

A Continuous Approach to Metaphorically Motivated Regular Polysemy in Language Models

Anonymous ACL submission

Abstract

Linguistic accounts show that a word’s polysemy structure is largely governed by systematic sense alternations that form overarching patterns across the vocabulary. While psycholinguistic studies confirm the psychological validity of regularity in human language processing, in the research on large language models (LLMs) this phenomenon remains largely unaddressed. Revealing models’ sensitivity to systematic sense alternations of polysemous words can give us a better understanding of how LLMs process ambiguity and to what extent they emulate representations in the human mind. For this, we employ the measures of surprisal and semantic similarity as proxies of human judgment on the acceptability of novel senses. We focus on two aspects that have not received much attention previously – metaphorically motivated patterns and the continuous nature of regularity. We find evidence that surprisal from language models represents regularity of polysemic extensions in a human-like way, discriminating between different types of senses and varying regularity degrees, and overall strongly correlating with human acceptability scores.

1 Introduction

Polysemy, a linguistic phenomenon whereby a word is associated with multiple related senses, is fundamental to language. As most lexical words are polysemes to varying degrees (Zipf, 1945; Durkin and Manning, 1989; Haber and Poesio, 2024), this form of ambiguity remains a challenge for NLP. However, recent studies show that current language models (LMs) based on Transformers are able to reveal the degree of a word’s polysemy, meaningfully cluster word senses, distinguish homonymy from polysemy or perform superior word sense disambiguation (see Garf Soler and Apidianaki, 2021; Li and Joannis, 2021; Nair et al., 2020; Wiedemann et al., 2019 for each of the

above).

We focus on the topic that received less attention in LM research – the regularity dimension of polysemy and its continuous nature. The definition and scope of regular polysemy vary depending on the linguistic theory. The widely cited definition has been proposed by Apresjan (1974, p. 16) and states that “Polysemy of the word A with the meanings a_i and a_j is called regular if [...] there exists at least one other word B with the meanings b_i and b_j , which are semantically distinguished from each other in exactly the same way as a_i and a_j [...]”. Pustejovsky’s (1991) approach, also frequently adopted, frames regular polysemy as an ability of words that belong to one semantic type to act as members of another, behaving predictably, unlike irregular (accidental) polysemes.

Regular polysemy forms patterns of meaning structure across vocabulary. Some of the widely used examples of such patterns are ANIMAL - MEAT pattern (instantiated by *chicken* or *salmon*) or CONTAINER - CONTENT (e.g. *cup*, *glass*). These examples are an instance of metonymy – a sense extension device that is based on contiguity (association, referential co-existence) of two concepts. The theoretical approaches mentioned above largely attribute regular polysemy to this figure, and so do the researchers in computational linguistics who adopt these theories (see Section 2 for their overview).

There is, however, another cognitive tool that structures polysemy – metaphor. Unlike metonymy, it is based on analogy, or referential disjunction (Lombard et al., 2023). Regular polysemy by metaphor can be exemplified by such polysemes as *antenna* (insect’s organ, signal transmission device) or *leg* (limb, table support) instantiating the pattern BODY PART - OBJECT PART. The two figures are based on different cognitive mechanisms, have different processing profiles in our brain (Klepousiotou et al., 2012), but, as recent psycholinguistic studies show, they equally govern polysemous

sense extensions (Lombard et al., 2023, 2024).

Another important aspect of regular polysemy is its continuous nature. In a recent study, Lombard et al. (2024) introduce a method to extract regular polysemes (including metaphors) from WordNet and suggest metrics to measure the degree of regularity of the patterns they are governed by (Table 1). Their findings are in contrast with the widely applied categorical approach to polysemy, where a sense extension of a polyseme is labeled in a binary way, i.e. as either regular or irregular.

