

What does memory retrieval *leave on the table*? Modelling the Cost of Semi-Compositionality with MINERVA2 and sBERT

Sydelle de Souza^{1,3,4*} Ivan Vegner^{1,3*} Francis Mollica^{2‡}
Leonidas A.A. Doumas^{4‡}

¹School of Informatics, University of Edinburgh ²University of Melbourne

³Centro de Linguística da Universidade de Lisboa

⁴School of Philosophy, Psychology & Language Sciences, University of Edinburgh

Abstract

Despite being ubiquitous in natural language, collocations (e.g., *kick+habit*) incur a unique processing cost, compared to compositional phrases (*kick+door*) and idioms (*kick+bucket*). We confirm this cost with behavioural data as well as MINERVA2, a memory model, suggesting that collocations constitute a distinct linguistic category. While the model fails to fully capture the observed human processing patterns, we find that below a specific item frequency threshold, the model’s retrieval failures align with human reaction times across conditions. This suggests an alternative processing mechanism that activates when memory retrieval fails.

1 The Curious Case of Collocations

From *killing time* and *playing dead* to *running baths* and *making beds*, word combinations with semi-compositional meanings are ubiquitous in human language (Cowie, 1998). Often referred to as *collocations*, these idiosyncratic lexical elements comprise one word used in its literal sense and another in its figurative sense, constrained by an arbitrary restriction on substitution (Mel’čuk, 2003; Howarth, 1998). Thus, one can *raise questions* or *lift bans*, but neither *#lift questions* nor *#raise bans*. Collocations are syntactically well formed, but deviate from or violate the expected semantic representation (Culicover et al., 2017). To illustrate, the verb *kill* prototypically requires an animate object, so one can *kill bugs* and *kill trees*, but not **kill books*. Yet one can *kill time*, *hope*, and *dreams*. Collocations are the largest subset of formulaic language (Barfield and Gyllstad, 2009) with many being cross-linguistically attested (Yamashita, 2018). It is hardly surprising, then, that proper knowledge and use of such units provides fluency and idiomaticity to the language user (Pawley and Syder,

1983; Durrant and Schmitt, 2009). Yet, they pose an enormous hurdle to second-language learners and machines.

According to Howarth (1998), human language lies on a theoretical continuum of semantic compositionality—the degree to which the meaning of a phrase can be derived from the meaning of its constituent parts and their syntactic relations (Frege, 1892). Fully compositional combinations (e.g., *chase rabbits*, *chase thieves*, etc.) and fully non-compositional figurative idioms (e.g., *chase one’s tail*)¹ lie on extreme ends of the spectrum. Semi-compositional collocations (e.g., *chase dreams*, *chase money*, etc.) lie in between. The psychological validity of this continuum has been tested with the expectation that a decrease in compositionality is directly proportional to a decrease in processing time (Gyllstad and Wolter, 2016). However, empirical evidence from both first (L1) and second (L2) language speakers shows that collocations are processed slower and less accurately than fully compositional combinations (Gyllstad and Wolter, 2016; de Souza et al., 2024), and fully opaque and non-compositional figurative idioms (e.g., *break the ice*) are processed faster and more accurately than compositional combinations (e.g., *break the cup*) (Carrol and Conklin, 2020; Tabossi et al., 2008).

These disparities are also seen in evidence from L2 acquisition². Research shows that colloca-

¹It is important to note that (Howarth, 1998) also specifies a fourth category called "pure idioms" (e.g., *blow the gaff*, *take a leak*, *shoot the breeze*). These do not possess well-specified literal meanings (see Mueller and Gibbs, 1987, for further reading) and comprise a very small subset of formulaic language occurring quite infrequently (Grant, 2005). Furthermore, most of the studies in this area focus on figurative idioms that have an additional literal reading (e.g., *kick the bucket*). Therefore, in order to constrain the scope of this paper, we limit our discussion to figurative idioms.

²Note that collocations as a distinct linguistic class have largely been ignored in L1 acquisition research and little is known about the mechanisms behind a child’s acquisition of

*Joint first authors ‡ Joint senior authors
Correspondence to: sydelle.desouza@ed.ac.uk

tions are a major hurdle for second language (L2) learners—be they early sequential bilinguals (Nishikawa, 2019; Riches et al., 2022) or adults (Yamagata et al., 2023; Sonbul et al., 2024), even at high proficiency levels (Wolter and Gyllstad, 2013; Tsai, 2020). In contrast, idioms are learned better and used more accurately than collocations (Fioravanti et al., 2021). Cast under the broader term of *conceptual metaphor* (Lakoff and Johnson, 1980), collocations are also found to be challenging for NLP systems (Liu et al., 2022; Zayed et al., 2018; Czinczoll et al., 2022) despite the fact that the last decade has seen immense progress (see Tong et al., 2021; Wang et al., 2025).

The coalescing picture suggests that idioms are processed the fastest, followed by compositional units, and collocations the slowest. However, this processing hierarchy has not yet been directly tested in a within-participants design—a gap the present study addresses. Next, we discuss how this empirical gap is compounded by a theoretical gap.

2 Accounting for Collocation Processing

It is generally agreed upon in the language processing literature that idioms are stored and retrieved from memory holistically (Carroll and Conklin, 2014; Noveck et al., 2023). Although there are several theories concerning the processing of compositional language, there is less consensus on the matter. Being not too compositional and not too idiomatic, collocations have occupied a theoretical grey zone in mainstream psycholinguistics wherein they have been conveniently ignored in favour of a binary distinction between rules and exceptions. As a result, researchers in second language acquisition and applied psycholinguistics have drawn on the (in)famous Past Tense Debate (Seidenberg and Plaut, 2014) in morphological processing and resorted to single- versus dual-route models to explain processing at the multi-word level (Wray, 2002). We explore these models and consider their ability to account for collocational processing.

Assuming a domain-general hypothesis space, **single-route models** posit that all linguistic forms are stored in and retrieved from a single massive associative memory system³ based on frequency of input and use (Bybee, 2012; Ambridge and Lieven, 2011). The more often a unit is encountered and/or

collocations over development (see Handl and Graf, 2010).

³Or that all forms are processed equally as in a connectionist network (see McClelland and Rumelhart, 1985).

used, the better it is entrenched in memory (Divjak, 2019; Langacker, 1987). Eventually, this leads to automatization—pure retrieval from memory⁴ (Bybee, 2006) which makes processing fast and effortless. Positing such a homogenous mechanism makes for a parsimonious theoretical account of our language abilities, in particular, and our cognition in general. However, human memory is not only limited in capacity (Christiansen and Chater, 2008) but is also unstable (Kornell and Bjork, 2009). We do not store everything we encounter, nor do we remember everything we do store. More importantly, recall that behavioural evidence points to collocations incurring a processing cost versus compositional units even when frequency-matched (see de Souza et al., 2024). While memory undoubtedly plays an important role in language processing, it does not provide a satisfactory account for the processing cost of collocations which occur quite frequently (Barfield and Gyllstad, 2009).

The **dual-route model** assumes a domain-specific hypothesis space, differentiating between words and rules (Pinker, 1991). Regular word forms are thought to be computed analytically (e.g., *walk* → *walk* + *ed*, *scratch* → *scratch* + *ed*) by way of rules, while irregular word forms (e.g., *run* → *ran*, *think* → *thought*) are processed via holistic storage and retrieval from memory (Pinker, 2013).

This theoretical distinction between computation and storage is a practical trade-off between two independent cognitive processes—procedural computation and declarative memory (Pinker and Ullman, 2002). More rule-based computation means less storage. More storage means less computation. Positing such a heterogenous mechanism makes for a persuasive theoretical account of how human language can be infinitely compositional despite our limited cognitive capacities (O'Donnell et al., 2009; Galke et al., 2024). The dual-route explanation is used to account for formulaic language processing as a whole, i.e., it does not distinguish between the various subsets of multi-word units such as idioms, phrasal verbs, binomials, etc. (see Wray, 2002, 2008; Sidtis, 2020). All formulaic language is thought to be stored, while compositional language is computed on the fly. Memory retrieval is faster than analytic processing (Logan, 1997; Dasgupta and Gershman, 2021), therefore formulaic language is thought to be processed faster

⁴See Logan and Etherton (1994) for a domain-general cognitive account of automatization.

than non-formulaic language (Carrol and Conklin, 2014; Vilkaite and Schmitt, 2019). This is empirically consistent across a variety of tasks only in the case of fully non-compositional units like idioms (Noveck et al., 2023). However, dual-route hypotheses make a binary distinction between compositional and formulaic language and ignores the effect of frequency on computation and retrieval. If collocations are frequent and retrieved from memory, the processing cost they incur remains unaccounted for.

3 The Present Study

We begin by addressing the empirical gap laid out in Section 1 and test whether collocations incur a processing cost relative to idioms and fully compositional phrases, as suggested by prior literature. To this end, we ask: *Do collocations take longer to process than idioms and compositional items?* We extend de Souza et al. (2024) by testing L1 English speakers on an acceptability judgement task (AJT) using stimuli from all three conditions and analyse reaction times (RTs) and accuracy. We consider three competing predictions:

Under a *single-route* account, frequency effects should dominate: idioms should be retrieved fastest, followed by collocations and then compositional items. Under a *dual-route* account, idioms and collocations—both stored, familiar units—should be processed similarly and faster than compositional phrases, which require computation. However, given prior findings, we predict that human participants will process idioms fastest, followed by compositional phrases, with collocations being the slowest—even though collocations are often more frequent than compositional items in our dataset (see Appendix A).

Furthermore, based on the review in Section 2, it would be uncontroversial to say that memory is critical to all forms of language processing (see also Divjak, 2019; Divjak et al., 2022; Corballis, 2019). It encapsulates single-route processes and is an integral component of dual-route models. Building on this foundational role of memory, we simulate memory retrieval using a well-established frequency-based mechanistic model of memory—MINERVA2 (Hintzman, 1984), modified to incorporate two key factors known to influence collocational processing: frequency (Wolter and Gyllstad, 2013) and semantics (Gyllstad and Wolter, 2016; Fioravanti et al., 2021). We adopt a distribu-

tional semantic framework (Landauer and Dumais, 1997; Mikolov et al., 2013) and use contextualized embeddings from Sentence-BERT (Reimers and Gurevych, 2019). We modify MINERVA to simulate RTs and load its memory according to the frequency of the stimuli in the corpus. We explore successful and failed retrievals to assess their influence on the processing signatures of different item conditions under a pure memory-based model. Our central research question is:

To what extent can pure memory retrieval account for processing differences observed in L1 English speakers across idioms, collocations, and compositional phrases?

