0.907 | 0.963 | 0.993 |
| type_II_Lex | 0.457 | 0.490 | 0.383 |
| type_III_Fun | 0.963 | 0.983 | 0.990 |
| type_III_Lex | 0.407 | 0.477 | 0.403 |
| | Trainig on Lex | | |
| | type_I_Lex | type_II_Lex | type_III_Lex |
| type_I_Fun | 0.460 | 0.457 | 0.497 |
| type_I_Lex | 0.733 | 0.763 | 0.967 |
| type_II_Fun | 0.450 | 0.450 | 0.457 |
| type_II_Lex | 0.680 | 0.717 | 0.937 |
| type_III_Fun | 0.540 | 0.523 | 0.450 |
| type_III_Lex | 0.877 | 0.927 | 0.963 |

Table 3: BLM-COS: Results as averaged F1 over three runs, for three training set-ups: joint training (training using both Fun and Lex instances), training on Fun instances, training on Lex instances. For all set-ups we use 2000 training instances. For the joint training these are evenly split between Fun and Lex. Standard deviation is less that 1e-3, so we do not include it.

| test on | train on | | |
|---|---|---|---|
| | Joint training | | |
| | type_I | type_II | type_III |
| type_I_Fun | 0.983 | 0.987 | 0.997 |
| type_I_Lex | 0.763 | 0.723 | 0.833 |
| type_II_Fun | 0.857 | 0.897 | 0.957 |
| type_II_Lex | 0.690 | 0.680 | 0.787 |
| type_III_Fun | 0.920 | 0.967 | 0.970 |
| type_III_Lex | 0.837 | 0.887 | 0.913 |
| | Train on Fun | | |
| | type_I_Fun | type_II_Fun | type_III_Fun |
| type_I_Fun | 1.000 | 1.000 | 1.000 |
| type_I_Lex | 0.510 | 0.553 | 0.410 |
| type_II_Fun | 0.907 | 0.963 | 0.993 |
| type_II_Lex | 0.457 | 0.490 | 0.383 |
| type_III_Fun | 0.963 | 0.983 | 0.990 |
| type_III_Lex | 0.407 | 0.477 | 0.403 |
| | Training on Lex | | |
| | type_I_Lex | type_II_Lex | type_III_Lex |
| type_I_Fun | 0.460 | 0.457 | 0.497 |
| type_I_Lex | 0.733 | 0.763 | 0.967 |
| type_II_Fun | 0.450 | 0.450 | 0.457 |
| type_II_Lex | 0.680 | 0.717 | 0.937 |
| type_III_Fun | 0.540 | 0.523 | 0.450 |
| type_III_Lex | 0.877 | 0.927 | 0.963 |

Table 4: BLM-OD: Results as averaged F1 over three runs, for three training set-ups: joint training (training using both Fun and Lex instances), training on Fun instances, training on Lex instances. For all set-ups we use 2000 training instances. For the joint training these are evenly split between Fun and Lex. Standard deviation is less that 1e-3, so we do not include it.