Here we adopt this continuous view, aligning with recent work that argues that word meaning, polysemy regularity, and productivity form a continuum rather than discrete representations (Trott and Bergen, 2023; Li, 2024). To the best of our knowledge, no experimental design has previously targeted the graded aspect of regularity in LLMs, although researchers have noted that some patterns seemed more regular or productive than others (Li and Armstrong, 2024). We also contribute by focusing on metaphorically motivated regular polysemy. Only a handful of works in computational linguistics include regular metaphor in their experiments, and even less in the experiments with LMs in particular.

In order to investigate the effect of graded regularity on models’ representation of metaphorically motivated polysemes, we rely on datasets compiled for psycholinguistic studies on human polysemy processing in French and English (Lombard et al., 2023, 2024). The datasets feature semantic neologisms – novel senses of existing words created using polysemy patterns of varying regularity degrees. These are compared against attested, existing polysemes and nonsensical derivations (refer to Table 2 for the examples). Human acceptability assessment confirmed the psychological validity of graded regularity for human processing: the more regular the polysemy pattern, the more acceptable its novel senses. Using surprisal and semantic similarity measures, we aim to find out how closely language model processing of semantic neologisms aligns with human processing, and whether the degree of regularity plays a role in it. With this in mind, we outline the following research questions:

RQ1. Which of the two measures (surprisal or semantic similarity) would be a better proxy for human behaviour in our task? As discussed in Methods section, both proved to have psycholinguistic predictive power, despite operating at different levels of language structure.

RQ2. Are the results consistent across model types and sizes? Oh and Schuler (2022) show, e.g., that larger models do not necessarily deliver more human-like linguistic representations.

RQ3. Do models distinguish between the novel senses based on existing regular polysemy patterns and the senses created using the patterns that do not exist? To match human behaviour, models should be able to discriminate between these groups.

RQ4. Are LMs sensitive to the varying degrees of regularity of polysemy patterns? If their processing matches human ratings, we should expect the models to be less surprised by neologisms from highly regular patterns and vice versa.

RQ5. What type of regularity metrics (as defined in Table 1) are models more sensitive to: count-based or consistency-based? Do word frequency and word length play a role, and how does this compare with data from human evaluators?

In the case of LLMs, evaluating novel senses allows us to test their ability to generalize beyond previously seen material and avoid data contamination. Additionally, on a higher level, we can assess their sensitivity to the polysemy patterns abstracted from concrete, previously seen words.

Our results show that LLMs could discriminate between different sense types and regularity gradations in a human-like way, and overall correlated well with human sense plausibility judgment.

In the following sections we will briefly discuss the existing work on regular polysemy (§2), justify our methodology (§3), present the experiments (§4) and discuss their results (§5).

2 Related Work

Aside from the theoretical frameworks cited in the Introduction, regular polysemy is studied in several areas dealing with language processing.

Psycholinguistics. In psycholinguistics, regular polysemy is addressed in the discussion about the meaning representation in human mind and the nature of restrictions that govern polysemy patterns in language. Many authors defend hybrid approaches to these problems. Rabagliati and Snedeker (2013) suggest that irregular senses are stored separately, while senses that follow regular patterns form core meanings. Analyzing co-predication acceptability and sense similarity of polysemes and homonyms, Haber and Poesio (2020) suggest that senses form groups according to their similarity (in line with Ortega-Andrés and Vicente, 2019), and reject the

idea of a fully underspecified representation. In contrast, [Vicente \(2024\)](#) analyses regular and irregular polysemy along several dimensions and defends the one-representation hypothesis.

In the discussion on whether linguistic conventions or an underlying conceptual structure restrict polysemy patterns, [Srinivasan and Rabagliati \(2015\)](#) propose the “conventions-constrained-by-concepts” model. Their study across 15 languages suggests that while the conceptual structure governs the patterns, the language-specific conventions define senses that instantiate them. A hybrid approach is also supported by the investigations in language learning: [Zhu \(2021\)](#) studies how preschoolers acquire regular metonymies, highlighting their ability to quickly grasp semantic generalizations without extensive prior exposure. Children rely on an early-emerging conceptual structure, although at later stages linguistic generalizations also play a crucial role in word learning.

