



**HAL**  
open science

# Compositionality as an Analogical Process: Introducing ANNE

Giulia Rambelli, Emmanuele Chersoni, Philippe Blache, Lenci Alessandro

► **To cite this version:**

Giulia Rambelli, Emmanuele Chersoni, Philippe Blache, Lenci Alessandro. Compositionality as an Analogical Process: Introducing ANNE. Workshop on Cognitive Aspects of the Lexicon, Proceedings of the Workshop on Cognitive Aspects of the Lexicon , pp.78-96, 2022, 978-1-959429-01-2. hal-04052104

**HAL Id: hal-04052104**

**<https://hal.science/hal-04052104>**

Submitted on 30 Mar 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# Compositionality as an Analogical Process: Introducing ANNE

**Giulia Rambelli**

University of Pisa  
giulia.rambelli@phd.unipi.it

**Emmanuele Chersoni**

The Hong Kong Polytechnic University  
emmanuelechersoni@gmail.com

**Philippe Blache**

Aix-Marseille University/CNRS  
philippe.blache@univ-amu.fr

**Alessandro Lenci**

University of Pisa  
alessandro.lenci@unipi.it

## Abstract

Usage-based constructionist approaches consider language a structured inventory of *constructions*, form-meaning pairings of different schematicity and complexity, and claim that the more a linguistic pattern is encountered, the more it becomes accessible to speakers. However, when an expression is unavailable, what processes underlie the interpretation? While traditional answers rely on the principle of compositionality, for which the meaning is built word-by-word and incrementally, usage-based theories argue that novel utterances are created based on previously experienced ones through *analogy*, mapping an existing structural pattern onto a novel instance.

Starting from this theoretical perspective, we propose here a computational implementation of these assumptions. As the principle of compositionality has been used to generate distributional representations of phrases, we propose a neural network simulating the construction of phrasal embedding as an analogical process. Our framework, inspired by word2vec and computer vision techniques, was evaluated on tasks of generalization from existing vectors.

## 1 Introduction

While the generative tradition has dominated linguistic research for over half a century, the last decades have seen the emergence of an alternative paradigm in linguistics and cognitive sciences, which goes under the name of *usage-based models of language* (Langacker, 1987; Croft, 1991, 2001; Givón, 1995; Tomasello, 2009; Bybee, 2010), a variety of approaches grounded on the idea that linguistic structures emerge and are shaped through the use of language. Their claim is that language is not different from any other cognitive domain: Linguistic structures are not the result of a specific-language function but are explainable as the implementation of domain-general processes (Ibbotson,

2013). The usage-based position shares the fundamental assumption of Construction Grammar (Hoffman and Trousdale, 2013): language consists of meaningful and symbolic form–meaning mappings, called *constructions*. Words, idiomatic expressions (e.g., *kick the bucket* —“to die” or *jog <someone’s> memory* —“to refresh <someone’s> memory”) and highly general and productive syntactic patterns (e.g., ditransitive structures) are all constructions varying along a continuum of schematicity and complexity.

In particular, usage-based constructionist approaches emphasize the notion of frequency: combinations that are more frequently encountered become more accessible (perhaps because they are stored in memory) and are preferred. Indeed, if the language system derives from language use, it follows that how often a speaker encounters a particular linguistic expression will affect the system itself. This assumption implies that any sequence of words – if used frequently enough – can be a construction, even if there are no idiosyncrasies of form and meaning (Goldberg, 2006). However, it is impossible to store any possible word combinations a speaker has or will ever produce. The traditional answer relies on the *principle of compositionality*: the meaning of a complex expression is entirely determined by its structure and the meanings of its constituents – once we specify what the parts mean and how they are put together, there is no more leeway regarding the meaning of the whole (Partee, 2004). Usage-based theories favor a different explanation: novel utterances are created based on previously experienced utterances thanks to the cognitive process of *analogy*.

The ability to make analogies – that is, to map familiar relations from one domain of experience to another – is a fundamental ingredient of human intelligence and creativity (Hofstadter, 2001). In the linguistic domain, analogy depends on similarity in form and meaning between constructions,

whether these constructions are of a concrete type or an abstract type: a novel instance is compared to those stored in our long-term memory to infer the new representation. In this perspective, the acceptability of a novel item is a gradient that depends on the extent of similarity to prior uses of a construction (Bybee, 2010). In a more radical stance, Ambridge (2020) proposed disregarding completely abstraction: unwitnessed forms are produced and comprehended “by on the fly analogy” across multiple stored exemplars. Without denying the existence of abstract representations, we also assume that analogical mechanisms play a key role in explaining systematic processes of language productivity.

This paper aims to articulate the hypotheses introduced above in computational terms. We address two interconnected questions: How can we represent (lexicalized) constructions? Is it possible to replicate the interpretation-as-analogy mechanism in computational terms? Specifically, we investigate how to model constructions as well as analogy-based compositionality using Distributional Semantic Models (DSMs). DSMs represent the lexicon in terms of vector spaces, where a lexical target is described in terms of a vector (also known as embedding) built by identifying in a corpus its syntactic and lexical contexts (Lenci, 2018).

As a first approximation, we decided to consider constructions any kind of frequent pairs of words linked by a syntactic relation. Traditionally, building distributional representations beyond individual words, such as phrases and sentences, is the focus of *Compositional Distributional Semantic Models*. Their proposed methodologies try to derive the meaning of an expression from the meanings of the sentence’s constituents (Baroni et al., 2014): the simplest CSDMs represent words as vectors and obtain sentence vectors with sum or product operations between constituent vectors (Mitchell and Lapata, 2010), while more complex models represent predicates with matrices and tensors (Baroni and Zamparelli, 2010; Coecke et al., 2010; Baroni et al., 2014; Paperno et al., 2014) or reproduce the compositionality operation by means of a neural architecture learning so-called sentence embeddings (Socher et al., 2012; Cheng and Kartsaklis, 2015). It is interesting to notice that most distributional models for phrases/constructions/sentences assume more or less explicitly the principle of compositionality, while the idea that units above the word level

could be stored and retrieved via analogy/similarity mechanisms has rarely been explored.<sup>1</sup>

The experiment presented here distances itself from these approaches, following a more usage-based perspective. Suppose frequently experienced word sequences are, to some extent, stored in memory, and the organization and productivity of language are understood as the result of analogical processes between form and meaning in this structured inventory of constructions. In that case, new phrases could be constructed by analogy with stored linguistic patterns. We propose a neural network model to infer a distributional representation of a new syntactic phrase by preserving the structural information encoded in the embeddings representing previously stored, high-frequency phrases.

As the main contributions of the paper, i) we introduce a new DSM in which both lemmas and syntactic relations in the form of  $\langle head, dependent, syntactic\ role \rangle$  triples have a unique distributional representation; ii) we propose an analogical model to create the distributional embeddings of new relations by applying deep-learning techniques, and evaluate different architectures in terms of generalization and systematicity; iii) we discuss the implications of our analogical model from a theoretical and computational perspective.

## 2 Relational Embeddings

The first step consisted in developing a DSM for lexicalized constructions. We represent the meaning of phrases following a holistic approach (Turney, 2012): as a numeric vector can represent nouns like *space* and *race*, in the same way, phrases like *space race* are associated with a unique embedding. For our goal, we built embeddings corresponding to triples  $\langle head, dependent, role \rangle$ , assuming that these vectors should keep track of the syntactic relation between words. For this reason, we called these **Relational Embeddings** (RelEmbs), and we assume they represent the meaning of lexicalized constructions.

We built our semantic space using `word2vecf` (Levy and Goldberg, 2014), a modification of the skip-gram model introduced by Mikolov et al. (2013a). While the original implementation assumes bag-of-words contexts, i.e., the model keeps

<sup>1</sup>Some partial exceptions are instance-based distributional models (Jones and Mewhort, 2007; Jamieson et al., 2018; Crump et al., 2020) and distributional models of event knowledge that store event occurrences in the form of syntactic graphs (Chersoni et al., 2019, 2021).

track of word counts and disregards the grammatical details and the word order, word2vecf allows us to use arbitrary context features. In detail, we extracted <target, context> occurrences from the concatenation of ukWaC and a 2018 dump of English Wikipedia, parsed using CoreNLP (Manning et al., 2014): targets are both words and <head, dependent, role> triples (e.g., <bark, dog, nsubj>), while context is always an open-class word (noun, verb, adjective) occurring with the target in the sentence within a window  $\pm 10$  (ten words before and ten words after the head of the relation excluding the dependent). Word2vecf parameters are reported in Appendix A. We built our DSM considering only words and relations with a frequency equal to or larger than 100 and filtering out <target, context> pairs with a frequency less than 20; lastly, we kept only <head, dependent, role> triples with a frequency  $\geq 1,000$ , where both the head and the dependent lemmas have a frequency  $\geq 10,000$ . This strategy is consistent with the idea that holistic representations of complex constructions are stored only for substantially frequent items. The final space contains 127,739 word embeddings and 173,496 RelEmbs, for a total of 301,235 items.

**Semantic space evaluation** We tested the quality of the semantic space over some most common benchmarks for the intrinsic evaluation of word and phrase embeddings. It is worth mentioning that we are not aiming at beating traditional DSMs, but rather at carrying out a general evaluation of the goodness of our distributional representations of lexicalized constructions.

For word embeddings, we ran the standard Word Similarity/Relatedness task using the well-known **WordSim-353** (Finkelstein et al., 2001) and **MEN** (Bruni et al., 2014). The task is to compute the cosine similarity between two words (e.g., *cup* and *mug*) and verify how their score correlates with the similarity rate given by humans. We also evaluated the DSM against **FAST** (Evert and Lapesa, 2021), a free associations dataset. The goal of this multiple-choice task is to determine the most frequent associate for a given stimulus among three candidates (e.g., which word between *neck*, *apple*, *wine* is most associated with *giraffe*?). As a baseline, we computed the performance of a DSM trained with the original word2vec Skip Gram model (Mikolov et al., 2013a) on the same concatenation of corpora.

Results are reported in Table 1. Considering the first task, we observe that Spearman’s correlation

scores for the baseline are a bit higher than our DSM in all settings, except for the MEN dataset. However, the differences are not statistically significant.<sup>2</sup> It is worth noticing that similarity results are better than relatedness results, showing the same trend reported in Agirre et al. (2009). We observe an opposite performance for the classification task: our space consistently beats the baseline, and the difference is statistically significant.