Here, we expect MINERVA2 to show differences across conditions primarily as a function of frequency, with no added processing cost for collocations. In any case, if retrieval alone were sufficient to account for human processing trends, the model should mirror the human patterns consistent with a single-route account. If not, the model will allow us to probe what aspects of human performance can be explained by memory alone, and where memory-based retrieval may fall short.

4 Collocations Incur a Processing Cost

4.1 Methodology

Stimuli de Souza et al. (2024) introduced a stimulus set consisting of 100 Verb-Noun collocations (e.g., *spill secrets*) and 100 compositional Verb-Noun combinations containing the same verb as the collocation (e.g., *spill water*). We attempted to augment this stimulus set with a matching figurative idiom (e.g., *spill the beans*) for each verb with the help of the ‘word sketch’ function in The Sketch Engine’s enTenTen21 corpus (Kilgarriff et al., 2024). However, we were only able to identify idioms for 82 verbs in the dataset resulting in a final dataset of 246 target items (1 collocation, one composition, and one idiom for each of the 82 verbs). 82 baseline items, nonsense Verb-Noun combinations (*fry knob*), were created to use as distractors in the experiment. The dataset was divided into 3 folds of 82 items wherein no two items had the same verb. As expected, there are statistically significant differences between the mean frequencies of all three constructions with idioms being the most frequent, followed by collocations and compositional items being the least frequent group (see Appendix A for more details). We account for this discrepancy by including frequency as a

covariate in our statistical models.

Participants & Task A total of 186 L1 English speakers ($F = 112$; $M = 71$; $NB = 3$) were recruited using Prolific. They were remunerated £1.50 for their participation. The mean age of the sample was 38.6 years ($SD = 10.81$). They were asked to judge whether or not the word combination presented to them sounded acceptable (i.e., would they as L1 English speakers use this word combination in their everyday speech). They were asked to respond as quickly and accurately as possible, by pressing the ‘y’ key for yes or the ‘n’ key for no. During testing, each participant saw 164 items: 82 target items and 82 distractors. Items were presented in an individualized random order. A fixation cross with an inter-stimulus interval of 350 ms was presented between trials. Trials timed out at 8,000 ms if no decision was taken.

Data Pre-processing Data pre-processing was carried out using R version 4.4.1 "Race for Your Life" (R Core Team, 2024). Due to an error in data collection, data of four participants were replaced. We also remove all incorrect trials for reaction time analyses (2, 752; see Appendix B). We then eliminated responses below 450 ms and responses over 3.5 standard deviation from the grand mean including time-outs. These outliers accounted for 1.484% of the total data ($n = 30, 504$ including distractors). In terms of accuracy, all participants scored above 50%. However, we found 4 items with a mean accuracy of less than 50%. We eliminated those items along with other items that comprised the same verbs from our analyses. We do not analyse distractors (15, 252). All reaction time (RT) analyses are conducted on this final dataset ($n = 13, 369$).

4.2 Statistical Modelling

We first specified a maximal model as “justified by the design” (Barr et al., 2013). The main dependent variable was the reaction times (RTs) from the acceptability judgement task while the main predictor variable was Condition (Compositional, Collocation, Idiom; treatment coded, with idiom as the reference level). Phrasal Frequency (scaled) was included as a covariate. The maximal converging random effect structure included intercepts for Participant and Verb. The analysis model in R syntax specified using the ‘lme4’ (Bates et al., 2015) package is: $RT \sim Condition + Phrasal\ Frequency + (1 | ID) + (1 | Verb)$.

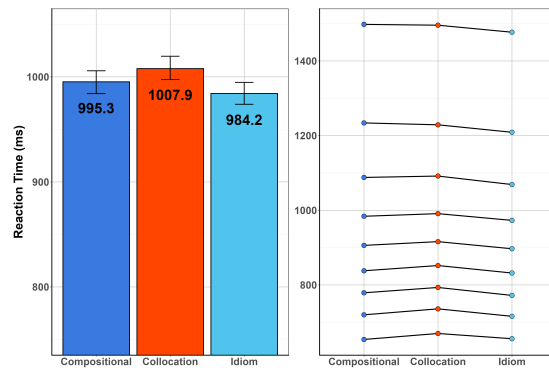


Figure 1: **Left:** mean reaction times (ms) by condition. Error bars indicate bootstrapped confidence intervals. **Right:** decile plot of reaction times by condition. Note the differences in the y-axes.

4.3 Results

Figure 1 shows the mean reaction times (RTs) by condition, as well as a breakdown by decile. Collocations have the slowest responses with a mean of 1007.87 ms ($SD = 370.84$ ms) compared to compositional items (995.32 ms, $SD = 375.76$ ms) and idioms (984.20 ms, $SD = 365.39$ ms).

Our statistical results showed a small, significant difference in RTs between compositional items and idioms ($\beta = 4.69$; $SE = 2.240$; $p = 0.037$), suggesting that compositional units were processed slower than idioms. A larger difference was found between collocations and idioms ($\beta = 13.80$; $SE = 1.760$; $p < 0.001$), replicating the processing costs predicted by the literature. Unsurprisingly, Phrasal Frequency also has a significant effect on RTs ($\beta = -18.50$; $SE = 1.640$; $p < 0.001$), corresponding to a 18.5 ms decrease in RT for every 1 standard deviation increase in phrasal frequency. In terms of accuracy, we found no significant difference between idioms and compositional items, but we do see a marginal difference ($p = 0.04$) between idioms and collocations. This is expected as all stimuli are highly frequent and should be familiar to adult L1 speakers. See Appendix B for detailed results.

5 Failures in Memory Retrieval Capture Behavioural Trends

As a first step toward elucidating the cognitive mechanisms underlying the processing trend that humans display across the compositionality continuum, we investigate the extent to which we can account for the trend with memory retrieval alone.

MINERVA is an instance-based model of

episodic memory that has been successfully applied to many cognitive phenomena from frequency judgements (Hintzman, 1988) to false memory (Arndt and Hirshman, 1998). It has also been used to model artificial grammar learning (Jamieson and Mewhort, 2009) and, recently, to metaphor recognition (Nick Reid and Jamieson, 2023).

MINERVA’s core assumptions are: (i) every item encountered leaves a memory trace, represented as a distributed set of features, and (ii) similar items have similar traces. Similarities between present and past encounters drive item-specific and parallel memory retrieval. As a global memory model, it encapsulates both episodic and semantic memory which communicate with each other. On encountering a stimulus, the episodic memory sends a probe to the semantic memory to retrieve traces from past encounters. The familiarity of the probe is then calculated as the sum of the values of a similarity measure between the probe and each stored trace.

MINERVA is instantiated in a linear algebra system. The MINERVA memory \mathbf{M} is an $n \times d$ matrix, each row of which contains a d -dimensional memory trace vector. When cued for retrieval with a probe $p \in \mathbb{R}^d$, MINERVA retrieves the representation of the probe iff the probe’s familiarity f is greater than a threshold $K \in [0, 1)$. Familiarity is calculated by taking the cosine similarity s of the probe to all instances stored in memory, scaling s to reflect activation (weighting) of memory items a over elapsed time τ , and linearly combining instances in memory to compute a memory echo e . The familiarity score at timestep τ is the cosine similarity of the echo to the probe, following this system of equations:

$$s = \text{sim}(p, \mathbf{M}) \quad (1)$$

$$a_\tau = s^\tau \text{sign}(s) \quad (2)$$

$$e_\tau = a_\tau \mathbf{M} \quad (3)$$

$$f_\tau = \text{sim}(e_\tau, p) \quad (4)$$

Modelling AJT Responses with Taus (τ) The free parameter τ is used to accentuate differences in similarity values (Hintzman, 1988; Nick Reid and Jamieson, 2023). By raising the value of τ , higher-similarity memory traces will elicit exponentially more activation, allowing those traces to play a larger role in the overall activation profile versus pooling a potentially large number of low-similarity items.

Following Nick Reid and Jamieson (2023), we depart from prior work wherein τ is kept constant

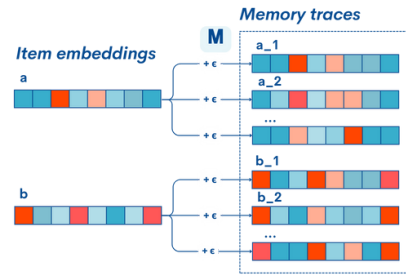


Figure 2: Illustration of how embeddings are noised and loaded into MINERVA’s memory matrix \mathbf{M} . Colors depict values within a vector. Note that the noise vectors ϵ are independently sampled for each memory trace.

for a particular experiment and model reaction times by dynamically increasing τ for a particular probe p until a desired threshold of familiarity $K \in [0, 1)$ is reached. At this point, we take the final value of τ as a proxy for the time required to recognize p from memory, i.e. a proxy for reaction time (RT). We set a time-out at $\tau = 300$ after which the next probe is presented.

In human acceptability judgements, reaction times serve as a proxy for processing difficulty. We implicitly model acceptability judgements in MINERVA as a function of whether the familiarity threshold K is reached within the allowable time window. If the familiarity score surpasses K before the time-out, i.e., successful recognition, we treat this as a "yes". Conversely, if familiarity remains below the threshold when $\tau = 300$, we treat the failure to retrieve as a "no" response.

5.1 Motivations & Assumptions

Collocational processing is known to be driven by two factors: *semantic transparency* and *frequency* (see Gyllstad and Wolter, 2016; Fioravanti et al., 2021). Our model captures semantic transparency by means of distributional semantics, i.e. vector embeddings, while frequency is captured by means of phrasal frequency in a dynamic web corpus. We demonstrate the effect of both factors in our ablations (see Section 5.4).

Semantics of Memory Traces Using distributed vector representations as memory traces for MINERVA is well-established in the literature (Chubala and Jamieson, 2013; Jamieson et al., 2018; Nick Reid and Jamieson, 2023). Given that the figurative idioms (e.g., *spill the beans*) also have a compositional reading, we need a contextualized, fine-grained vector representation to capture the

semantics of each word combination. Therefore, we rely on Sentence-BERT (sBERT) which provides semantically meaningful vector embeddings for sentences (Reimers and Gurevych, 2019). To derive the vector embedding for each of the 246 target stimuli, we follow Vulić et al. (2020). First, we collect a set of 100 sentences of the word combination⁵ from the enTenTen21 corpus, in which the noun occurs as the direct object of the verb. We feed each sentence to sBERT obtaining a set of contextualized word embeddings representing each word in the sentence (we perform mean pooling over sub-words). Given that the higher layers of BERT architectures are the most sensitive to lexical semantics (Reif et al., 2019), we take our embeddings from the last hidden layer of the model. From each of the 100 sentences, we extract the embeddings corresponding to the verb and the noun and average across them separately, resulting in the mean contextualized representation of the verb when paired with the noun, and of the noun when paired with the verb. Finally, we concatenate the mean embedding for the verb with the mean embedding for the noun to form the vector representation of our stimulus⁶.