Mental processing of ambiguous words is affected by the degree of relatedness of meanings in memory. This is demonstrated by [Brocher \(2016; 2018\)](#), who report increased processing effort associated with disambiguation of unrelated meanings.

Computational Linguistics. In this field, regular polysemy is addressed in a variety of works, such as [Boleda et al. \(2012a,b\)](#); [Lopukhina and Lopukhin \(2016\)](#), who model systematic polysemy, or [Del Tredici and Bel \(2015\)](#), exploring the representations of polysemous and monosemous words in static word embeddings. A number of researchers propose methods of sense annotation for regular polysemy ([Nimb and Pedersen, 2000](#); [Freihat et al., 2013](#); [Martinez Alonso, 2013](#)), while other authors use WordNet to automatically extract regular polysemes ([Peters and Peters, 2000](#); [Barque and Chaumartin, 2009](#); [Lombard et al., 2024](#)). Interestingly, the latter authors recognize metaphoric extensions as types of systematic polysemy patterns, in contrast to most of the previously mentioned studies. [Peters and Peters \(2000\)](#) depart from an assumption that metaphoric alternations are irregular, but after applying their extraction method, “stumble upon” the instances of metaphoric sense extensions that can only be described as regular. Only a few more works mentioned in this section fully recognize that regular polysemy by metaphor is possible: [Nimb and Pedersen, 2000](#); [Freihat et al., 2013](#); [Lopukhina and Lopukhin, 2016](#); [Lombard et al., 2023](#) and [Lombard et al., 2024](#).

Language models. Regarding regular polysemy

and neural language models, [Haber and Poesio \(2021\)](#) test BERT’s ability to predict human assessment of sense similarity degree. They report that BERT_{LARGE} captures distinctions between polysemic, homonymic and same-sense samples in a human-like way. BERT delivers sensible results in sense clustering, suggesting that this model is sensitive to polysemy patterns. [Sørensen et al. \(2023\)](#) explore BERT sense clustering as a guidance tool for annotation of systematic polysemy in lexical resources. Similarly to [Haber and Poesio \(2021\)](#), they got mixed results but see potential: for one of the patterns, BERT discovered a sense that the authors overlooked when creating the dataset. Finally, [Li and Armstrong \(2024\)](#) use sense analogy questions to investigate how regular polysemy is represented in BERT embeddings. The authors observe that the pattern of BERT’s sense similarity score distribution reflects differences not only in the processing of regular polysemes and irregular/homonymous controls, but also of distinct polysemy patterns. They also note on the scalar nature of regularity, an observation that contributes to Li’s (2024) comprehensive approach to polysemy as continuous in its sense individuation, regularity, and productivity.

The present paper adopts the recent insights about the graded nature and metaphoric motivation of regular patterns and incorporates them in the experimental design.

3 Materials and Methods

3.1 Data

To answer our research questions, we evaluated two datasets compiled by [Lombard et al. \(2023\)](#) and [Lombard et al. \(2024\)](#)¹. Both data sets were created for psycholinguistic experiments investigating the effect of graded regularity on the human perception of neology in English and French².

The stimuli. The datasets contain sentences with target words of three types:

1. Semantic neology: words used in a novel, unattested sense. The derived metaphoric sense, together with the base sense, represent a polysemy pattern that a given word has never developed, unlike other words from its semantic field. To ex-

¹Licensed under Creative Commons Attribution 4.0 International (CC-BY-4.0)

²The more recent study is in English and focuses solely on regular metaphor, whereas the earlier one is in French and involves both metaphor and metonymy. Since the present research focuses on metaphoric polysemy, we only evaluate the part of the French dataset containing metaphors.