Dataset	RelEmbs.w	baseline	Coverage
WS353-all	0.684	0.721	333/353
WS353-sim	0.734	0.75	195/203
WS-353-rel	0.628	0.675	236/252
MEN	0.774	0.735	3000/3000
FAST-EAT	0.786***	0.737	5877/7610
FAST-USF	0.725***	0.719	4057/4719

Table 1: Word embeddings evaluation. On top: Spearman’s correlation scores for Word Similarity/Relatedness task. Bottom: Accuracy scores for Free Association task. \*\*\* =  $p < 0.01$  using McNemar test.

Moving to the relational embeddings, we used the Mitchell et al. (2010) Phrase Similarity dataset (**ML10**), which includes 324 English phrase pairs, tripartite in noun phrases, verb phrases, and adjective phrases. Given two expressions (e.g., *general principle* and *basic rule*), the task consists in comparing the cosine similarity between the two corresponding vectors and then correlating the score with the human similarity rating. As a baseline, we represented the phrases as the sum of the word2vec vectors used for word embedding evaluation. Table 2 reveals that correlation scores are not homogeneous among the different sets: the noun phrase subset achieves a higher score (0.635) compared to the other two sets, whose score is lower than 0.5. Moreover, baseline results are consistently better than our model and are statistically significant for the AN subset.

Dataset	RelEmbs	baseline	Coverage
ML-vo	0.499	0.599	99/108
ML-nn	0.635	0.716	99/108
ML-an	0.462	0.683**	102/108

Table 2: Relational embeddings evaluation. Spearman’s correlation scores for Phrase Similarity task. \*\* =  $p < 0.01$  using Fisher r-to-z transformation test.

<sup>2</sup> $p > 0.1$ , the  $p$ -value is computed with Fisher’s r-to-z transformation, one-tailed test

**Qualitative analysis** Results in Table 2 suggest that RelEmbs perform worst than the baseline in the phrase similarity test. To gain more insight, we selected some problematic pairs from the ML10 dataset and manually inspected the  $k$ -nearest neighbors, i.e., the most similar words by cosine similarity. Let us look at the pair *reduce amount* and *cut cost*: the two expressions are judged very similar (6.55), but their cosine similarity is just 0.41. However, their distributional neighbors are coherent and somehow systematic in the sense that they are similar to relational embeddings in which the same head or dependent word occurs. So, *reduce amount* is mostly similar to *increase amount*, *reduce waste*, *a person reduce*, *large amount*, *high amount*; on the other hand, the neighbors of *cut cost* are *reduce cost*, *improve efficiency*, *increase profit*, *lower cost*, *save money*. Similar observations are for nominal phrases, like *government leader* and *health minister*. While ML10 reports a high score (4.95), the cosine similarity between the two is quite low (0.43). However, this is explainable by observing their neighbors. In the first case, *health minister* is similar to other types of ministers (*health secretary*, *transport minister*, *environment minister*, *minister for health*); conversely, *government leader* is more associated with situations (*invite a leader*, *include a leader*) or other offices (*chief whip*, *head of the committee*, *regional leader*) associated to leaders. In other words, while the phrases refer to government members, the two roles are not the same (and functions also differ).

To sum up, the qualitative analysis of the neighbors reveals that RelEmbs form a semantically coherent space, even though they do not outperform the baseline in the phrase similarity task.

### 3 Analogical Neural Network for Embeddings

Usage-based theories of language assume that systematic processes of language productivity can be explained mainly by analogical inferences rather than by sequential compositional operations. In this perspective, we present a system to expand the coverage of the RelEmbs space simulating the construction of phrasal meaning as an analogical process via deep learning techniques.

#### 3.1 Architecture

We aim to infer a distributional representation of a new syntactic phrase (ANALOGICAL TARGET) by

preserving the structural information encoded in an existing relational embedding (ANALOGICAL BASE). For simplicity, we represent this process using the familiar four-term formalism.<sup>3</sup> Approximately, solving the analogy A:B::C:? requires a system that generates an appropriate embedding to make a valid analogy: if we need to infer an embedding for the target phrase *drink cider* using *drink water* as the base, we can reformulate the analogy as: *water* : *drink\_water*<sub>dobj</sub> :: *cider* : ?

We framed the problem of analogy completion as a regression task: the aim is to build a phrasal vector given the embeddings of the other expressions in the analogy. While word embeddings have been widely employed to perform analogy by addition and subtraction of word embeddings (Mikolov et al., 2013b; Gladkova et al., 2016), we argue that directly training a deep neural network on the task of analogy completion could provide better results, as already proposed by Reed et al. (2015) for visual analogy-making. We named our novel neural network model as **Analogical Neural Network for Embeddings (ANNE)**.

In detail, we implemented a feed-forward neural network architecture with one hidden layer: the model is trained to learn a function  $f : R^{2D} \rightarrow R^D$  that maps an input vector  $x$  to a generated embedding  $y$  of dimension  $D$  (where  $D=300$ ), preserving the structural properties of the selected base. The input vector  $x$  should incorporate the analogical base (e.g., *drink water*), and the new argument (e.g., *cider*). We tested two possible combinations: i) the input vector is the concatenation between the analogical base and the new argument (CONCAT, Figure 1a); ii) we compute the difference between the analogical base and the argument in the same relation; the resulting vector is concatenated to the new argument vector (DIFF; Figure 1b). The intuition below the DIFF input representation is that we apply some aspects of Mikolov’s analogical operation with the nonlinearities and supervision offered by a neural network.

We developed several variants of this network, each with a distinct objective function. The basic architecture (SIMPLE) is trained to maximize the cosine similarity between the original and predicted embedding. However, ANNE should not simply create a vector similar to the actual instance in the DSM but also learn the relational structure

<sup>3</sup>It is, however, doubtful that linguistic analogies are computed in this way at the brain level (Bybee, 2010).

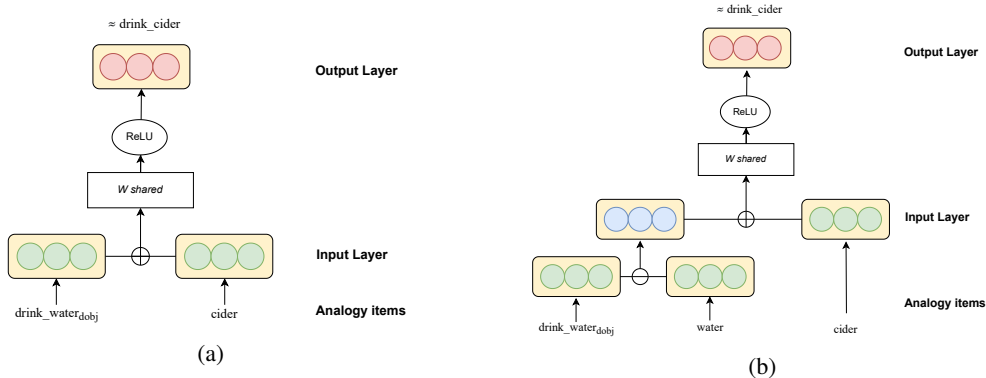


Figure 1: ANNE architecture with CONCAT (a) and DIFF (b) input. The  $\oplus$  indicates a vector concatenation, while  $\ominus$  indicates vector subtraction.

of the base and transpose it to the generated embedding. To this end, we implemented a multiple losses function, which combines the SIMPLE loss with a new loss aimed at minimizing the difference of the similarity between the relational vector and the embeddings of its words computed for the base and the target. For instance, if the similarity between *drink water* and *drink* is 0.60 and the similarity between *drink water* and *water* is 0.49, similar scores should be obtained by computing the similarity of the output vector with the vectors of *drink* and *cider*. Therefore, the network is trained to generate ReLEmb that preserve the same relations with its components as the ones in the analogical base. To compare the similarity scores we tested three functions (cf. Appendix B): the Mean Squared Error (MULTIMSE), the mean of the scores difference (MULTIAVG), or the hinge loss function (MULTIHINGE).

**Training the network** We selected analogical base-target pairs from the Relational Embeddings attested in space to train the neural network. We assembled the dataset as follows: given a relation  $\langle head, dependent, syntactic\ relation \rangle$ , we selected all ReLEmb with the same head and syntactic relation and calculated the similarity between the arguments of each pair, keeping only those pairs with the cosine similarity between arguments  $\geq 0.4$  and the cosine similarity between ReLEmb  $\geq 0.6$ . We chose the filter on similarity heuristically: the idea is that the candidate targets should be somehow similar with respect to their analogical bases but not the exact synonyms. The final dataset consisted of 350,404 items and was divided into **Train** and **Test** parts (respectively, 95% and 5%).

To verify the analogy-solving capability of the

network, i.e., its ability to generalize from the base, we kept some analogical pairs out of the training step. The resulting data (named **Test-unseen**), comprises 3,201 pairs (cf. Appendix C). This dataset should verify the network’s performance when encountering new relations, which is to say, evaluate the model’s generalization ability. The training setup configurations are reported in Appendix D.

### 3.2 RSA Evaluation

A preliminary evaluation of ANNE consisted of computing the similarity between relational embeddings attested in the DSM and embeddings analogically generated from ANNE attested in the Test and the Test-unseen datasets. We applied the Representational Similarity Analysis (RSA; Kriegeskorte et al. (2008); Kriegeskorte and Kievit (2013)), a computational technique that allows us to compare heterogeneous representations in higher-order spaces. The core idea is simple: instead of directly correlating representations of stimuli in different representation spaces, we compute how similar representations are between pairs of stimuli in each space, and the resulting similarity matrices are then compared. As we are interested in understanding how similar the original and generated embeddings are, we created a pair of matrices where rows are the vectors representing the analogical targets from a test set and columns correspond to a subset of the ReLEmb vocabulary.<sup>4</sup> Following Lenci et al. (2022), we randomly sampled 100 disjoint sets of 1,000 lexemes, ran RSA analyses on each sample, and then computed the average score.

Table 3 reports Spearman’s  $\rho$  between the similarity matrix computed with the original ReLEmb

<sup>4</sup>A matrix with 301,235 columns would be computationally too expensive.

and the matrix with vectors generated with ANNE. We can observe that the models reach similar results for the two test data, even if the Test-unseen scores are always slightly lower than those for the Test set. Overall, the SIMPLE model reaches the best scores (0.851 for Test and 0.835 for Test-unseen), while MULTIAVG performs the worst (reaching just 0.739 for Test-unseen with DIFF input). However, the average correlation of all models is significantly high. As a baseline, we also performed the vector offset method (Mikolov et al., 2013a). RSA correlation scores significantly drop (0.734 and 0.71 for the Test and Test-unseen, respectively). The worst architecture (MULTIAVG diff) is still better than the baseline for Test ( $p < 0.1$ ), but not for Test-unseen. The best architecture (SIMPLE<sub>concat</sub>) is different from the baseline with  $p < 0.001$ .<sup>5</sup> This result corroborates our assumption that the ANNE architecture is better at generating analogical vectors than a simple vector operation.