Memory Frequencies & Forgetting In accordance with the instance theory, MINERVA’s retrieval time is inversely proportional to the number of memory traces that strongly respond to a particular probe (Nick Reid and Jamieson, 2023). Therefore, we populate MINERVA’s memory matrix using 10,000 items sampled proportionally to their phrasal frequency. Following prior work, we simulate forgetting by adding zero-centered Gaussian noise to each memory trace vector such that each dimension of each trace has an independent probability $F \in [0, 1)$ of being corrupted with noise. The more frequent a particular item, the more traces it will have in memory, averaging out the noise and making high-frequency items easier to retrieve.

5.2 Simulations

To explore the extent to which simple memory retrieval is sufficient to reproduce processing trends for each condition, we load the memory matrix

⁵Distractor items were not included in the simulations as they are nonsense combinations, have no context sentences and would have very low frequency in MINERVA’s memory.

⁶We use concatenation instead of mean pooling as our stimuli are all Verb + Direct Object and concatenation preserves word order and therefore, syntactic role information. However, see Appendix E.

as described above (see Figure 2) and test MINERVA’s recognition capabilities using a noiseless vector embedding of the target stimulus as the probe. To simulate N different participants who are exposed to different samplings of items from the same environmental distributions, as well as different patterns of forgetting, we run each simulation $N=300$ times with different random seeds, re-sampling and re-noising the memory matrix each time. We perform a thorough hyperparameter sweep of activation threshold K and forgetting probability F . We discuss results for hyperparameter values $K=0.99$ and $F=0.8$, although our results are robust across many hyperparameter combinations (see Figure 10).

We use the same statistical model described in Section 4 to analyse the effect of semantics and frequency on retrieval (i.e., Tau).

5.3 Results

The results of our computational experiment are shown in Figure 3. As MINERVA was not presented with any baseline items and as all items were in MINERVA’s memory, it should have succeeded at recognizing all items (Figure 3, left panel). Thus, we first considered only successful retrievals. Despite being provided with meaningful embeddings and frequencies, the model failed to capture human processing trends. Collocations were retrieved faster than idioms ($\beta = -0.41$; $SE = 0.004$; $p < 0.001$) while compositional items were retrieved slower than idioms ($\beta = 0.62$; $SE = 0.004$; $p < 0.001$). See Appendix C for more details. Given the surprising results, we analyzed the model’s failures to retrieve, i.e., timeouts, on every item (see Figure 3, right panel). MINERVA timed out on 50% of the retrievals for collocations, followed by compositional items (38.6%), with idioms timing out the least (33.8%). A mixed-effects logistic regression confirmed all differences between conditions to be statistically significant⁷ (see Table 6).

⁷To rule out the possibility that these results are a quirk of the MINERVA architecture, we also ran memory retrieval simulations using the same memory matrix on the Modern Hopfield Network (Ramsauer et al., 2021, MHN). The MHN is a generalization of the classical Hopfield network (Hopfield, 1982) adapted to work with continuous states, and is formally connected to the QKV attention mechanism in Transformers. We find that the MHN displays the same characteristic pattern of failures in at least one configuration of hyperparameters (see Appendix D). Although our MHN results are a proof-of-concept, given the robustness of the MINERVA findings with respect to hyperparameters (see Figure 10) they suggest that elevated failure rates on collocations may be a property of associative memory retrieval writ large.

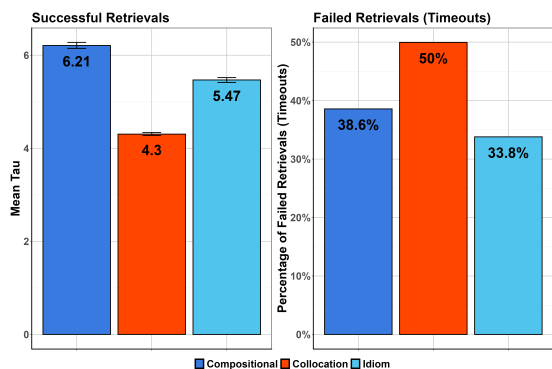


Figure 3: **Left:** mean Tau (τ) by condition for successful retrievals in MINERVA. The y-axis represents mean Tau, the model’s output which acts as a proxy for reaction times. Error bars indicate bootstrapped confidence intervals. **Right:** percentage of failed retrievals, i.e., timeouts, per condition. Note that while the pattern of Taus on *successful retrievals* is different from the pattern of human RTs, the pattern of *timeouts* per condition matches the pattern of human RTs.

Unlike the pattern of Taus on successful retrievals, the pattern of retrieval failures in MINERVA appears to capture the trend in human RTs across the three conditions.

Additionally, we found that MINERVA always succeeds at retrieving items above a high frequency threshold (Figure 4, black line). We find a similar frequency boundary in humans (Figure 4, green line), which lies very close to the MINERVA threshold. On items above this threshold⁸, participants did not show a significant difference in RT by condition, while still showing a significant effect of frequency.

5.4 Ablations

Semantics-only In the semantics-only ablation, the model was loaded with all instances being equally frequent. Thus, the only distinguishing factor between the memory traces were their semantics. The results are shown in Figure 5. We visually observe that the trends for collocations match those of the main experiment—they time out the most, but are retrieved the fastest on successful retrievals. However, unlike in the main experiment, we observe that compositional items are retrieved slightly faster and time out less frequently than idioms. Investigating the cause of this discrepancy is an interesting avenue for future work.

⁸16 compositional, 18 collocations, 17 idioms

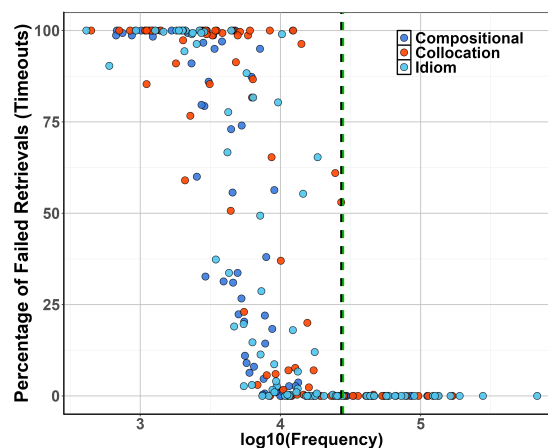


Figure 4: Percentage of failed retrievals (i.e., timeouts) in MINERVA per stimulus item, as a function of the frequency of the item. The x-axis is displayed in log scale. The black line indicates the frequency threshold ($f = 27123$) above which MINERVA times out less than 1% of the time. The green line ($f = 28000$) indicates the frequency threshold above which condition stops being a significant predictor of human RTs.

Frequency-only In the frequency-only ablation, the model was loaded with embeddings comprised of Gaussian noise⁹. However, each noise-item was sampled according to correct frequency information. The results are shown in Figure 6. For successful retrievals, we visually observed that idioms and collocations were retrieved equally quickly, whereas compositional items were retrieved slower. This pattern also persists in the timeouts. Given that frequency drives MINERVA’s retrieval mechanism, this pattern of Taus and timeouts is not surprising. Idiomatics—the most frequent subset—are retrieved most easily, followed by collocations, and finally compositional units which are the least frequent.

The results of these ablations suggest that it is the semantics of the item traces that drive the unique processing cost for collocations in MINERVA. Additionally, as shown with the trends for idioms vs compositional items, frequency dynamics mitigate the effects of semantic dynamics, producing the overall behaviour observed in the main experiment.

6 Discussion & Analysis

Our behavioural results confirm the processing cost for collocations surmised from the literature. This effect occurs despite collocations and compositional items being very close in frequency (with the balance in favour of collocations), and the par-

⁹We calibrate the noise to the mean and standard deviation of the embeddings in the main condition.

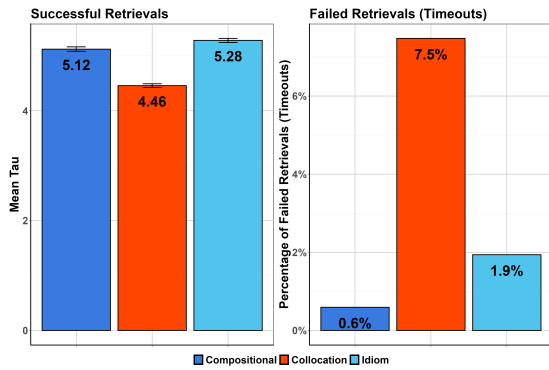


Figure 5: **Left:** mean Tau (τ) by condition for successful retrievals in Ablation 1, wherein frequency information was eliminated. The y-axis represents mean Tau, the model’s output which acts as a proxy for reaction times. Error bars indicate bootstrapped confidence intervals. **Right:** percentage of failed retrievals, i.e., timeouts, per condition in Ablation 1. Note that just as in the main MINERVA experiment, collocations time out much more frequently than the other conditions.

Participants as adult L1 English speakers being highly familiar with the items. The result is mirrored in our computational findings. These stark differences in processing patterns for collocations compared to idioms and compositional items suggest that they must be treated as a separate class of linguistic items, and not be cast under the broad umbrella of formulaic language.

To recap our simulation results, MINERVA’s successful retrievals failed to reproduce the processing trend observed in humans and also exhibited many more incorrect responses, i.e., unsuccessful memory retrievals, than humans. However, these retrieval failures do appear to capture the key asymmetries in human processing. Again, this is especially noticeable for collocations on which MINERVA timed out the most. We also found that above a certain frequency threshold, MINERVA matches human patterns. This suggests that simple memory retrieval, as implemented in a frequency based model of memory, is only sufficient to explain human processing trends for highly frequent items. Below this threshold, retrieval starts to fail. Given that MINERVA does not have any processing mechanism beyond memory retrieval, it simply times out on these items. We conjecture that at this point humans invoke other processing mechanisms to facilitate interpreting of the stimulus, incurring a cost in reaction time.