Metr.	Definition	Formula
R1	Number of words having SENSE ₁ and SENSE ₂ in a given pattern.	$R_1 = N_{S2}$
R2	Ratio of R1 and the number of words with SENSE ₁ , whether or not they have SENSE ₂ .	$R_2 = \frac{N_{S2}}{N_{S1}}$
R3	R1 weighted by the log-frequency of occurrence of the word.	$R_3 = \sum_{w=1}^{N_{S2}} \log(f_w)$
R4	R2 weighted by the log-frequency of occurrence of the word	$R_4 = \frac{\sum_{w=1}^{N_{S2}} \log(f_w)}{\sum_{w=1}^{N_{S1}} \log(f_w)}$

Table 1: Regularity metrics as proposed by (Lombard et al., 2024, pp. 4–5). While R1 and R3 capture the number of pattern instantiations, R2 and R4 reflect the consistency with which words having a base sense (SENSE₁) also have a derived sense (SENSE₂) within a pattern.

Type	Pattern	Example	W.	S.
new	ANIMAL - ARTIFACT	My sister cleaned the porcupine of the brush.	35	70
	ANIMAL - PERSON	The chessplayer is always a cruel spider with his opponents.		
	ARTIFACT - MESSAGE	A mean spear slipped through her lips in an angry tone.		
	BODY PART - OBJECT PART	We can see the knee of the chair getting damaged.		
	NATURAL EVENT - HAPPENING	There was a huge tornado of claps at the final of the challenge.		
	PERSON - ANIMAL	Some zoos are trying to protect the doctor from extinction.		
	PHYS. PROP. - PSYCHOL. PROP.	She said that the density of the project was an issue.		
illegal		My brother painted the curry of the controller in blue.	40	80
existing		My dog chewed the tongue of my new shoes	40	40
all			115	190

Table 2: Sentence examples of each sense type, labeled in the original dataset as *new*, *illegal*, and *existing*. *New* senses include 7 polysemy patterns (5 words per pattern). *Illegal* and *existing* senses are not annotated with patterns in the original dataset. The column *W.* lists the number of words per sense type, while *S.* – the number of sentences.

emply, the word *knee* represents a pattern BODY PART-OBJECT PART and is used in the sentence *We can see the knee of the chair getting damaged*. For comparison, some of the words that actually developed both senses are *leg*, *heart*, *artery*, *vein*, *antenna*, *wing*, *head*, *skeleton*, *brow*, *tongue* etc.

2. Non-sensical derivation: semantic neologisms that follow a non-existent pattern in each language. For instance, *curry* in *My brother painted the curry of the controller in blue* represents an unattested pattern FOOD-OBJECT PART.

3. Existing polysemy: words used in an attested sense of a valid, existing polysemy pattern. For example, *tongue* in *My dog chewed the tongue of my new shoes* is used in an attested sense of an OBJECT PART. An overview of the English dataset with sentence examples is presented in Table 2.

The dataset is annotated with human acceptability scores, regularity degree of polysemy patterns, word frequency and word length.

Human acceptability rating. Human acceptability scores are derived from the initial psycholinguistic experiment. They reflect how plausible the annotators found each sentence on a scale from ‘no

sense at all’ (0) to ‘completely acceptable’ (100).

Regularity. Each target word is annotated with a score reflecting the degree of regularity of a polysemy pattern it instantiates. For the two languages, this metric has been calculated using different procedures. For English, the authors developed an automatic extraction technique using WordNet and proposed several formulas to calculate the regularity degree of a pattern based on the extracted data. These regularity metrics are summarized in Table 1. For French, the authors relied on the judgment of experts in French lexicology to assess the degree of regularity for each pattern. The methodological differences in the compilation of both datasets seem to affect our results, which will be discussed in more detail in Section 4.2.

3.2 Methods

To answer our research questions, we explore two common methods in NLP and computational psycholinguistics – surprisal and semantic similarity from large language models.

Surprisal. Surprisal is the negative log-probability of a token given its immediate con-

text. Surprisal theory (Hale, 2001; Levy, 2008) assumes that the processing difficulty of the word is based on its predictability. This information-theoretic measure is typically used in the studies on human reading, where it proved to predict reading times and, consequently, cognitive processing difficulty in multiple languages (for recent work, see de Varda and Marelli, 2022; Nair and Resnik, 2023; Wilcox et al., 2023; Xu et al., 2023). It is also used to assess the models’ ability to predict linguistic acceptability (grammaticality) of sentences (Noh et al., 2024). In our study, we use surprisal from language models as a proxy of human acceptability judgment of novel word senses: we assume that higher surprisal values assigned to a target word by an LM correspond to lower acceptability scores obtained from human evaluation.