	TEST	TEST-UNSEEN
simple <sub>concat</sub>	<u>0.851</u>	<u>0.835</u>
simple <sub>diff</sub>	0.848	0.834
multiHinge <sub>concat</sub>	0.819	0.805
multiHinge <sub>diff</sub>	0.806	0.788
multiAVG <sub>concat</sub>	0.782	0.754
multiAVG <sub>diff</sub>	0.77	0.739
multiMSE <sub>concat</sub>	0.835	0.82
multiMSE <sub>diff</sub>	0.824	0.804
baseline	0.734	0.71

Table 3: Average Spearman’s correlation between original and analogically generated semantic spaces computed with RSA on 100 random samples of 1, 000 words for Test and Test-unseen datasets.

## 4 Compositionality vs. Idiomaticity

Finally, we present a series of analyses to evaluate the meaning encoded in analogically-generated embeddings. We hypothesize that the best-generated embedding should keep the same relationship among components as the base (*systematicity*). As a counterproof, we also generated embeddings from idiomatic expressions. In this case, we expect analogies with idiomatic bases to give odd results in the semantic space because of their reduced compositionality and systematicity. The results should answer the following questions: What are the char-

<sup>5</sup> $p$ -values for Fisher’s  $r$ -to- $z$  transformation, one-tailed test.

acteristics of analogically-generated embeddings? How does the type of input (concatenation or difference) affect the final representation? What loss functions are better at retaining the same structural relation of the base, while at the same time generalizing from the original embedding?

**Data** The analogical bases employed are 44 verbal phrases (22 idioms from Libben and Titone (2008) + 22 compositional manually picked from frequent relations) and 24 nominal compounds (12 idiomatic + 12 compositional) selected from the Noun Compound Senses dataset (Cordeiro et al., 2019) and the dataset by Reddy et al. (2011).

For each phrase, we manually chose a relation similar to the base but not attested in the vocabulary space, with the same head and syntactic role. For example, given the relation  $\langle market_N, fish_N, compound \rangle$  (“a fish market”), we replaced the noun *fish* (i.e., the dependent) with the noun *shrimp*; expressly, the relation  $\langle market_N, shrimp_N, compound \rangle$  (“a shrimp market”) is not attested in RelEmbs vocabulary. The final dataset consists of 68 analogical pairs, half with an idiomatic base and half with the compositional counterpart.

	Idiomatic $\rightarrow$ Target	Compositional $\rightarrow$ Target
VN	<i>break ice</i> $\rightarrow$ <i>break chunk</i>	<i>break bone</i> $\rightarrow$ <i>break finger</i>
NC	<i>loan shark</i> $\rightarrow$ <i>credit shark</i>	<i>reef shark</i> $\rightarrow$ <i>atol shark</i>

Table 4: Examples of analogical pairs (the idiomatic/compositional base on the left, the target on the right of the arrow).

### 4.1 Analysis 1: Correlation of the Similarities with the Components

To evaluate if and how the ANNE configurations are generating embeddings systematically, we observed if the similarities between the relational embedding and those of the component words are similar for both the analogical base and the generated target. The assumption is that the embedding generated by ANNE should have the same internal structure as the base from which it is inferred: that is, the relationship between the phrase meaning and the meaning of its components should be systematically retained in the generated distributional vector. This idea can be approximated by the similarities between the RelEmbs and its parts: if the similarity between *break (a) bone* and *break* is 0.4 (*simHead* score) and the similarity with the dependent *bone* is 0.42 (*simDep* score), comparable scores should

be obtained computing the similarities of *break* (a) *finger* with *break* and *finger*, respectively.

We computed the cosine similarity scores for all ANNE implementations. We assume that the best architecture (i.e., the one that best fits our theoretical hypothesis) should be the one that has i) comparable similarity distributions for compositional bases and derived targets (for both word components), and ii) different (or incoherent) similarity distributions for targets generated from idiomatic bases. By looking closely at the plots in Appendix E, we observe that each architecture produces different outputs. Among all models, MULTIAVG is the one performing worse (plots in (c) and (d)): the generated embeddings have high similarities with the dependent component in both idiomatic and compositional cases, possibly because they retain too much distributional information from dependent words used to generate the new embedding. The MULTIMSE (plots in (e) and (f)) and the SIMPLE (plots in (a) and (b)) losses show a similar behavior: they give a high *simHead* and *simDep* to vector generated from idiomatic targets. This result shows that, when deriving a new literal phrase meaning from a figurative one should be impossible, the models largely rely on attribute similarity instead of truly learning a relation. In this sense, the MULTIMSE<sub>diff</sub> (plots in (f)) model is the only one that perfectly respects our hypothesis (distributions should be the same for targets from compositional bases but different for targets with idiomatic bases). Conversely, the MULTIHINGE model (plots in (g) and (h)) reduces the impact of the dependent word, as proved by the fact that the mean similarity of *simDep* is lower for the target (orange) than for the base (blue).

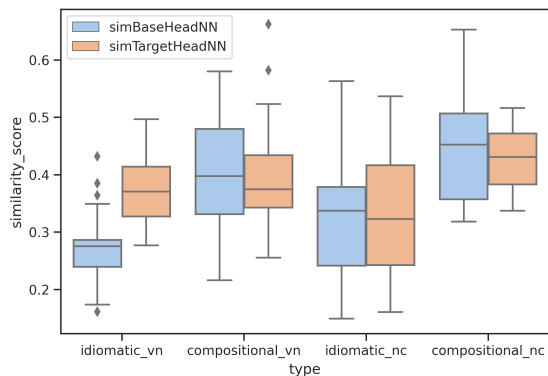


Figure 2: Distribution of the similarities between RelEmbs and their head for MULTIMSE<sub>diff</sub>. Blue boxplots refer to the base embedding, orange to the analogically-generated ones.

## 4.2 Analysis 2: Intersection of Neighbors

As a complementary measure to cosine similarity, we computed the intersection between the 50-nearest neighbors of i) the base and the generated target, and ii) the generated target and the respective head/dependent.<sup>6</sup> The first measure tells us how much information the analogical embedding retains from its base: the higher the value, the higher their similarity, so it could be that the network did not generalize from the input. The second measure should say how much the analogy moved the distribution towards the component meanings. Appendix F reports the results as a series of heatmaps.

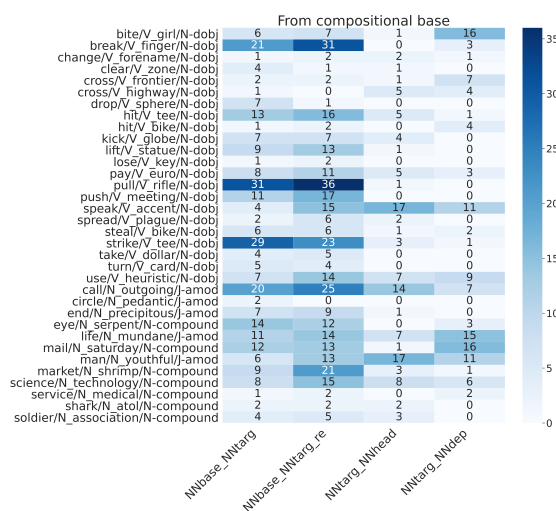


Figure 3: Heatmap for MULTIMSE<sub>concat</sub> shows the intersection between the neighbors of the analogically-generated embedding and the base ( $NN_{base}NN_{targ}$ , \*re only RelEmbs), the head ( $NN_{targ}NN_{head}$ ), and the dependent ( $NN_{targ}NN_{dep}$ ).

Embeddings generated from idiomatic bases have no shared neighbors with the vectors of their heads or dependents: as the network was not trained on this type of analogies (which are impossible), we expected the neural model to fail. What we can add, however, is that sometimes it generates a new embedding that has no common neighbors with either its base or components, sometimes it resolves the analogy by copying the distributional signature of the base. The only exception to this trend is MULTIAVG (subfigures (c) and (d)): we notice that some items, such as *shark credit*, *cockroach market*, and *gastropod mail*, partake many neighbors with their dependent, revealing that this

<sup>6</sup>For head and dependent embeddings, we only considered neighbors that are RelEmbs to limit the variability.



model is not generalizing correctly.

A more complex scenario appears if we consider the targets generated from compositional bases. As noticed above, ANNE with CONCAT input has more shared neighbors between the target and the base (first two columns), while this is not the case for the ANNE with DIFF input (see heatmaps in (a) and (b) as example). This finding is further proof that a neural network that takes as input the concatenation of vectors for the base and the target argument attempts to generate an embedding as close as possible to the input relational embedding. In other words, this type of input could negatively impact ANNE in learning the correct inference.

### 4.3 Architectures' comparison

Previous analyses reveal that some parameter configurations are better than others. ANNE models that take as input the SIMPLE concatenation of the RelEmb base and the word embedding generate vectors too similar to the base, while modifying the base with an operation similar to Mikolov's vector offset produces better results. Overall, it seems that ANNE trained with MULTIHINGE and MULTIMSE losses (with DIFF input) induce more consistent and explainable results, while MULTIAVG is sub-optimal for its tendency to generate embeddings similar to the target's dependent.

## 5 General Discussion

An open issue in DSMs is how distributional representations can be projected from the lexical level to the sentence or even discourse level. Most previous approaches have tried to solve this issue by explicitly relying on the classic principle of compositionality. Given the Fregean assumption that phrase meaning is a function of the meanings of its constituents, different computational strategies have been proposed to derive vectors for phrases by taking word embeddings as inputs.