The fact that collocations incur a higher processing cost despite being more frequent than com-

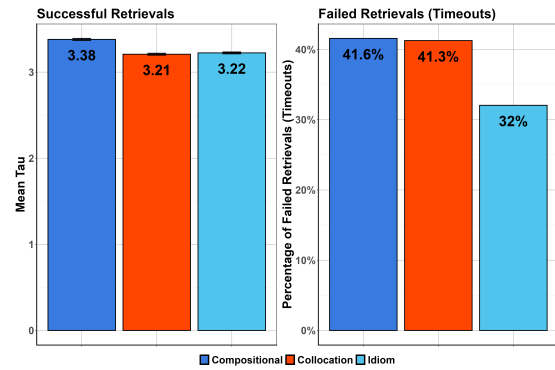


Figure 6: **Left:** mean Tau (τ) by condition for successful retrievals in Ablation 2, wherein semantic information was eliminated while leaving the correct item frequency distribution. The y-axis represents mean Tau, the model’s output which acts as a proxy for reaction times. Error bars indicate bootstrapped confidence intervals. **Right:** percentage of failed retrievals, i.e., timeouts, per condition. Note that the trends in timeouts follow the frequency distribution across the conditions.

positional items shows that single-route accounts provide an incomplete picture. They further demonstrate that dual-route accounts with a binary distinction between formulaic versus compositional language are also insufficient to account for the processing of this large and frequent subset of language. This underscores the need for a model which can account for a more fine-grained representation of semantic compositionality. One such plausible mechanism is analogical reasoning (Eddington, 2000; Ambridge, 2020). Like single-route models, this domain-general approach posits that all linguistic units are processed by a single mechanism (Skousen, 1990). However, in addition to memory retrieval, it posits on-the-fly analogy without resorting to any rule-based mechanisms. On receiving an input, a memory search is undertaken to find analogous exemplars previously experienced. The input is then evaluated based on the degree of similarity in order to find the most frequent category within the found set of most similar exemplars (Gentner and Namy, 2006).

Memory retrieval is the first step in analogical processing (Gentner and Colhoun, 2010). Thus, processing a sufficiently frequent item via analogy will simply resort to memory retrieval. Such a mechanism would be invariant to the semantic compositionality of the item in question, as we have seen in humans. Below this threshold, however, proper analogical machinery comes into play.

In compositional items, both the verb and the

noun play a prototypical role. Thus, even though the language user may not recall this exact verb-noun pairing from memory, it is relatively easy to map the verb and noun to similar instances of the same, due to the high semantic overlap between compositional uses of the verb and the noun. In collocations, however, the verb is not used in its prototypical sense. Resolving the meaning of the verb requires a much “farther” mapping, which may involve increased search over possible abstractions of the verb or extensive structure-mapping. Engaging such machinery inevitably incurs a processing cost with respect to compositional items (Gentner and Namy, 2006), as reflected in RTs. Finally, idioms, which cannot be processed analytically, must be retrieved holistically irrespective of frequency.

Moreover, there is a body of evidence for the role of analogy in metaphor comprehension (see Morsanyi et al., 2022, for a review), child language acquisition (see Raynal et al., 2024), and processing of novel verb metaphors (King and Gentner, 2022), which are, in essence, unconventionalized collocations. We posit that that an analogical account of language processing may provide a more complete explanation of these findings, and that further work should explore this proposal.

The retrieval failures for idioms may stem from a limitation of our dataset—the fact that we only consider figurative idioms which have a compositional reading. We were unable to ascertain the relative frequency of idiomatic versus literal readings in the context sentences of every idiom in our stimuli that we use to generate embeddings. It is also unknown to what precise extent sBERT can accurately represent idiomatic meanings, or whether our human participants interpreted idiomatic stimuli in a figurative sense. Combined, these factors suggest that the semantics of our set of idioms are somewhat akin to our set of compositional items, and some of the processing trends which pertain to compositional items are inadvertently present in the trend of responses to idioms. In line with the holistic retrieval hypothesis, we surmise that idioms for which the literal reading is much less frequent than the idiomatic one (e.g., *kick the bucket*) will tend to be processed faster and with fewer timeouts than more ambiguous ones (e.g., *hold the key*). Future work will attempt to investigate this prediction and further augment our understanding of idiomatic processing by including pure idioms, i.e., those without a literal reading, in the dataset, and employing other behavioural tasks which in-

volve presentation of items within context (e.g., self-paced reading).

One intriguing implication of our computational experiment may be of interest to the NLP community. Switching Equation 2 for $a_\tau = \text{softmax}(\tau s)$, MINERVA’s retrieval mechanism becomes identical to query-key-value attention in Transformers (Vaswani et al., 2017). Here, the probe plays the role of the query and the memory items the role of keys and value. Increasing Tau can be considered as a mechanism for dynamically weighting the output toward the keys which bear most similarity to the query, similar to increasing the softmax scale parameter over time. Under this formulation, MINERVA is also a variant of the Modern Hopfield Network with no learnable parameters (Ramsauer et al., 2021). This wealth of connections suggests that our findings may apply more broadly to all attention-based methods. Given the prevalence of collocations in language, if neural embeddings of semi-compositional language are particularly prone to failures in attention-based retrieval, this may significantly impair language understanding and generation in Transformer-based models. Future work will attempt to mechanistically diagnose the underlying reasons for the increased failure rates in collocations and ascertain whether these issues impact the performance of general-purpose Transformer-based language models.

On a more applied level, large language models (LLMs) as writing assistants have gained popularity (Boisson et al., 2024). This has drawn attention to how these models handle figurative language like metaphors, of which collocations are a prominent subset. While LLMs can produce metaphors, users often note shortcomings such as clichéd phrasing or a lack of creativity in metaphor generation (Chakrabarty et al., 2024), perhaps reflecting an underlying lack of capability in interpreting these linguistic units. These observations further underscore the importance of understanding collocational processing not only in human cognition, but also in NLP systems.

Overall, we show that semi-compositional units are a bigger “pain in the neck” (Sag et al., 2002) than other subsets of the semantic compositionality continuum: too complex for rote retrieval, yet too idiosyncratic for rule-based computation. As it stands, memory retrieval does *leave something on the table*, underscoring the need for theories that capture the graded nature of meaning and structure in language.

7 Limitations

- Our approach relies on contextual embeddings to capture semantic information. However, these embeddings do not always differentiate clearly between compositional and idiomatic readings. Given that our idiomatic stimuli also have a productive reading, the same embedding may be used for both literal and figurative interpretations. Similarly, we cannot ensure that our task is eliciting an idiomatic reading in humans as human listeners disambiguate based on context.
- The current dataset was not built from scratch with frequency-matching criteria for idioms. Frequency is a well-established predictor of language processing and an ideal dataset would equate or carefully control the frequency distributions of idioms relative to other word types.
- Our study exclusively examined verb–noun (VN) collocations. While these are a critical class of multiword expressions, little is known about other collocational structures (e.g., adjective–noun, phrasal verbs, etc.) which are also prevalent in natural language and may be processed differently. Extending our investigation to these additional types will be important for assessing the generalizability of our findings across the broader spectrum of semi-compositional linguistic units.
- MINERVA provides a parsimonious framework for modelling memory retrieval, yet it inherently simplifies many aspects of human cognitive processing. The model does not integrate attentional mechanisms or dynamic contextual cues beyond the static embeddings provided, and it does not account for developmental changes in memory and language processing. These simplifications may limit the model’s ability to capture the full complexity of human language processing, particularly in cases where retrieval failures (time-outs) interact with other cognitive processes. Our simulations relied on specific hyperparameter settings (e.g., activation threshold $K=0.99$ and forgetting probability $F=0.8$) that were chosen based on qualitative assessments. Although results were robust across a range of parameter values, the possibility remains that

different parametrizations could yield different patterns.

Ethics

The study received ethics approval in accordance with the University of Edinburgh’s School of Philosophy, Psychology and Language Sciences Research Ethics Process (RT number: 339-2122/4).

Data and Code

Our data and code will be made available at the following URL: <https://github.com/psydelle/minerva-release>

Acknowledgements

This work was primarily supported by the UKRI Centre for Doctoral Training in Natural Language Processing, funded by the UKRI (grant EP/S022481/1) and the University of Edinburgh, School of Informatics and School of Philosophy, Psychology & Language Sciences.

SdS is also grateful for the support of the Fundação para a Ciência e a Tecnologia, under project UID/00214: Centro de Linguística da Universidade de Lisboa, of Center of Linguistics of the University of Lisbon, from School of Arts and Humanities of the University of Lisbon.

We are also grateful to Sivan Milton for checking our stimuli set, and Mattia Opper and Dr Edoardo Ponti for advice on adapting sBERT embeddings for our experiment.

References

- Ben Ambridge. 2020. *Against stored abstractions: A radical exemplar model of language acquisition*. 40(5-6):509–559.
- Ben Ambridge and E. Lieven. 2011. *Child language acquisition: Contrasting theoretical approaches*. Cambridge University Press, Cambridge.
- Jason Arndt and Elliot Hirshman. 1998. *True and False Recognition in MINERVA2: Explanations from a Global Matching Perspective*. *Journal of Memory and Language*, 39(3):371–391.
- Andy Barfield and Henrik Gyllstad, editors. 2009. *Researching Collocations in Another Language*. Palgrave Macmillan UK, London.
- Dale J. Barr, Roger Levy, Christoph Scheepers, and Harry J. Tily. 2013. *Random effects structure for confirmatory hypothesis testing: Keep it maximal*. *Journal of Memory and Language*, 68(3):255–278. Publisher: Elsevier Inc.

- Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. 2015. [Fitting Linear Mixed-Effects Models Using lme4](#). *Journal of Statistical Software*, 67(1):1–48.
- Joanne Boisson, Asahi Ushio, Hsuvas Borkakoty, Kiamehr Rezaee, Dimosthenis Antypas, Zara Siddique, Nina White, and Jose Camacho-Collados. 2024. [How Are Metaphors Processed by Language Models? The Case of Analogies](#). In *Proceedings of the 28th Conference on Computational Natural Language Learning*, pages 365–387, Miami, FL, USA. Association for Computational Linguistics.
- Joan Bybee. 2006. From usage to grammar: The mind’s response to repetition. *Language*, 82(4).
- Joan Bybee. 2012. [A usage-based perspective on language](#). In *Language, Usage and Cognition*, pages 1–13. Cambridge University Press, Cambridge.
- Gareth Carrol and Kathy Conklin. 2014. [Getting your wires crossed: Evidence for fast processing of L1 idioms in an L2](#). *Bilingualism: Language and Cognition*, 17(4):784–797. Publisher: Cambridge University Press.
- Gareth Carrol and Kathy Conklin. 2020. [Is All Formulaic Language Created Equal? Unpacking the Processing Advantage for Different Types of Formulaic Sequences](#). *Language and Speech*, 63(1):95–122.
- Tuhin Chakrabarty, Vishakh Padmakumar, Faeze Brahman, and Smaranda Muresan. 2024. [Creativity Support in the Age of Large Language Models: An Empirical Study Involving Emerging Writers](#). *arXiv preprint*. ArXiv:2309.12570 [cs].
- Morten H. Christiansen and Nick Chater. 2008. [Language as shaped by the brain](#). *Behavioral and Brain Sciences*, 31(5):489–509.
- Chrissy M. Chubala and Randall K. Jamieson. 2013. [Recoding and representation in artificial grammar learning](#). *Behavior Research Methods*, 45(2):470–479.
- Michael C. Corballis. 2019. [Language, Memory, and Mental Time Travel: An Evolutionary Perspective](#). *Frontiers in Human Neuroscience*, 13. Publisher: Frontiers.
- Anthony Paul Cowie. 1998. *Phraseology : theory, analysis, and applications*. Clarendon Press, Oxford. Series Title: Oxford studies in lexicography and lexicology.
- Peter W. Culicover, Ray Jackendoff, and Jenny Audring. 2017. [Multiword Constructions in the Grammar](#). *Topics in Cognitive Science*, 9(3):552–568.
- Tamara Czinczoll, Helen Yannakoudakis, Pushkar Mishra, and Ekaterina Shutova. 2022. [Scientific and Creative Analogies in Pretrained Language Models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2094–2100, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Ishita Dasgupta and Samuel J. Gershman. 2021. [Memory as a Computational Resource](#). *Trends in Cognitive Sciences*, 25(3):240–251.
- Sydelle de Souza, Francis Mollica, and Jennifer Culbertson. 2024. [What can L1 speakers tell us about killing hope? A Novel Behavioral Measure for Identifying Collocations](#). *Proceedings of the Annual Meeting of the Cognitive Science Society*, 46(0).
- Mete Demircigil, Judith Heusel, Matthias Löwe, Sven Uppgang, and Franck Vermet. 2017. [On a Model of Associative Memory with Huge Storage Capacity](#). *Journal of Statistical Physics*, 168(2):288–299.
- Dagmar Divjak. 2019. *Frequency in Language: Memory, Attention and Learning*. Cambridge University Press, Cambridge.
- Dagmar Divjak, Petar Milin, Srdan Medimorec, and Maciej Borowski. 2022. [Behavioral Signatures of Memory Resources for Language: Looking beyond the Lexicon/Grammar Divide](#). *Cognitive Science*, 46(11):e13206.
- Philip Durrant and Norbert Schmitt. 2009. [To what extent do native and non-native writers make use of collocations?](#) *IRAL - International Review of Applied Linguistics in Language Teaching*, 47(2):157–177.
- David Eddington. 2000. [Analogy and the dual-route model of morphology](#). *Lingua*, 110(4):281–298.
- Stefan Falkner, Aaron Klein, and Frank Hutter. 2018. [BOHB: robust and efficient hyperparameter optimization at scale](#). In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 1436–1445. PMLR.
- Irene Fioravanti, Marco Silvio Giuseppe Senaldi, Alessandro Lenci, and Anna Siyanova-Chanturia. 2021. [Lexical fixedness and compositionality in L1 speakers’ and L2 learners’ intuitions about word combinations: Evidence from Italian](#). *Second Language Research*, 37(2):291–322.
- Gottlob Frege. 1892. *Über Sinn und Bedeutung*, 1. auflage edition. Zeitschrift für Philosophie und philosophische Kritik, Neue Folge. Pfeffer, Leipzig.
- Lukas Galke, Yoav Ram, and Limor Raviv. 2024. [Deep neural networks and humans both benefit from compositional language structure](#). *Nature Communications*, 15(1):10816. Publisher: Nature Publishing Group.
- Dedre Gentner and Julie Colhoun. 2010. [Analogical Processes in Human Thinking and Learning](#). In Britt Glatzeder, Vinod Goel, and Albrecht Müller, editors, *Towards a Theory of Thinking: Building Blocks for a Conceptual Framework*, pages 35–48. Springer, Berlin, Heidelberg.

- Dedre Gentner and Laura L. Namy. 2006. [Analogical Processes in Language Learning](#). *Current Directions in Psychological Science*, 15(6):297–301. Publisher: SAGE Publications Inc.
- Lynn E. Grant. 2005. Frequency of ‘core idioms’ in the British National Corpus (BNC). *International Journal of Corpus Linguistics*, 10(4):429–451. Publisher: John Benjamins.
- Henrik Gyllstad and Brent Wolter. 2016. [Collocational Processing in Light of the Phraseological Continuum Model: Does Semantic Transparency Matter?](#) *Language Learning*, 66(2):296–323.
- Susanne Handl and Eva-Maria Graf. 2010. [Collocation, anchoring, and the mental lexicon – an ontogenetic perspective](#). In Hans-Jörg Schmid and Susanne Handl, editors, *Cognitive Foundations of Linguistic Usage Patterns: Empirical Studies*, pages 119–148. De Gruyter Mouton.
- Douglas L. Hintzman. 1984. [MINERVA 2: A simulation model of human memory](#). *Behavior Research Methods, Instruments, & Computers*, 16(2):96–101.
- Douglas L. Hintzman. 1988. [Judgments of frequency and recognition memory in a multiple-trace memory model](#). *Psychological Review*, 95(4):528.
- J J Hopfield. 1982. [Neural networks and physical systems with emergent collective computational abilities](#). *Proceedings of the National Academy of Sciences of the United States of America*, 79(8):2554–2558.
- Peter Howarth. 1998. [Phraseology and second language proficiency](#). *Applied Linguistics*, 19(1):24–44.
- 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.
- Randall K. Jamieson and D. J.K. Mewhort. 2009. [Applying an exemplar model to the artificial-grammar task: Inferring grammaticality from similarity](#). *Quarterly Journal of Experimental Psychology*, 62(3):550–575.
- Adam Kilgarriff, Pavel Rychly, Pavel Smrz, and David Tugwell. 2024. The Sketch Engine.
- Daniel King and Dedre Gentner. 2022. [Verb Metaphoric Extension Under Semantic Strain](#). *Cognitive Science*, 46(5):e13141. Publisher: John Wiley & Sons, Ltd.
- Nate Kornell and Robert A. Bjork. 2009. [A stability bias in human memory: Overestimating remembering and underestimating learning](#). *Journal of Experimental Psychology: General*, 138(4):449–468.
- Dmitry Krotov and John J. Hopfield. 2016. Dense associative memory for pattern recognition. In *Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS’16*, pages 1180–1188, Red Hook, NY, USA. Curran Associates Inc.
- George Lakoff and Mark Johnson. 1980. *Metaphors We Live By*. University of Chicago Press, Chicago.
- Thomas K. Landauer and Susan T. Dumais. 1997. [A solution to Plato’s problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge](#). *Psychological Review*, 104(2):211–240.
- Ronald W Langacker. 1987. *Foundations of cognitive grammar*. Stanford University Press, Stanford, Calif.
- Emmy Liu, Chenxuan Cui, Kenneth Zheng, and Graham Neubig. 2022. [Testing the Ability of Language Models to Interpret Figurative Language](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4437–4452, Seattle, United States. Association for Computational Linguistics.
- Gordon D. Logan. 1997. [Automaticity and Reading: Perspectives from the Instance Theory of Automatization](#). *Reading & Writing Quarterly*, 13(2):123–146.
- Gordon D Logan and Joseph L Etherton. 1994. [What Is Learned During Automatization? The Role of Attention in Constructing an Instance](#). *Journal of Experimental Psychology: Learning Memory, and Cognition*, 20(5):1022–1050.
- James L McClelland and David E Rumelhart. 1985. On learning the past tense of English verbs. In J. L. McClelland and D. E. and the PDP Research Group Rumelhart, editors, *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, volume 2. Bradford Books/MIT Press, Cambridge, MA.
- Igor Mel’čuk. 2003. Collocations: définition, rôle et utilité. *Travaux et recherches en linguistique appliquée. Série E, Lexicologie et lexicographie.*, (1):23–31. Num Pages: 9 Place: Amsterdam Publisher: Editions ‘De Werelt’.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. [Distributed Representations of Words and Phrases and their Compositionality](#). In *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc.
- Kinga Morsanyi, Jayne Hamilton, Dušan Stamenković, and Keith J. Holyoak. 2022. [Linking metaphor comprehension with analogical reasoning: Evidence from typical development and autism spectrum disorder](#). *British Journal of Psychology*, 113(2):479–495.
- Rachel A. G. Mueller and Raymond W. Gibbs. 1987. [Processing idioms with multiple meanings](#). *Journal of Psycholinguistic Research*, 16(1):63–81.
- J. Nick Reid and Randall K. Jamieson. 2023. [True and false recognition in MINERVA 2: Extension to sentences and metaphors](#). *Journal of Memory and Language*, 129:104397. Publisher: Elsevier Inc.