Semantic relatedness. Semantic similarity between a word and its context is used along with surprisal to predict reading times, assess processing difficulty and explain brain activity during language processing (Leal et al., 2021; Salicchi et al., 2021; Kun et al., 2023). Specifically, we apply the cosine similarity between the vector of the target word and the vector of the sentence obtained by mean-pooling. Additionally, since a few rogue dimensions often dominate similarity measures in transformer models (Timkey and van Schijndel, 2021), we compare the original and normalized vectors (z -scoring) to assess their impact. We also use Spearman’s ρ as a similarity metric, another technique suggested by Timkey and van Schijndel (2021) and replicated by Lyu et al. (2023) and Salicchi et al. (2023).

In reading experiments, low similarity between a word and its context is associated with increased human reading difficulty. In our study, we expect to associate low similarity rating from LMs with low human acceptability of semantic neologisms.

As shown by Salicchi et al. (2023), both surprisal and semantic relatedness equally contribute to the prediction of reading difficulty, despite operating at different levels of language structure. While surprisal operates at the syntagmatic level and reflects how predictable the word is from its context, semantic relatedness reflects coherence of a word with its context modeling paradigmatic dimension. Both surprisal and semantic relatedness proved to predict brain activity during language comprehension and are associated with signals from distinct brain areas (Frank and Willems, 2017; Michaelov et al., 2023; Salicchi and Hsu, 2025).

3.3 Models

We used a set of masked language models and compared them with an autoregressive Llama.

For English, we use the monolingual BERT as well as RoBERTa. For French, we took the BERT-based FlauBERT, and the RoBERTa-based CamembERTv2. We also evaluate multilingual models on both languages: mBERT and XLM-RoBERTa.

Surprisal experiments typically use unidirectional decoder models (e.g., GPT), as they rely only on left-context to emulate human reading, avoiding access to future words. In our case, the experimental settings of the initial psycholinguistic study entail the choice of a masked model: the evaluators were first presented with the context on both sides before seeing the full sentence. We still include an autoregressive LM to compare the results and challenge our assumption about masked language modeling being more suitable for our task. For this, we chose Llama 3.1 8B and Llama 3.2 3B, which we oppose to BERT as more recent and significantly larger multilingual models that include English and French. For all models, weights were taken off HuggingFace. Additional information on these models is presented in Table 5 of Appendix A.

4 Experiments and Results

4.1 Experiments

We feed the sentences into each of the language models and compute³ the surprisal and semantic relatedness scores as described in Methods section. For this, we use the minicons library provided by Misra (2022). For bi-directional models, we rely on the ‘pseudo-log-likelihood’ proposed by Kauf and Ivanova (2023), which takes into account multi-token and out-of-vocabulary words.⁴

We then compute the Spearman correlation between the human acceptability scores and each of the measures (target word surprisal and the similarity between the target word and its context).

4.2 Results

Surprisal. Across models, we observe moderate to strong correlation with human judgment. As expected, models correlate negatively, showing that more acceptable senses elicit lower surprisal.

³The information on GPU use and computation time is reported in Appendix B.

⁴We also tested the standard scoring based on Salazar et al., 2020, and that of PsychoFormers (Michaelov and Bergen, 2022), obtaining generally lower results, as shown in Table 3.