In this paper, we proposed a new methodology grounded on a usage-based perspective: we tried to generate new distributional representations by implementing an analogical function in the form of a neural network. Word analogies have been used as a standard intrinsic evaluation task for measuring the quality of word (Mikolov et al., 2013c; Levy and Goldberg, 2014; Linzen, 2016) and sentence embeddings (Zhu and de Melo, 2020; Wang et al., 2021; Ushio et al., 2021b). However, the task is usually defined as a candidate retrieval: given an

analogical proportion, find the correct completion from a list of candidates to solve the analogy. On the contrary, our aim is to generate a completely new embedding, similarly to what is done in reasoning and computer vision (Reed et al., 2015; Sadeghi et al., 2015; Upchurch et al., 2016; Ichien et al., 2021): the task consists in training deep learning models to recognize a relationship among two images and generate a transformed query representation (in this case, an image) accordingly. We believe that future investigations in linguistic analogies should pick up from this literature, and ANNE is a first attempt along this direction.

Our ANNE approach is not without limitations. One controversial aspect of ANNE is the choice of building the target by simply changing the argument in the relation. While it is not too problematic for verbal phrases, it raises questions for adjective-noun phrases and noun compounds. Consider the expressions *blue car* and *fast car*. Many things can be blue and not be a car, but not everything can be fast (e.g., *\*fast carrot*) because *fast* constrains the possible realizations of its head. A similar observation could be shown for noun compounds: in some cases, their meaning is related to both components (e.g., *bank account*), but sometimes their meaning retains aspects of one component (e.g., *head teacher*). To take into account the specificities of each type of phrases, we could train different ANNE architectures for each type of phrases.

The main difficulty is to balance relational and attributional similarity. Indeed, the use of a new item in a construction requires a great deal of relational knowledge (Gentner and Markman, 1997); nonetheless, the importance of similarity or shared attributes to linguistic analogy is not less vital (Bybee, 2010). A qualitative evaluation of analogical inferred embeddings reveals that analogy is easier to compute if the similarity between the entities in the syntactic relations is high. For instance, most all architectures build a good representation of *science technology* generated from *earth science*, maybe because there are lots of "topic science" expressions (cf. Table 5). Conversely, if attribute similarity is lower (i.e., the words between the base and the target are somewhat dissimilar), the analogical model is challenged. The neighbors of *pedantic circle* (derived from *literary circle*, cf. Table 6) are odd and incoherent with the expected meaning, maybe because the adjective *literary* is usually associated with a work of literature (an inanimate

	concat	diff	concat	diff
SIMPLE	<i>earth science</i> <i>apply science</i> <i>marine science</i> <i>new science</i> <i>area of technology</i>	<i>area of technology</i> <i>apply technology</i> <i>focus ORGANIZ.</i> <i>include technology</i> <i>area of engineering</i>	<i>show (a) letter</i> <i>explain in letter</i> <i>(a) disciple PERSON</i> <i>(a) letter address</i> <i>refer in (the) letter</i>	<i>guess PERSON</i> <i>extol (the) virtue</i> <i>point_out PERSON</i> <i>complain about PERSON</i> <i>dismiss (an) idea</i>
MULTIHINGE	<i>earth science</i> <i>new science</i> <i>apply science</i> <i>relate to technology</i> <i>area of technology</i>	<i>focus ORGANIZ.</i> <i>area of technology</i> <i>apply technology</i> <i>electronic technology</i> <i>create technology</i>	LOCATION <i>scholar</i> <i>join on return</i> <i>accompany (an) expedition</i> <i>await (the) return</i> <i>(a) letter address</i>	<i>state for example</i> <i>extol (the) virtue</i> <i>join on return</i> <i>serve curacy</i> <i>say in july</i>
MULTIMSE	<i>earth science</i> <i>apply science</i> <i>area of technology</i> <i>new science</i> <i>area of engineering</i>	<i>area of technology</i> <i>include technology</i> <i>focus ORGANIZ.</i> <i>aspect of technology</i> <i>aspect of use</i>	<i>show letter</i> <i>explain in letter</i> <i>letter address</i> <i>enlist aid</i> <i>(a) PERSON demand</i>	<i>complain about PERSON</i> <i>guess PERSON</i> <i>extol (the) virtue</i> <i>point_out PERSON</i> <i>say according to PERSON</i>
MULTIAVG	<i>apply science</i> <i>information technology</i> <i>development in science</i> <i>role of technology</i> <i>area of technology</i>	<i>information technology</i> <i>apply technology</i> <i>area of technology</i> <i>apply science</i> <i>new technology</i>	<i>see before PERSON</i> <i>like (one's) style</i> <i>tell (a) girl</i> <i>tell about time</i> <i>everyone tell(s)</i>	<i>see before PERSON</i> <i>feel like PERSON</i> <i>tell (a) girl</i> <i>realize PERSON</i> <i>want (a) baby</i>

Table 5: 5-nearest neighbors of *technology science* (compound) generated *earth science*.

Table 6: 5-nearest neighbors of *pedantic circle* (amod) generated from *a literary circle*.

object), while *pedantic* collocates with a person. In these cases, different factors could contribute to the success or failure of the model, which should be further investigated.

The introduction of analogy as a strategy to derive meaning for novel expressions does not entail the entire suppression of compositional approaches. From a theoretical stance, not every expression can be built using analogical inference: if analogy fails, compositional operations switch over to guide interpretation. In this regard, the question should not be whether analogically-generated vectors are better than computationally-built ones, but when one mechanism is preferred to the other. Answering this question is challenging from both a psycholinguistic and computational stance. The issues related to computational models of analogy as a productive mechanism in language are theoretical before methodological. While it is true that the cognitive process of analogy represents a central mechanism in human cognition (Hofstadter, 2001), the problem in defining a linguistic theory that formalizes precisely what an analogy is and when it occurs is complex. In other words, it is hard to predict which analogies will actually be drawn and at what linguistic level (Behrens, 2017, p. 215). Ideally, future systems aiming at modeling language comprehension should be able to include this mechanism too. New benchmarks will have to be built with the aim of identifying analogical inferences. These datasets could also be valuable for behavioral analyses.

## 6 Conclusion and Future Works

In this paper, we presented a new approach that simulates the construction of phrasal meaning as an analogical process implemented with deep learning techniques. We proposed a distributional representation of constructional phrases and a model of generating new embeddings analogically rather than applying traditional compositional operations. We experimented with our analogical neural network to understand how it can generalize and be extendable to different scenarios. We argued that the proposed methodology could open the doors to new analyses in distributional semantics as well as in computational models of language processing.

The future research perspectives on ANNE are considerable. Firstly, we could build a more sophisticated phrasal representation using contextualized embeddings (Ethayarajh, 2019) based on Transformers (Vaswani et al., 2017; Devlin et al., 2019). Moreover, we should compare our Relemb with other phrasal representations representations (Shwartz, 2019; Alipoor and Schulte im Walde, 2020) and Relation Embeddings (Camacho-Collados et al., 2019; Ushio et al., 2021a). Moreover, while we performed analogy over pre-selected base-target pairs, we aim at investigating methods to automatically retrieve the best analogical candidate. Finally, we plan to evaluate ANNE’s ability to model human behavior on more complex tasks regarding compositionality and language productivity.

## References

- Eneko Agirre, Enrique Alfonseca, Keith Hall, Jana Kravalova, Marius Pasca, and Aitor Soroa. 2009. A Study on Similarity and Relatedness Using Distributional and WordNet-based Approaches. In *Proceedings of NAACL-HLT*.
- Pegah Alipoor and Sabine Schulte im Walde. 2020. [Variants of vector space reductions for predicting the compositionality of English noun compounds](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4379–4387, Marseille, France. European Language Resources Association.
- Ben Ambridge. 2020. Against Stored Abstractions: A Radical Exemplar Model of Language Acquisition. *First Language*, 40(5-6):509–559.
- Marco Baroni, Raffaella Bernardi, and Roberto Zamparelli. 2014. Frege in Space: A Program of Compositional Distributional Semantics. *LiLT (Linguistic Issues in Language Technology)*, 9.
- Marco Baroni and Roberto Zamparelli. 2010. Nouns Are Vectors, Adjectives Are Matrices: Representing Adjective-noun Constructions in Semantic Space. In *Proceedings of EMNLP*.
- Heike Behrens. 2017. The Role of Analogy in Language Processing and Acquisition. *The Changing English Language: Psycholinguistic Perspectives*, pages 215–239.
- Elia Bruni, Nam-Khanh Tran, and Marco Baroni. 2014. Multimodal Distributional Semantics. *Journal of Artificial Intelligence Research*, 49:1–47.
- Joan L. Bybee. 2010. *Language, Usage and Cognition*. Cambridge University Press.
- Jose Camacho-Collados, Luis Espinosa Anke, and Steven Schockaert. 2019. [Relational word embeddings](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3286–3296, Florence, Italy. Association for Computational Linguistics.
- Jianpeng Cheng and Dimitri Kartsaklis. 2015. Syntax-aware Multi-sense Word Embeddings for Deep Compositional Models of Meaning. In *Proceedings of EMNLP*.
- Emmanuele Chersoni, Enrico Santus, Alessandro Lenci, Philippe Blache, and Chu-Ren Huang. 2021. Not All Arguments Are Processed Equally: A Distributional Model of Argument Complexity. *Language Resources and Evaluation*, 55(4):873–900.
- Emmanuele Chersoni, Enrico Santus, Ludovica Panitto, Alessandro Lenci, Philippe Blache, and C-R Huang. 2019. A Structured Distributional Model of Sentence Meaning and Processing. *Natural Language Engineering*, 25(4):483–502.
- Bob Coecke, Mehrnoosh Sadrzadeh, and Stephen Clark. 2010. Mathematical Foundations for a Compositional Distributional Model of Meaning. *arXiv preprint arXiv:1003.4394*.
- Silvio Cordeiro, Aline Villavicencio, Marco Idiart, and Carlos Ramisch. 2019. Unsupervised Compositionality Prediction of Nominal Compounds. *Computational Linguistics*, 45:1–57.
- William Croft. 1991. *Syntactic Categories and Grammatical Relations: The Cognitive Organization of Information*. University of Chicago Press.
- William Croft. 2001. *Radical Construction Grammar: Syntactic Theory in Typological Perspective*. Oxford University Press.
- Matthew Crump, Randall Jamieson, Brendan T Johns, and Michael N Jones. 2020. Controlling the Retrieval of General vs Specific Semantic Knowledge in the Instance Theory of Semantic Memory. In *CogSci*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of NAACL*.
- Kawin Ethayarajh. 2019. How Contextual Are Contextualized Word Representations? Comparing the Geometry of BERT, ELMo, and GPT-2 Embeddings. In *Proceedings of EMNLP*.
- Stefan Evert and Gabriella Lapesa. 2021. FAST: A Carefully Sampled and Cognitively Motivated Dataset for Distributional Semantic Evaluation. In *Proceedings of CONLL*.
- Lev Finkelstein, Evgeniy Gabrilovich, Yossi Matias, Ehud Rivlin, Zach Solan, Gadi Wolfman, and Eytan Ruppin. 2001. Placing Search in Context: The Concept Revisited. In *Proceedings of International Conference on World Wide Web*, pages 406–414.
- Dedre Gentner and Arthur B Markman. 1997. Structure Mapping in Analogy and Similarity. *American Psychologist*, 52(1):45.
- Talmy Givón. 1995. *Functionalism and Grammar*. John Benjamins Publishing.
- Anna Gladkova, Aleksandr Drozd, and Satoshi Matsuo. 2016. Analogy-based Detection of Morphological and Semantic Relations with Word Rmbeddings: What Works and What doesn’t. In *Proceedings of the NAACL Student Research Workshop*.
- Adele E. Goldberg. 2006. *Constructions at Work: The Nature of Generalization in Language*. Oxford University Press on Demand.
- Thomas Hoffman and Graeme Trousdale, editors. 2013. *The Oxford Handbook of Construction Grammar*. Oxford University Press, Oxford.