- Tomomi Nishikawa. 2019. [Non-nativelike outcome of naturalistic child L2 acquisition of Japanese: The case of noun–verb collocations](#). *International Review of Applied Linguistics in Language Teaching*, (Lenneberg 1967).
- Ira A. Noveck, Nicholas Griffen, and Diana Mazzarella. 2023. Taking stock of an idiom’s background assumptions: an alternative relevance theoretic account. *Frontiers in Psychology*, 14.
- Timothy J O’Donnell, Noah D Goodman, and Joshua B Tenenbaum. 2009. Fragment Grammars : Exploring Computation and Reuse in Language Fragment Grammars. *Computer Science and Artificial Intelligence Laboratory Technical Report*.
- Andrew Pawley and Frances Hodgets Syder. 1983. [Two puzzles for linguistic theory: Nativelike selection and nativelike fluency](#). *Language and Communication*, pages 191–226. ISBN: 9781317869634.
- Steven Pinker. 1991. Rules of Language. In *Science*, volume 253, pages 530–535. American Association for the Advancement of Science. Issue: 5019.
- Steven Pinker. 2013. [Learnability and Cognition, New Edition: The Acquisition of Argument Structure](#). MIT Press, Cambridge, MA.
- Steven Pinker and Michael Ullman. 2002. [The past and future of the past tense](#). *Trends in Cognitive Sciences*, 6(11):456–463.
- R Core Team. 2024. [R: a language and environment for statistical computing](#). manual, R Foundation for Statistical Computing, Vienna, Austria.
- Hubert Ramsauer, Bernhard Schöfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Lukas Gruber, Markus Holzleitner, Thomas Adler, David P. Kreil, Michael K. Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. 2021. [Hopfield Networks is All You Need](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Lucas Raynal, Evelyne Clément, Louise Goyet, Pia Rämä, and Emmanuel Sander. 2024. [Neural correlates of unconventional verb extensions reveal preschoolers’ analogical abilities](#). *Journal of Experimental Child Psychology*, 246:105984.
- Emily Reif, Ann Yuan, Martin Wattenberg, Fernanda B Viegas, Andy Coenen, Adam Pearce, and Been Kim. 2019. [Visualizing and Measuring the Geometry of BERT](#). In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Nils Reimers and Iryna Gurevych. 2019. [SentenceBERT: Sentence Embeddings using Siamese BERT-Networks](#). *arXiv preprint*. ArXiv:1908.10084 [cs].
- Nick Riches, Carolyn Letts, Hadeel Awad, Rachel Ramsey, and Ewa Dąbrowska. 2022. [Collocational knowledge in children: a comparison of English-speaking monolingual children, and children acquiring English as an Additional Language](#). *Journal of Child Language*, 49(5):1008–1023.
- Ivan A Sag, Timothy Baldwin, Francis Bond, Ann Copestake, and Dan Flickinger. 2002. [Multiword expressions: A pain in the neck for NLP](#). *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2276:1–15. ISBN: 3540432191.
- Mark S. Seidenberg and David C. Plaut. 2014. [Quasiregularity and Its Discontents: The Legacy of the Past Tense Debate](#). *Cognitive Science*, 38(6):1190–1228.
- Diana Sidtis. 2020. [Familiar Phrases in Language Competence](#). In Alexander Haselow and Gunther Kaltenböck, editors, *Grammar and Cognition : Dualistic Models of Language Structure and Language Processing*, pages 38–67. John Benjamins Publishing Company, Amsterdam, Philadelphia.
- Royal Skousen. 1990. [Analogical Modeling of Language](#). Springer Netherlands, Dordrecht.
- Suhad Sonbul, Dina Abdel Salam El-Dakhs, and Rezan Alharbi. 2024. [Rendering natural collocations in a translation task: The effect of direction, congruency, semantic transparency, and proficiency](#). *International Journal of Applied Linguistics*, 34(1):117–133.
- Patrizia Tabossi, Rachele Fanari, and Kinou Wolf. 2008. [Processing idiomatic expressions: Effects of semantic compositionality](#). *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 34(2):313–327.
- Xiaoyu Tong, Ekaterina Shutova, and Martha Lewis. 2021. [Recent advances in neural metaphor processing: A linguistic, cognitive and social perspective](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4673–4686, Online. Association for Computational Linguistics.
- Mei-Hsing Tsai. 2020. [Teaching L2 collocations through concept-based instruction: The effect of L2 proficiency and congruency](#). *International Journal of Applied Linguistics*, 30(3):553–575.
- 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*, 2017-Decem(Nips):5999–6009. ArXiv: 1706.03762.
- Laura Vilkaite and Norbert Schmitt. 2019. [Reading collocations in an L2: Do collocation processing benefits extend to non-adjacent collocations?](#) *Applied Linguistics*, 40(2):329–354.

Ivan Vulić, Edoardo Maria Ponti, Robert Litschko, Goran Glavaš, and Anna Korhonen. 2020. [Probing Pretrained Language Models for Lexical Semantics](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7222–7240, Online. Association for Computational Linguistics.

Dian Wang, Yang Li, Suge Wang, Xin Chen, Jian Liao, Deyu Li, and Xiaoli Li. 2025. [CKEMI: Concept knowledge enhanced metaphor identification framework](#). *Information Processing & Management*, 62(1):103946.

Brent Wolter and Henrik Gyllstad. 2013. [Frequency of Input and L2 Collocational Processing: A Comparison of Congruent and Incongruent Collocations](#). *Studies in Second Language Acquisition*, 35:451–482.

Alison Wray. 2002. *Formulaic Language and the Lexicon*, volume 80. Cambridge University Press, Cambridge. Publication Title: Language ISSN: 1535-0665.

Alison Wray. 2008. *Formulaic language: pushing the boundaries*. Oxford University Press, Oxford. Series Title: Oxford Applied Linguistics.

Satoshi Yamagata, Tatsuya Nakata, and James Rogers. 2023. [Effects of distributed practice on the acquisition of verb-noun collocations](#). *Studies in Second Language Acquisition*, 45(2):291–317.

Junko Yamashita. 2018. [Possibility of semantic involvement in the L1-L2 congruency effect in the processing of L2 collocations](#). *Journal of Second Language Studies*, 1(1):60–78.

Omnia Zayed, John Philip McCrae, and Paul Buitelaar. 2018. [Phrase-Level Metaphor Identification Using Distributed Representations of Word Meaning](#). In *Proceedings of the Workshop on Figurative Language Processing*, pages 81–90, New Orleans, Louisiana. Association for Computational Linguistics.

A Dataset Statistics

Table 1: Descriptive statistics of phrasal frequency by condition

Condition	Mean	SD	N
Compositional	19374.47	30671.53	78
Collocation	21528.21	30971.42	78
Idiom	36784.68	87468.40	78

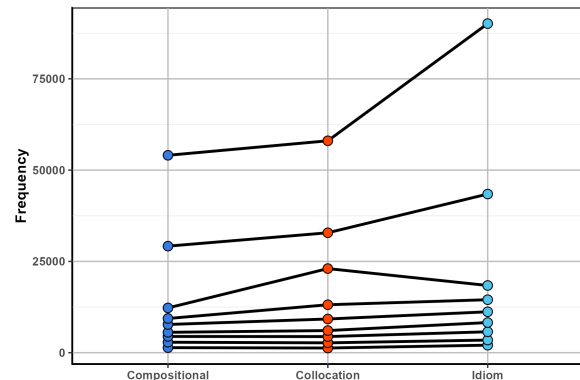


Figure 7: Item frequencies across conditions, by decile

B Human Data

Table 2: Descriptive statistics of human reaction times (ms) by condition

Condition	Mean	SD	N
Idiom	984.20	365.39	4462
Compositional	995.32	375.76	4423
Collocation	1007.87	370.84	4484

Table 3: Descriptive statistics of human accuracy by condition

Condition	Mean	SD	N
Idiom	0.93	0.25	4785
Compositional	0.92	0.27	4791
Collocation	0.94	0.24	4772

Table 4: Number of incorrect trials by condition

Condition	n
Compositional	400
Collocation	464
Idiom	433
Baseline	1455

C GLMM Results for Main Simulation

Table 5: Generalized mixed-effects regression results for human AJT reaction times (left), and Tau, a proxy for reaction times, simulated in MINERVA (right). MINERVA is run with $K = 0.99$, $F = 0.8$. Only correct responses and successful retrievals are analysed.

	<i>Dependent variable:</i>	
	RT	Tau
	<i>Human</i>	MINERVA
Compositional	4.690** (2.240)	0.624*** (0.004)
Collocation	13.800*** (1.760)	-0.410*** (0.004)
Frequency	-18.500*** (1.640)	-0.541*** (0.004)
Constant	1,047.0*** (2.140)	5.900*** (0.004)
N	13,369	43,708
<i>Note:</i>	* $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$	

D Modern Hopfield Network Experiments

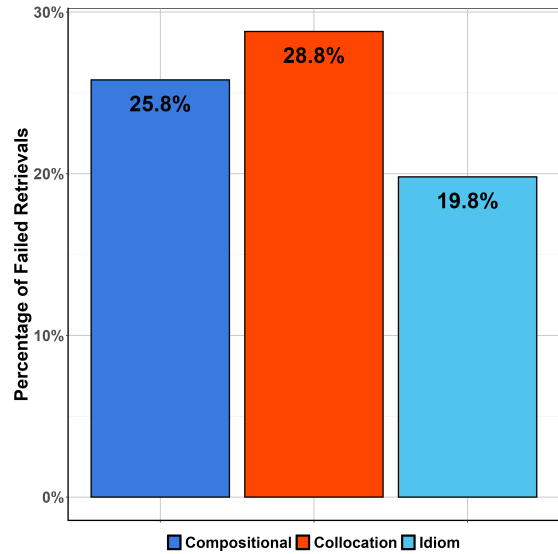


Figure 8: Percentage of failed retrievals in the Modern Hopfield Network. Collocations fail the most, followed by compositional items, and idioms fail the least. This matches the trends for the failures by condition in MINERVA, and the trend of reaction times in humans.

We suggest that our computational results pertain to the algorithmic implications of memory retrieval and are not based on a spurious quirk of the MINERVA architecture. To reinforce this claim, we present a proof-of-concept simulation of the same behaviour with the Modern Hopfield Network (Ramsauer et al., 2021, MHN).

The classical Hopfield network (Hopfield, 1982) is a model of associative memory based on binary states, designed for pattern completion and associative memory retrieval. MHN generalizes the Hopfield network and the recent iterations thereof (Krotov and Hopfield, 2016; Demircigil et al., 2017) with continuous states, and introduces a corresponding update rule which is closely connected to the query-key-value attention operation in Transformers.

The MHN can be configured in a number of different ways in order to simulate different functions, such as pattern matching, sequence pooling, and attention. In this simulation, we focus on the configuration denoted as `HopfieldLayer` in Ramsauer et al. (2021). Specifically, the model learns a static memory matrix of size $M \in \mathbb{R}^{n \times c}$, where n is the number of memory items and c is the dimensionality of the memory’s latent space. It also learns projection matrices $A \in \mathbb{R}^{d \times c}$, $B \in \mathbb{R}^{c \times d}$ which

Table 6: Logistic regression results for retrieval failures. Reference Level is Idiom. **Left column:** MINERVA is run with $K = 0.99$, $F = 0.8$. **Right column:** Modern Hopfield Network for parameters presented in Appendix D. Both models converged but with singular fits. This was due to the (1 | ID) random intercept (where ID is the random seed for the model run) accounting for nearly 0 variance.

<i>Dependent variable:</i>		
Retrieval Failures		
	MINERVA	MHN
Compositional	0.68*** (0.63, 0.73)	0.817*** (0.71, 0.942)
Collocation	4.074*** (3.78, 4.39)	1.192** (1.022, 1.39)
Frequency	1.61e−08*** (8.68e−09, 2.98e−08)	1.75e−25*** (1.22e−26, 2.53e−24)
Constant	1,047.0*** (1.02e−03, 3.58e−03)	5.900*** (1.03e−10, 8.25e−10)
N	70200	23400
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01	

project inputs into and out of the latent space, respectively. Matching between inputs and memory items is done via query-key attention in the latent space, so a lower value for c forces the MHN to compress more strongly.