Models	Default	K & I	PF
BERT _{BASE}	-0.65	-0.63	-0.61
BERT _{LARGE}	-0.68	-0.65	-0.64
RoBERTa _{BASE}	-0.67	-0.76	-0.68
RoBERTa _{LARGE}	-0.72	-0.78	-0.70
XLM-RoBERTa _{BASE}	-0.38	-0.56	-0.39
XLM-RoBERTa _{LARGE}	-0.44	-0.63	-0.43
mBERT _{BASE}	-0.26	-0.47	-0.47
Llama 3.1 8B	-0.65	-	-
Llama 3.2 3B	-0.65	-	-

Table 3: Results of the surprisal experiment in English. The column *Default* reports results obtained from the default implementation of minicons (Misra, 2022), the column *K & I* reports the results from the method by Kauf and Ivanova (2023), and *PF* – from the PsychFormers application (Michaelov et al., 2023). All results are statistically significant. Bold formatting points to the strongest correlation achieved by each model.

Among masked models, the strongest correlation was achieved by RoBERTa_{LARGE} at -0.78, $p < .001$. It is followed by BERT_{LARGE} showing moderate negative correlation of -0.68, $p < .001$. Multilingual models demonstrated poorer results with correlation coefficients of -0.63 for XLM-RoBERTa_{LARGE} ($p < .001$), as well as -0.47 for mBERT ($p < .001$). As mentioned previously, the method of Kauf and Ivanova (2023) yielded the best results, except for the BERT models which performed slightly better using the standard metric. See Table 3 for a complete overview of the different models and metrics.

As for autoregressive models, Llama 3.1 8B and Llama 3.2 3B achieved correlation of -0.65 ($p < .001$ for both), yielding the best results among the multilingual models but exhibiting a lower correlation than the smaller monolingual encoders.

In French, none of the models gave statistically significant correlation at the word level. We attribute this to the much smaller dataset size (42 sentences). However, we could still obtain usable results by changing the experimental settings: we checked correlation of sentence-wise surprisal with human judgment (obtained by sum and mean) and received statistically significant results for XLM-RoBERTa_{LARGE}, at -0.32, $p = .039$ (sum). We compared this result with the sentence surprisal of the English version from XLM-RoBERTa, and curiously, for English, this was the only model that showed stronger correlation when computing sentence surprisal instead of the target word surprisal (-0.65 vs. -0.63, $p < .001$ in both cases). Table 6 in Appendix C presents all scores obtained from

the sentence-wise correlation experiment. Additionally, it reports correlation of sentence surprisal with the acceptability of polysemy patterns, where XLM-RoBERTa_{LARGE} achieved moderate significant correlation.

Semantic relatedness. The results of the experiment with semantic relatedness are more difficult to summarize, as the data does not allow to discern clear trends. In different models, the highest correlation was achieved across varying layers, model sizes and normalization approaches. Moreover, some models show positive correlation with human judgment, while others correlate negatively. This is not expected, as usually we assume a better word/context coherence to elicit higher acceptability scores. We will highlight the best results to give an idea of the final picture, while Tables 7 to 10 of Appendix D offer an additional illustration of correlation scores distribution within several selected models: masked RoBERTa and FlauBERT for English and French, as well as a significantly bigger multilingual autoregressive Llama 3.1 8B.

The strongest correlation was reached by Llama 3.1 8B (32 layers) in the layer 4 using Spearman’s ρ instead of cosine, the correlation being positive (0.66, $p < .001$). RoBERTa_{BASE} (12 layers) follows with coefficient of -0.58, $p < .001$ in the ninth layer without applying any normalization techniques. Finally, BERT_{BASE} (12 layers) achieved the correlation of 0.52 in the last layer when applying Spearman’s ρ instead of the cosine ($p < .001$). Multilingual models score 0.5 and below, their best achieved correlation coefficients being scattered across different experimental settings.

For French, the scores lie in the same range, but with the strongest correlation achieved by a smaller Llama 3.2 3B (-0.53, $p < .001$, in the last 28th layer, non-normalized).

Lyu et al. (2023) report similar outcome of their study of lexical stylistic features in language models: although normalization generally improves the results (especially for the multilingual models), it is hard to single out the best technique for all models and experimental settings. As for Salicchi et al. (2023), they do not notice any effect of BERT’s embedding anisotropy on reading times prediction.

Overall, semantic relatedness results show no clear interpretable trend across models and settings.