- Douglas R Hofstadter. 2001. Analogy as the Core of Cognition. *The Analogical Mind: Perspectives from Cognitive Science*, pages 499–538.
- Paul Ibbotson. 2013. The Scope of Usage-based Theory. *Frontiers in Psychology*, 4:255.
- Nicholas Ichien, Qing Liu, Shuhao Fu, Keith J Holyoak, Alan Yuille, and Hongjing Lu. 2021. Visual Analogy: Deep Learning versus Compositional Models. *arXiv preprint arXiv:2105.07065*.
- Randall K Jamieson, Johnathan E Avery, Brendan T Johns, and Michael N Jones. 2018. An Instance Theory of Semantic Memory. *Computational Brain & Behavior*, 1(2):119–136.
- Michael N Jones and Douglas JK Mewhort. 2007. Representing Word Meaning and Order Information in a Composite Holographic Lexicon. *Psychological Review*, 114(1):1.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980*.
- Nikolaus Kriegeskorte and Rogier A Kievit. 2013. Representational Geometry: Integrating Cognition, Computation, and the Brain. *Trends in Cognitive Sciences*, 17(8):401–412.
- Nikolaus Kriegeskorte, Marieke Mur, and Peter A Bاندettini. 2008. Representational Similarity Analysis—connecting the Branches of Systems Neuroscience. *Frontiers in Systems Neuroscience*, page 4.
- Ronald W Langacker. 1987. *Foundations of Cognitive Grammar: Theoretical Prerequisites*, volume 1. Stanford University Press.
- Alessandro Lenci. 2018. Distributional Models of Word Meaning. *Annual Review of Linguistics*, 4(1):151–171.
- Alessandro Lenci, Magnus Sahlgren, Patrick Jeuniaux, Amaru Cuba Gyllensten, and Martina Miliani. 2022. A Comprehensive Comparative Evaluation and Analysis of Distributional Semantic Models. *Language, Resources and Evaluation*.
- Omer Levy and Yoav Goldberg. 2014. Dependency-based Word Embeddings. In *Proceedings of ACL*.
- Maya R Libben and Debra A Titone. 2008. The Multi-determined Nature of Idiom Processing. *Memory & Cognition*, 36(6):1103–1121.
- Tal Linzen. 2016. Issues in evaluating semantic spaces using word analogies. In *Proceedings of the 1st Workshop on Evaluating Vector-Space Representations for NLP*, pages 13–18.
- Christopher D Manning, Mihai Surdeanu, John Bauer, Jenny Rose Finkel, Steven Bethard, and David McClosky. 2014. The Stanford CoreNLP Natural Language Processing Toolkit. In *Proceedings of ACL: System Demo*.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013a. Distributed Representations of Words and Phrases and their Compositionality. *Advances in Neural Information Processing Systems*, 26.
- Tomáš Mikolov, Wen-tau Yih, and Geoffrey Zweig. 2013b. Linguistic Regularities in Continuous Space Word Representations. In *Proceedings of NAACL-HLT*.
- Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. 2013c. Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 conference of the north american chapter of the association for computational linguistics: Human language technologies*, pages 746–751.
- Jeff Mitchell and Mirella Lapata. 2010. Composition in Distributional Models of Semantics. *Cognitive Science*, 34(8):1388–1429.
- Jeff Mitchell, Mirella Lapata, Vera Demberg, and Frank Keller. 2010. Syntactic and Semantic Factors in Processing Difficulty: An Integrated Measure. In *Proceedings of ACL*.
- Denis Paperno, Nghia The Pham, and Marco Baroni. 2014. A Practical and Linguistically-motivated Approach to Compositional Distributional Semantics. In *Proceedings of ACL*.
- Barbara H. Partee. 2004. *Compositionality in Formal Semantics: Selected Papers*. Explorations in Semantics. Blackwell Publishing Ltd.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An Imperative Style, High-performance Deep Learning Library. *Advances in Neural Information Processing Systems*, 32.
- Siva Reddy, Diana McCarthy, and Suresh Manandhar. 2011. An Empirical Study on Compositionality in Compound Nouns. In *Proceedings of IJCNLP*.
- Scott E Reed, Yi Zhang, Yuting Zhang, and Honglak Lee. 2015. Deep Visual Analogy-making. *Advances in Neural Information Processing Systems*, 28.
- Fereshteh Sadeghi, C Lawrence Zitnick, and Ali Farhadi. 2015. Visalogy: Answering Visual Analogy Questions. *Advances in Neural Information Processing Systems*, 28.
- Vered Shwartz. 2019. [A systematic comparison of English noun compound representations](#). In *Proceedings of the Joint Workshop on Multiword Expressions and WordNet (MWE-WN 2019)*, pages 92–103, Florence, Italy. Association for Computational Linguistics.
- Richard Socher, Brody Huval, Christopher D Manning, and Andrew Y Ng. 2012. Semantic Compositionality through Recursive Matrix-vector Spaces. In *Proceedings of EMNLP-CONLL*.

Michael Tomasello. 2009. The Usage-based Theory of Language Acquisition. In *The Cambridge Handbook of Child Language*, pages 69–87. Cambridge University Press.

Peter D Turney. 2012. Domain and Function: A Dual-space Model of Semantic Relations and Compositions. *Journal of Artificial Intelligence Research*, 44:533–585.

Paul Upchurch, Noah Snively, and Kavita Bala. 2016. From A to Z: Supervised Transfer of Style and Content Using Deep Neural Network Generators. *arXiv preprint arXiv:1603.02003*.

Asahi Ushio, Jose Camacho-Collados, and Steven Schockaert. 2021a. Distilling relation embeddings from pretrained language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.

Asahi Ushio, Luis Espinosa Anke, Steven Schockaert, and Jose Camacho-Collados. 2021b. BERT is to NLP what AlexNet is to CV: Can pre-trained language models identify analogies? In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Online. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. *Advances in Neural Information Processing Systems*, 30.

Yizhe Wang, Béatrice Daille, and Nabil Hathout. 2021. Caractérisation des relations sémantiques entre termes multi-mots fondée sur l’analogie (semantic relations recognition between multi-word terms by means of analogy). In *Actes de la 28e Conférence sur le Traitement Automatique des Langues Naturelles. Volume 1 : conférence principale*, pages 115–124. ATALA.

Xunjie Zhu and Gerard de Melo. 2020. Sentence analogies: Linguistic regularities in sentence embeddings. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3389–3400.

## A RelEmbs parameters

We use the skip-gram algorithm adopting the default configuration: no hierarchical softmax, 15 negative samples (how many negative contexts to sample for every correct one), and 300 as the vector dimension.

## B ANNE losses

The basic architecture (SIMPLE) uses the cosine similarity between original and predicted vectors

to make backpropagation. The CosineEmbeddingLoss<sup>7</sup> criterion from PyTorch library (Paszke et al., 2019) measures whether two inputs  $t$  and  $t'$  are similar or dissimilar using the cosine distance ( $cos$ ):

$$CEloss(t, t', y) = \begin{cases} 1 - cos(t, t') & \text{if } y = 1 \\ max(0, cos(t, t')) & \text{if } y = -1 \end{cases} \quad (1)$$

The loss function takes as inputs  $t$ ,  $t'$ , and a label tensor  $y$  containing values (1 or -1). For our purposes, we set  $y=1$ , so the loss is  $1 - cos(t, t')$ : The closer the cosine value to 1, the more the two inputs are similar, and then the loss is closer to 0. The optimization strategy is to minimize the cost function, that is, obtaining a loss value near 0 for all items in the training set.

The MULTI-criterion loss function is defined by the general formula:

$$loss_{multi} = CEloss(t, t') + g(CEloss(b, b_{head}), CEloss(t', t_{head})) + g(CEloss(b, b_{dep}), CEloss(t', t_{dep})) \quad (2)$$

where  $t$  stands for the vector originally attested in RelEmbs space and  $t'$  corresponds to the output vector generated by the network;  $b$  represents the analogical base vector,  $b_{head/dep}$  represents the vectors for the head and the dependent of the base (the same applies for  $t_{head/dep}$ ). Finally,  $g(\cdot)$  represents the function used to compare the phrase-argument similarity scores, which can be either the Mean Squared Error (equation 3), the mean of the scores difference (equation 4), or the hinge loss function (equation 5).

$$MSE(x, x') = (x - x')^2 \quad (3)$$

$$AVG(x, x') = mean(x - x') \quad (4)$$

$$HINGE(x, x') = max(0, x - x') \quad (5)$$

For each loss function, the cost derivative for the model’s parameters (weight matrices  $W_1$ , bias vector  $b_1$ ) is computed, and the appropriate parameters are updated through backpropagation.

<sup>7</sup><https://pytorch.org/docs/stable/generated/torch.nn.CosineEmbeddingLoss.html>

## C ANNE Test unseen preparation

We randomly selected 15 verbs, 15 nouns, and 15 adjectives attested in the RelEmbs vocabulary and we picked out from the original list all pairs in which one of these lemmas appeared. For instance, given the verb *study*, we saved in a separate file all pairs in which the verb occurs, such as  $\langle study_V, aspect_N, dobj \rangle$  (“to study the aspect”)  $\rightarrow \langle study_V, development_N, dobj \rangle$  (“to study the development”).