The experimental setup is similar to that for MINERVA, with the difference being the learnable nature of the model’s memory. As in Section 5.1, for each stimulus item $x \in \mathbb{R}^d$ we construct a number of noisy versions $\{x'_1, \dots, x'_i\}$ of its vector embedding, where i is proportional to the item’s corpus frequency. Analogously to the MINERVA memory matrix, our total training set for the MHN comprises 10,000 embeddings, with more frequent items being more represented. Given one such noisy embedding x' , the MHN was trained to output the un-noisy source embedding x . Specifically, it must maximize $s = \text{cosine_similarity}(\hat{x}, x)$, where \hat{x} model output. We conjecture that to succeed on the task, the model must optimize its limited memory to map multiple noisy versions of the same item to the item’s canonical representation.

As in the MINERVA experiments, we judge whether the model’s retrieval was successful based on the cosine similarity s . If s is above threshold K , we say the retrieval is a success. Otherwise, it

is a failure. Unlike in MINERVA, there is currently no analogue to RTs in MHN.

We ran a Bayesian hyperparameter sweep with Hyperband early stopping (Falkner et al., 2018) to find configurations in which the pattern of failures across the three conditions which match human trends. The results of one such configuration run over 100 seeds is presented in Figure 8. As in the MINERVA experiments, collocations fail the most, followed by compositional items, with idioms failing the least. All differences between conditions are statistically significant (see Table 6) when analyzed with the same model as the MINERVA failures. This matches the trends for the failures by condition in MINERVA, and the trend of reaction times in humans.

Please note that our parameter sweep was not exhaustive and we cannot make any claims with respect to hyperparameter robustness of the MHN simulations. However, the fact that the MHN can be shown to display the same behaviour as MINERVA in at least one case suggests that the elevated retrieval failure rate of collocations may be a property of associative memory retrieval writ large, rather than a quirk of the MINERVA architecture.

The model hyperparameters are as follows: $F =$

0.2, $K = 0.971$, $c = 625$, $n = 2205$. Scaling parameter β was set to $\frac{1}{\sqrt{c}} = 0.0016$. The model was trained for 300 epochs with batch size 16 on one A100 GPU.

F Hyperparameter Sweeps for Simulation Experiments

See next page.

E Averaging vs Concatenating sBERT Embeddings

In this ablation, we investigate the impact which concatenating verb and noun embeddings has on our modelling results. Instead of concatenating verb and noun embeddings, we perform mean-pooling across them, the same as we do for sub-word tokens. As shown in Figure 9, the trends exhibited by the model in the $K = 0.99$, $F = 0.8$ hyperparameter configuration are largely the same as those reported in the main text.

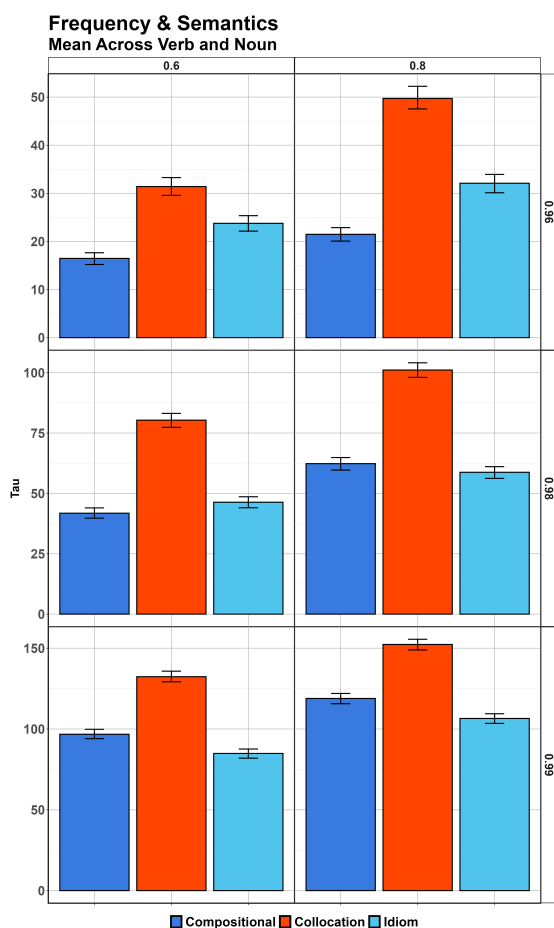


Figure 9: Reduced hyperparameter sweep showing the effects of mean-pooling the verb and noun embeddings before loading them into MINERVA, instead of concatenating them. Note that the hyperparameter combination reported in the main text is $K = 0.99$, $F = 0.8$.

Frequency & Semantics With Contextual Embeddings

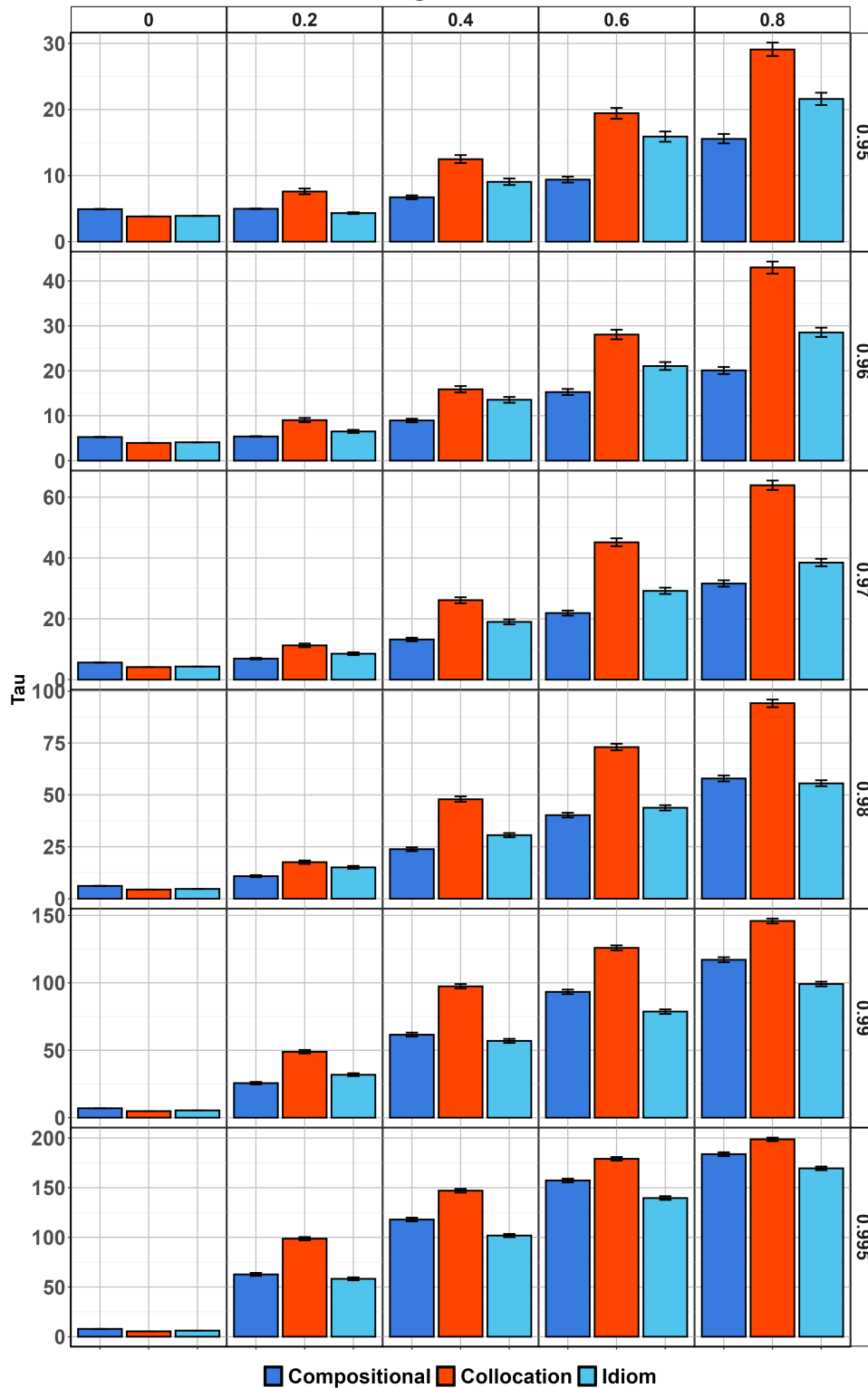


Figure 10: Results of the hyperparameter sweep for all values of activation threshold K and forgetting probability F for our main experiment. Error bars indicate bootstrapped confidence intervals. Note the difference in scales on the y-axis.