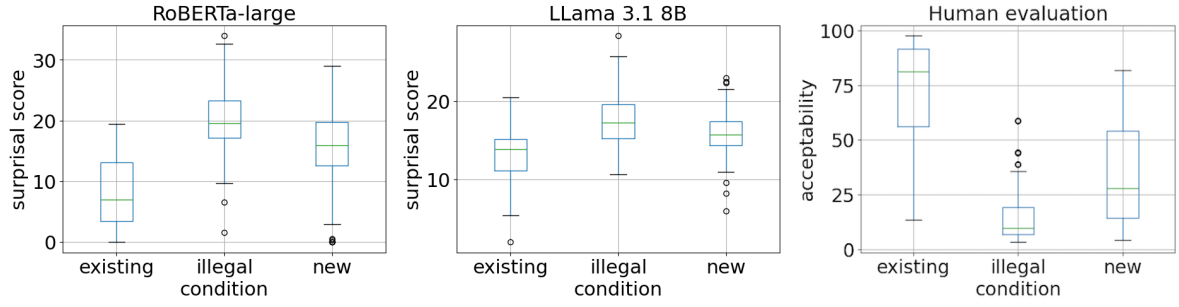


Figure 1: Distribution of English surprisal scores by condition labeled in the original dataset as *new*, *illegal* and *existing*. These correspond to the groups (1), (2) and (3) respectively, as described in the Section 3.1.

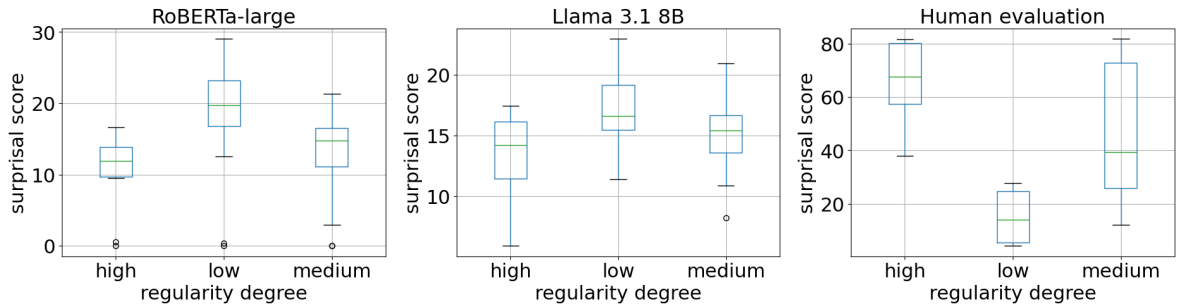


Figure 2: Distribution of surprisal scores by regularity degrees labeled as *high*, *medium* and *low* for the English data.

5 Discussion

In this section, we will address the research questions presented in the Introduction.

RQ1. Regarding the choice of measure (surprisal or semantic relatedness), the results suggest that surprisal is preferable. Not only because it achieved strong correlation (-0.78 for RoBERTa_{LARGE} in surprisal vs. -0.61 for Llama 3.2 3B in similarity setting), but also because it is consistent, more interpretable, and easier to obtain. While we can confirm the assumption that surprisal is in an inverse relationship with sense plausibility as assessed by humans, the semantic similarity scores correlated both negatively and positively depending on the model, its layer and the embedding normalization technique (see Tables in Appendix D). Finding the most suitable configuration thus demands running a considerable number of trials.

RQ2. In surprisal setting, masked LMs performed better, confirming our assumption that masked model scoring with its access to the bi-directional context would be more suitable for our task. Previous research has repeatedly shown that larger model size delivers a poorer prediction of processing difficulty (Oh and Schuler, 2022; Salicchi et al., 2023; Liu et al., 2024; Shain et al., 2024). In contrast, in our experiments, large varieties of the same models always performed above the base ones (see

Table 3). Interestingly, much larger Llama 3.1 8B and 3.2 3B did not outperform masked monolingual BERT and RoBERTa (330M and 355M respectively for large varieties). We attribute this to differences in model architecture, although it requires further investigation. In the case of semantic similarity, results are not consistent enough to draw conclusions on this topic, as explained in Results section.