## D ANNE Training Setup

Given the possible combinations of input type (CONCAT and DIFF) and losses functions (SIMPLE, MULTIMSE, MULTIAVG, and MULTIHINGE), we trained eight different versions of ANNE. All models were trained using 5-cross validation for 10 epochs with the Adam (Kingma and Ba, 2014) gradient descent, using a batch size of 25. Hyperparameter values equal for all models. The training was performed on a TITAN Xp GPU (12gb).

## E Task1-Correlation of the component similarities of the base and the generated target

In order to visualize how these measures differ among architectures, we plotted the similarity scores using boxplots (Figure 4). Each subfigure represents the similarities computed over embeddings generated from a specific model architecture. The plot on the left refers to the RelEmb-head similarities; the plot on the right illustrates the RelEmb-dependent similarities. In each plot, we grouped boxplots for the type of base (idiomatic or compositional) and the syntactic type of phrase (verbal—VN—or nominal—NC). Finally, similarities are computed for both the base embedding (blue) and target embedding (orange).

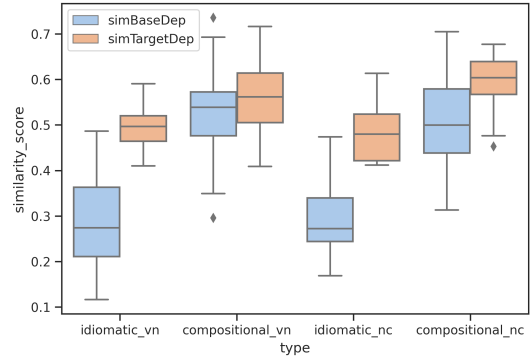
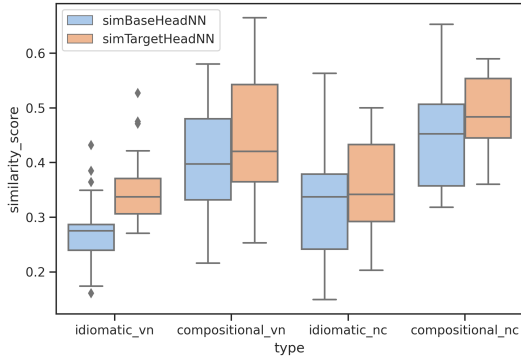
## F Task 2-Intersection of neighbors

We propose here a visual aid to investigate ANNE behavior. Figure 5 groups a series of heatmaps. In each heatmap, rows correspond to a specific item from the dataset, while columns represent the intersection between the neighbors of:

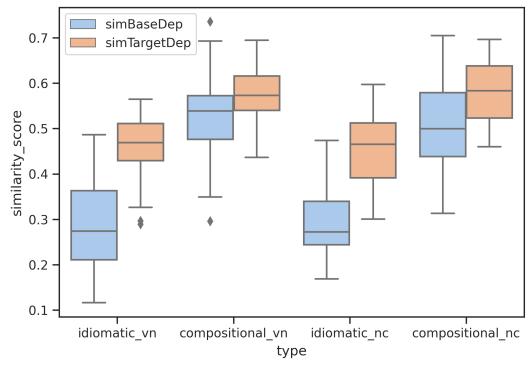
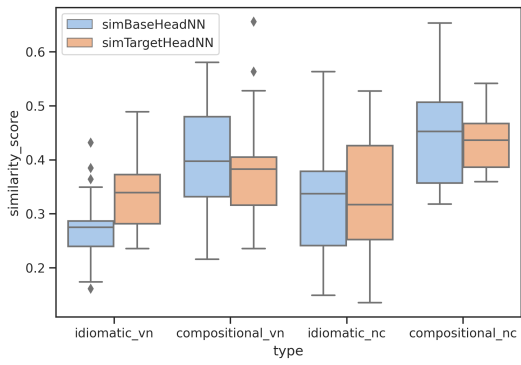
- The base and the generated target ( $NNbase\_NNtarg$ ;  $NNbase\_NNtarg_{re}$  considers only RelEmbs)

- The generated target and the respective head ( $NNtarg\_NNhead$ ) or dependent ( $NNtarg\_NNdep$ )—for these, we consider only RelEmbs neighbors.

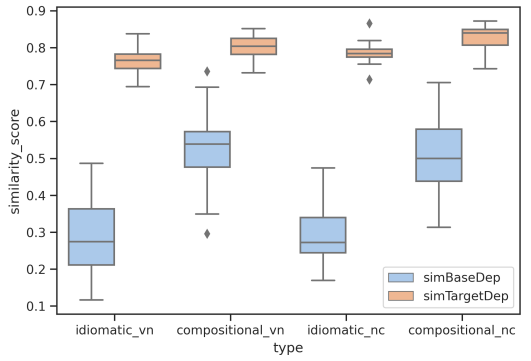
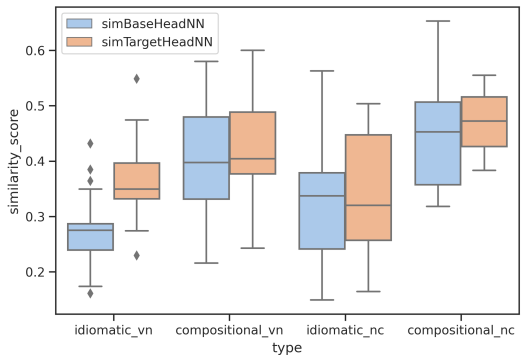
Numbers in the cells correspond to the number of neighbors retrieved. We present the results of analogical targets generated from a compositional (on the left) or idiomatic (on the right) base separately. Each subplot shows the results obtained for a specific model architecture.



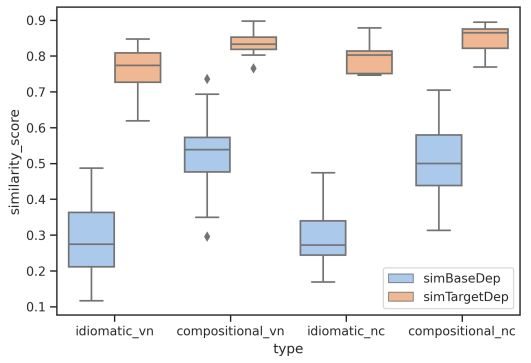
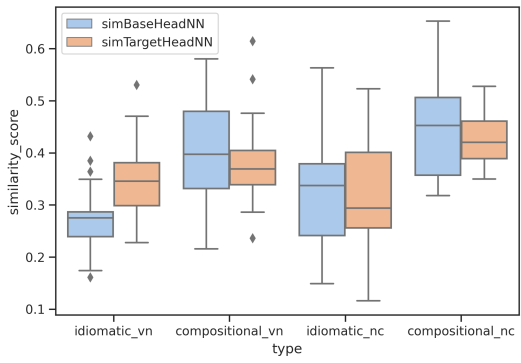
(a) SIMPLE with CONCAT input.