Semantics-only With Contextual Embeddings

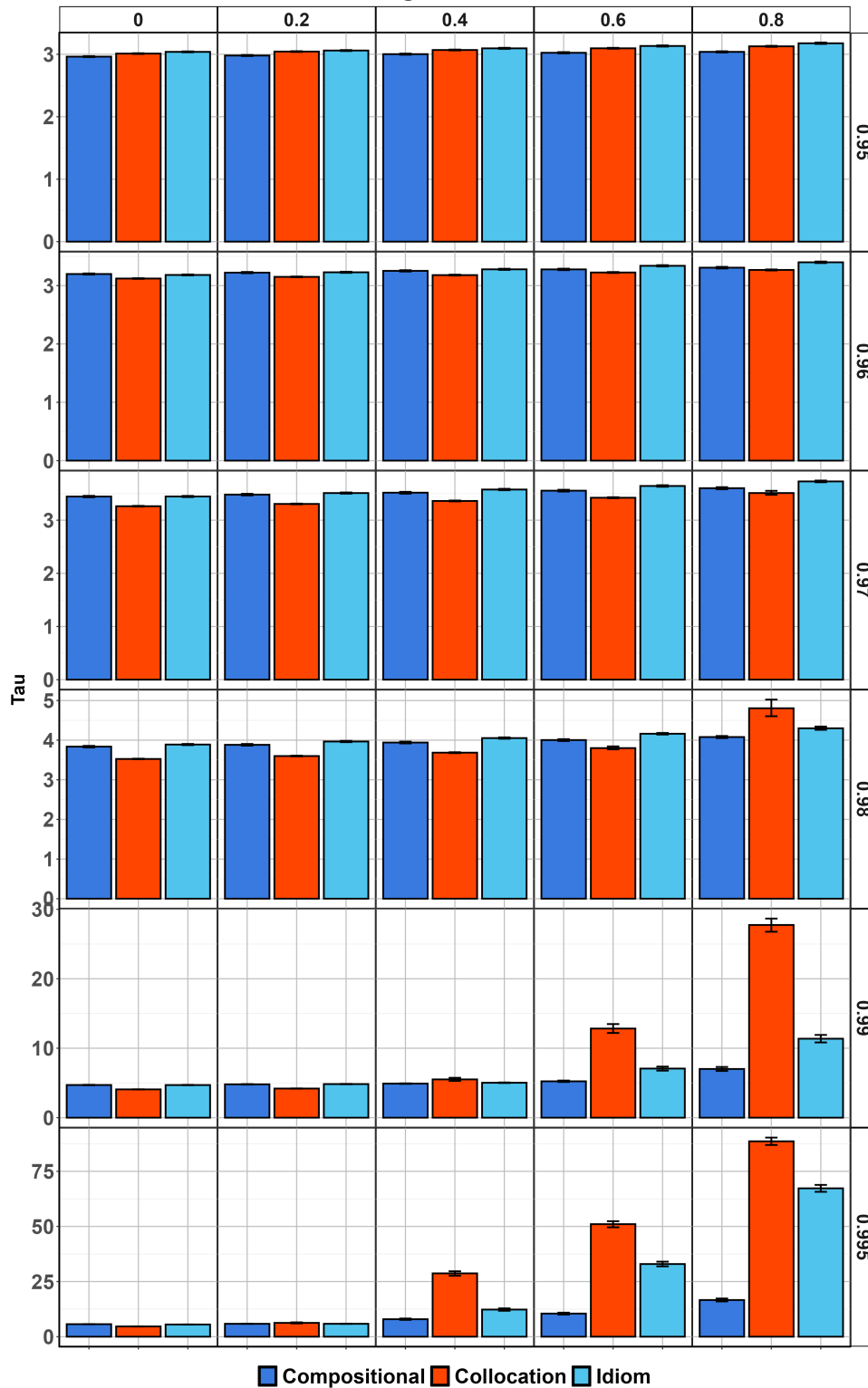


Figure 11: Results of the hyperparameter sweep for all values of activation threshold K and forgetting probability F for Simulation 2: Semantics-only wherein the matrix was loaded with all items having equal frequency. Error bars indicate bootstrapped confidence intervals. Note the difference in scales on the y-axis.

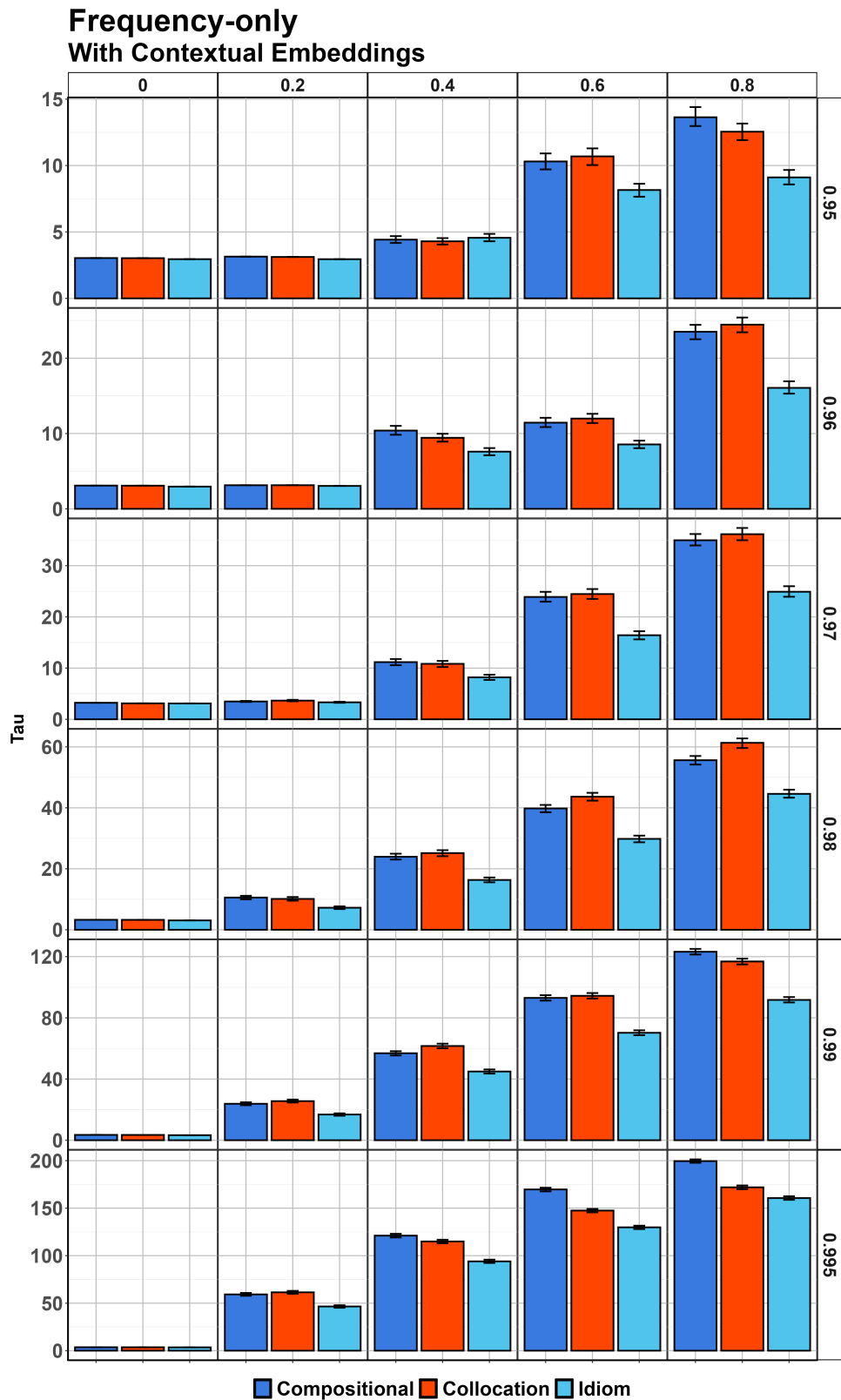


Figure 12: Results of the hyperparameter sweep for all values of activation threshold K and forgetting probability F for Simulation 2: Semantics-only wherein the matrix was loaded with noised embeddings but with the correct frequency. Error bars indicate bootstrapped confidence intervals. Note the difference in scales on the y-axis.

Null Model With Contextual Embeddings

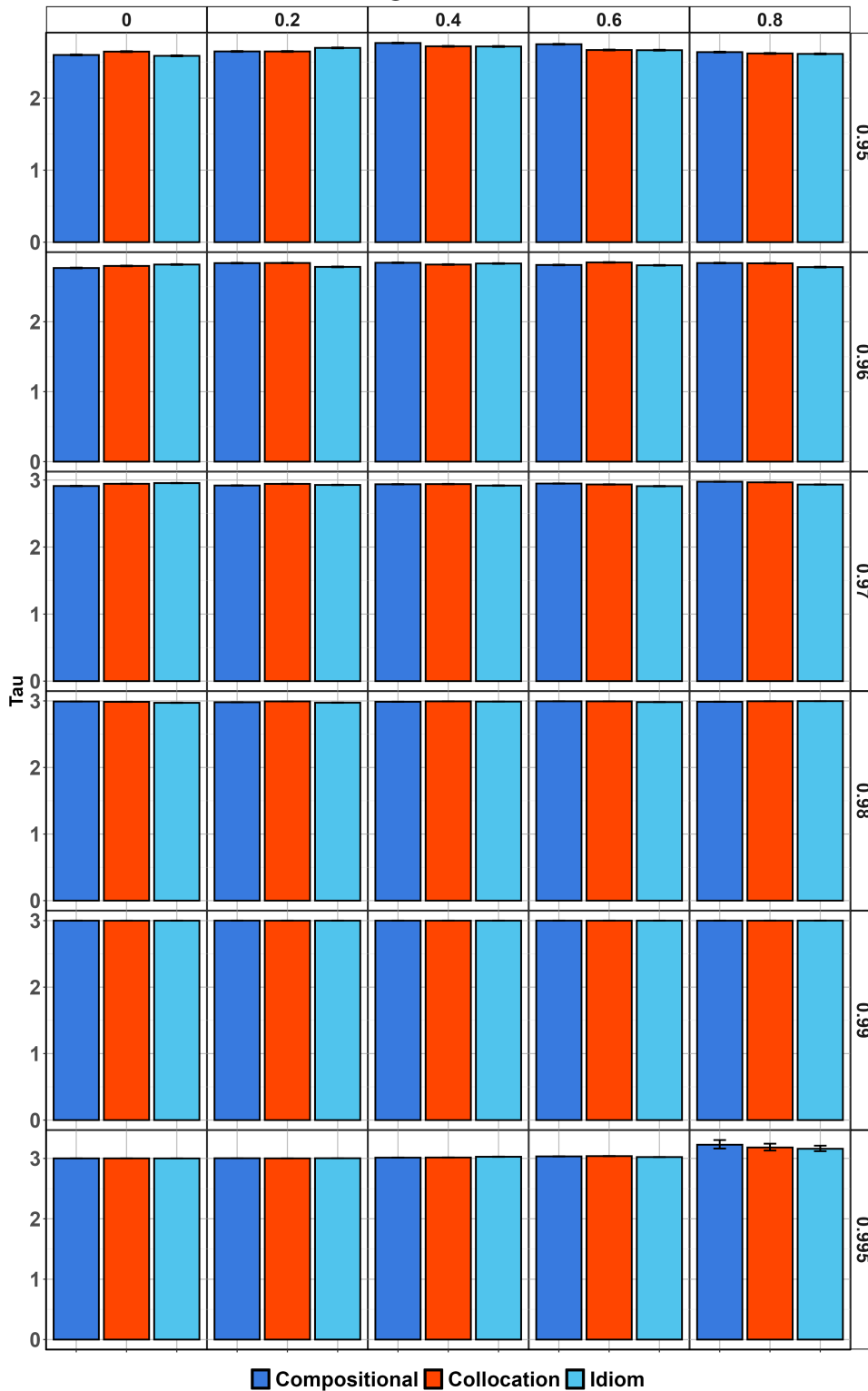


Figure 13: Results of the hyperparameter sweep for all values of activation threshold K and forgetting probability F for the Null Model wherein all the items in the matrix were loaded with noised embeddings and equal frequency. Error bars indicate bootstrapped confidence intervals. Note the difference in scales on the y-axis.