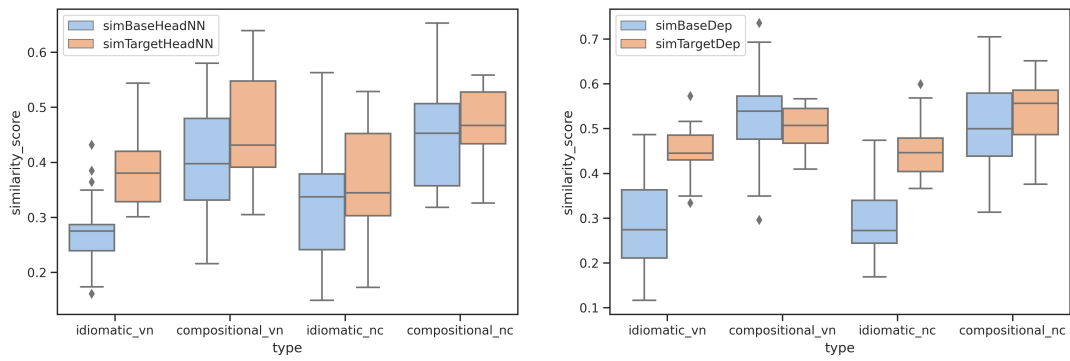
(b) SIMPLE with DIFF input



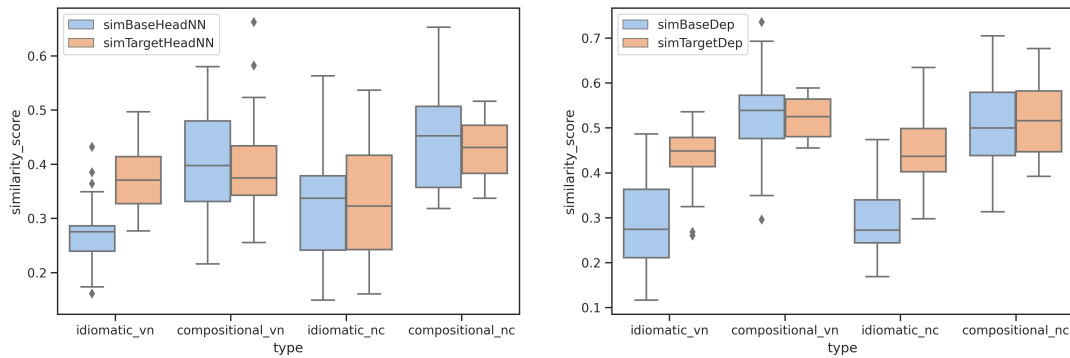
(c) MULTIAVG with CONCAT input



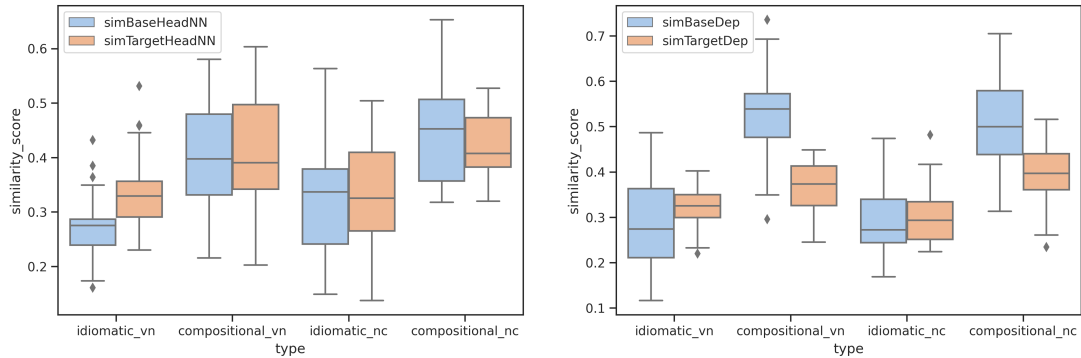
(d) MULTIAVG with DIFF input



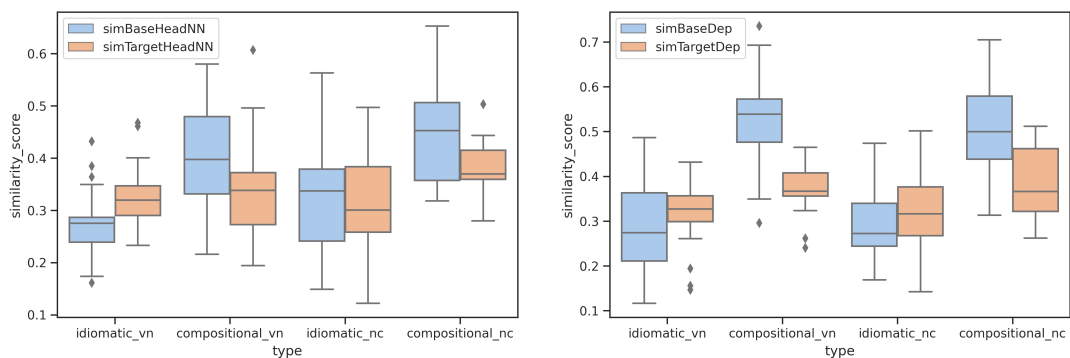
(e) MULTIMSE with CONCAT input



(f) MULTIMSE with DIFF input



(g) HINGE with CONCAT input

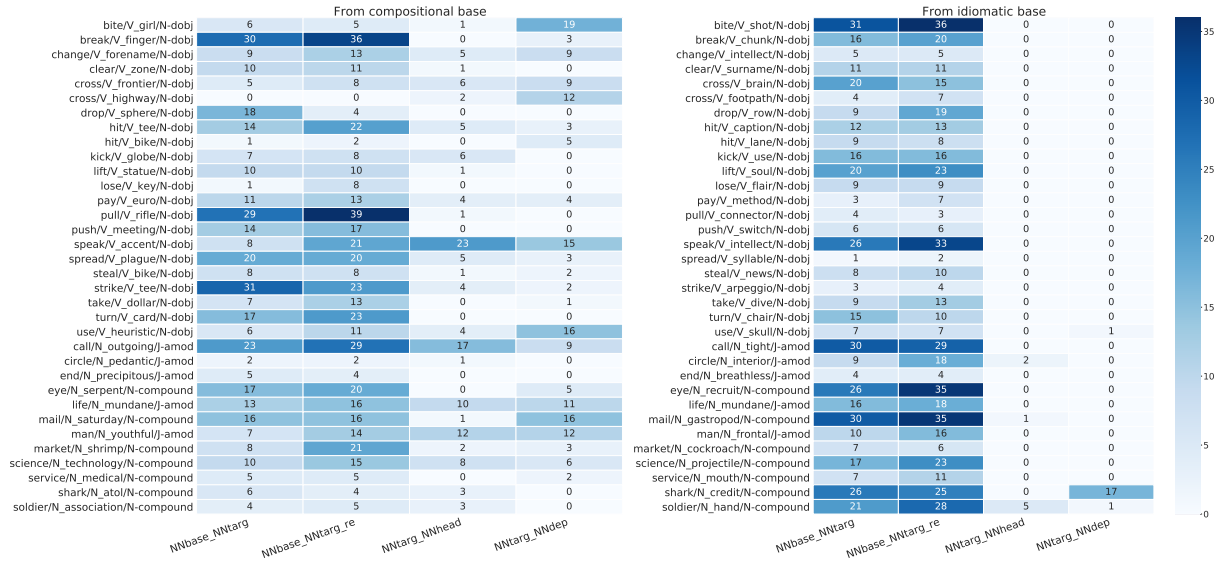


(h) HINGE with DIFF input

Figure 4: Distribution of the similarities between the ReLemb and its head (left), between the ReLemb and its argument (right). Data are grouped for syntactic type (nominal, NC, or verbal, VN) and if it is compositional (*compos*) or idiomatic (*idiom*). Similarities are computed for both the base embedding (blue) and target embedding (orange).

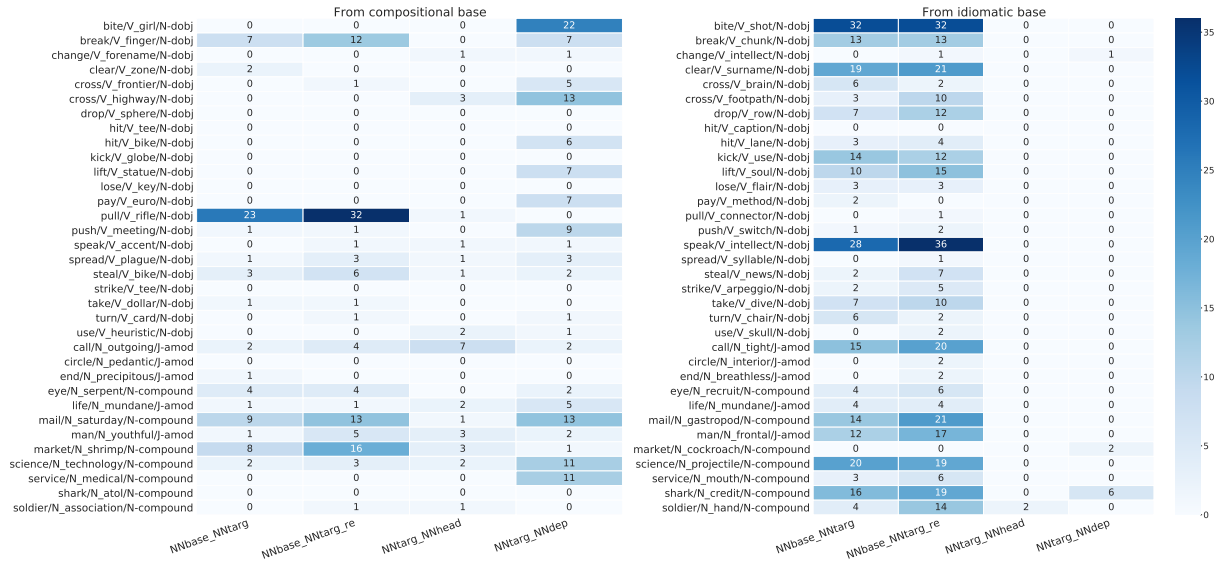


simple\_concat model



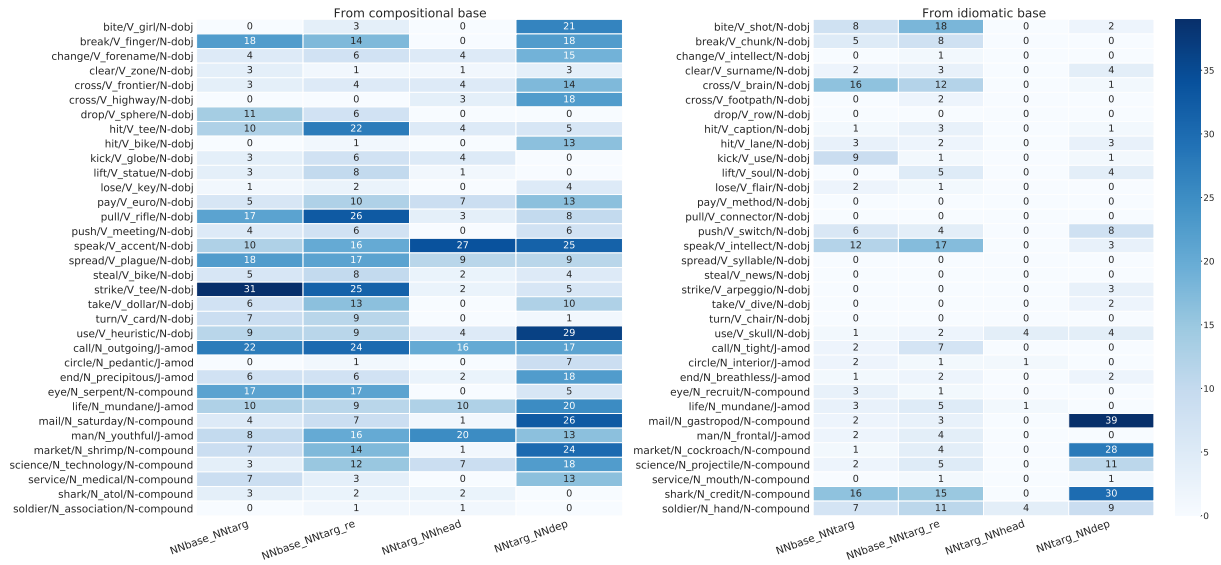
(a)

simple\_diff model



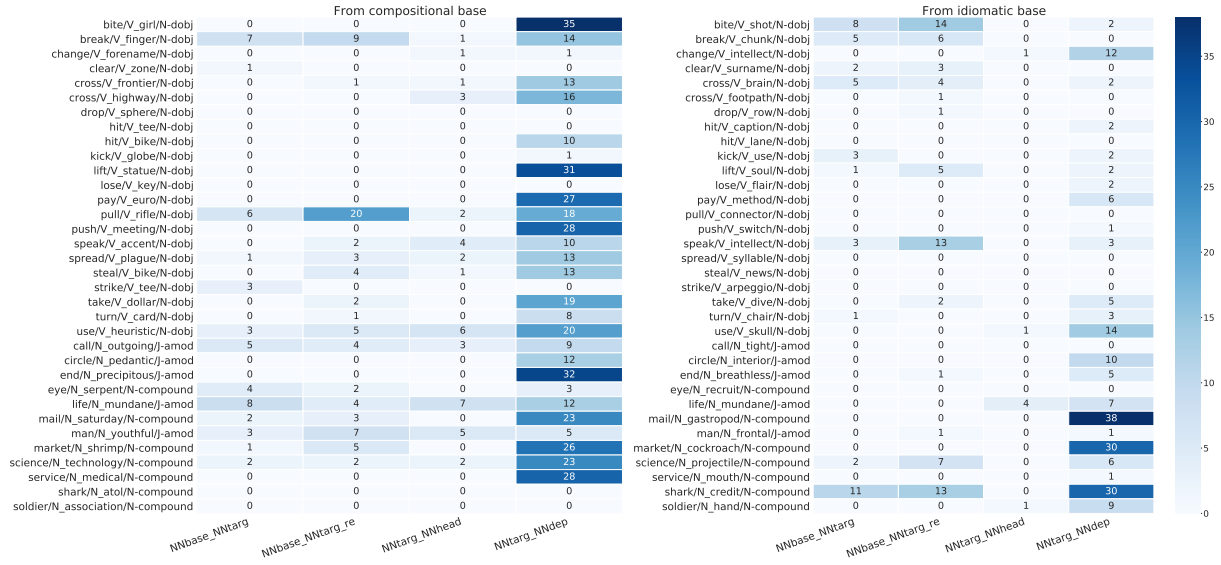
(b)

multiAVG\_concat model



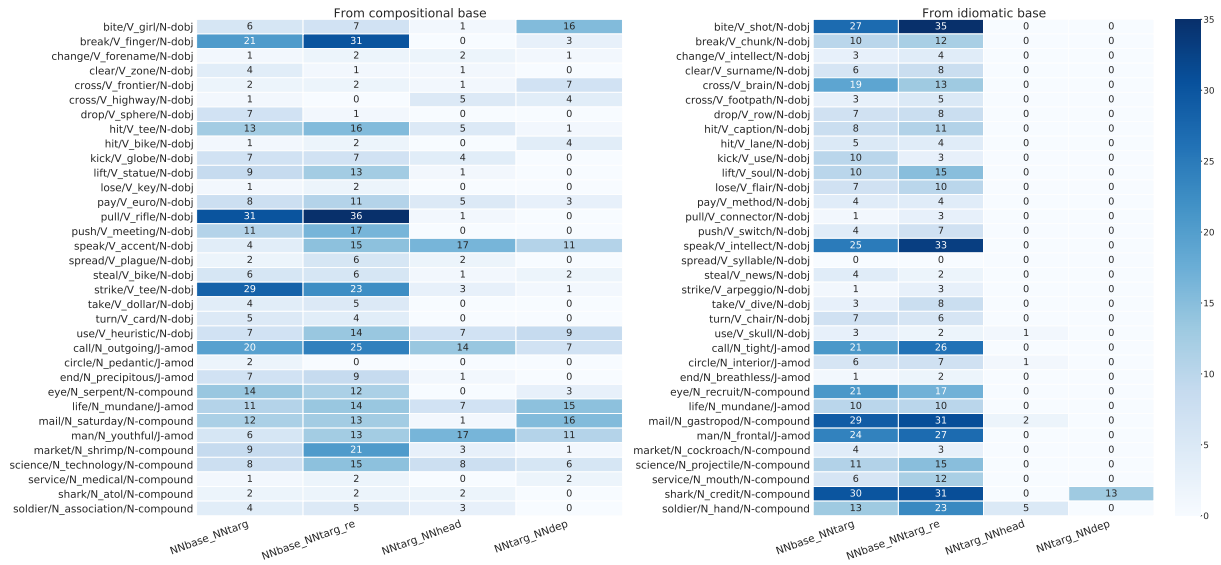
(c)

multiAVG\_diff model



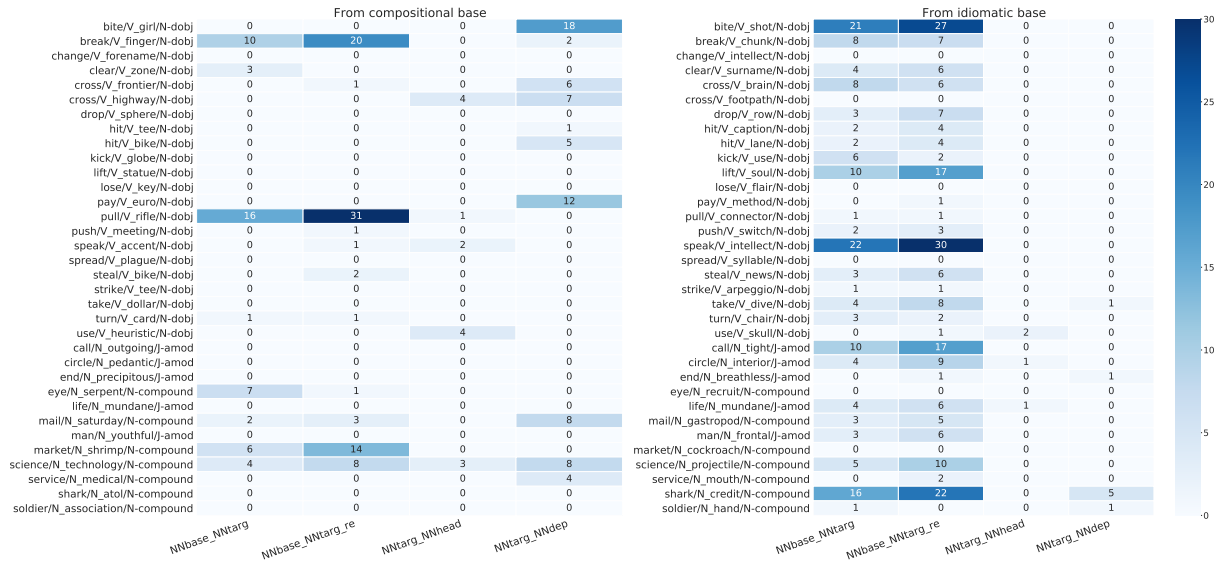
(d)

multimSE\_concat model

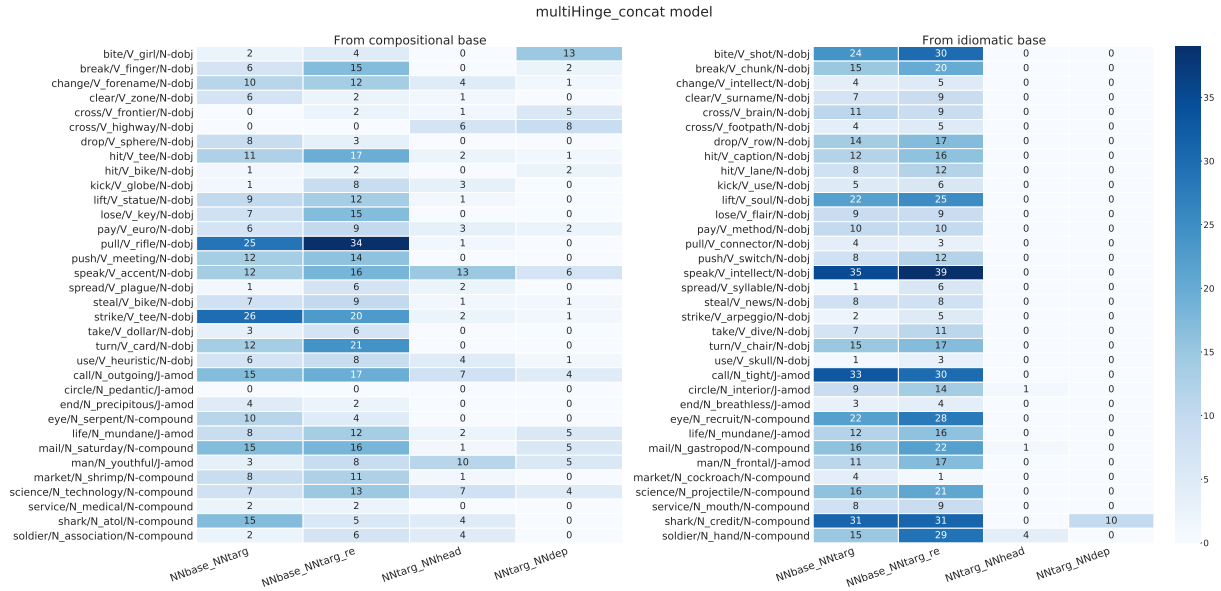


(e)

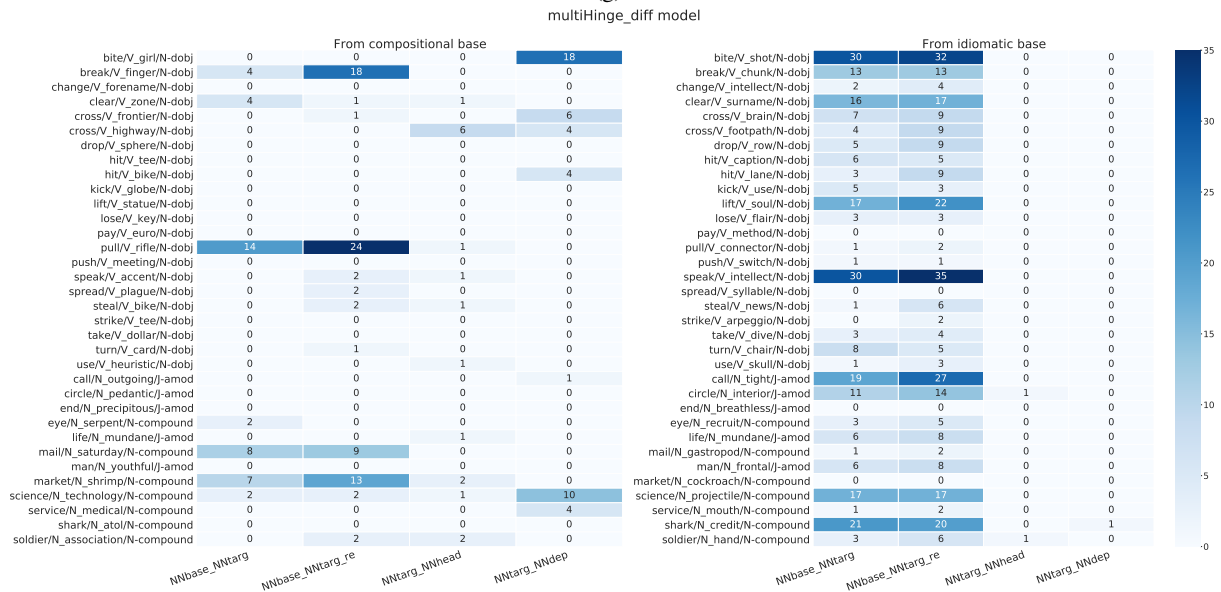
multimSE\_diff model



(f)



(g)



(h)

Figure 5: Heatmaps showing the intersection of common neighbors. Plot on the left refers to the target computed from a compositional base, plot on the right shows results for vectors generated from idiomatic base.