We further analyse the results to test the models' sensitivity to such features as sense types, regularity degrees, word frequency and word length.

RQ3. We run a series of tests to confirm whether the models discriminate between the senses derived using the existing and non-existing patterns, as well as to see if they are sensitive to the varying pattern regularity degrees. We took our best-performing masked model RoBERTa_{LARGE} and an autoregressive Llama 3.1 8B, for comparison. For French, we picked the same Llama model and FlauBERT_{LARGE}. A Mann-Whitney U test on two independent samples for the two sense types (two-sided, $p < 0.05$) shows that the difference is significant for RoBERTa, Llama and human evaluators. They could distinguish between all three groups of senses (Figure 1). For French, FlauBERT and Llama 3.1 8B did not yield significance (see Figure 4 in Appendix E for score distribution).

RQ4. In the same way, we established that language models were sensitive to the degrees of regularity of the polysemy patterns the senses instantiated, although not as fine-grained as humans: while the Mann-Whitney U test shows significance in the difference between low, medium and high regularity of patterns for humans, the models only discriminate between high/low and medium/low groups (Figure 2). Again, neither Llama nor FlauBERT reached statistical significance in French. The distribution plot can be found in Appendix E, Figure 5.

RQ5. We also establish whether there is a relation between the model scoring and such factors as the degree of polysemy pattern regularity, word frequency and word length. The latter two factors contribute to the cognitive processing load in humans since less frequent and longer words require more time to process (Pollatsek et al., 2008). Figure 3 illustrates the pattern of correlation (Spearman) between four regularity metrics and the measures of semantic similarity and surprisal from RoBERTa and Llama, as compared to human evaluators. As described in Table 1, we consider a count-based metric R1, consistency-based metric R2 and two metrics that weight them by the log-frequency of the occurrence of the word – R3 and R4. The measure of surprisal closely follows the line representing human judgment correlated with regularity, with Llama showing almost identical coefficients. The same as for human evaluators, for both language models, the regularity metrics that reflect how consistently words instantiate a polysemy pattern appeared to be more relevant than the sheer number of words having the SENSE₁ and the SENSE₂. Weighting R1 and R2 by word frequency generally did not improve the correlation coefficients (except for RoBERTa in R4 where it gains one point). Again, the correlation scores for the measure of semantic relatedness are generally low, with apparent preference for consistency-based and frequency-weighted metrics. All correlation coefficients are listed in Table 11, Appendix E.

As for the effect of the word frequency and word length, the models in general do not show high correlation, which is in line with the results from human evaluation. The two exceptions are RoBERTa in the similarity setting and FlauBERT in the surprisal setting relying on these features more and correlating moderately (see Table 4 for correlation scores and Figure 6 in Appendix E for visualization).

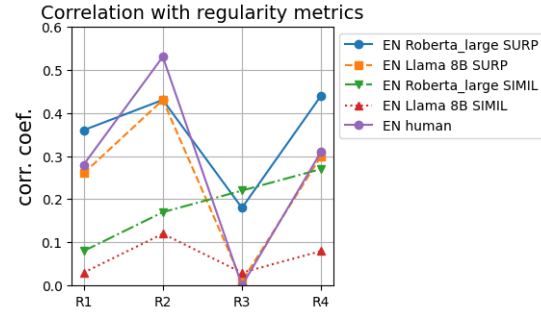


Figure 3: Spearman’s ρ for regularity metrics as described in Section 5, in absolute numbers.

English			
Models	W. freq.	W. length	
RoBERTa _{LARGE} SURP	0.16	0.11	
EN Llama 3.1 8B SURP	0.27*	0.01	
RoBERTa _{LARGE} SIMIL	0.4*	0.04	
Llama 3.1. 8B SIMIL	0.04	0.16	
Human acceptability rating	0.17*	0.18*	
French			
FlauBERT _{LARGE} SURP	0.38*	0.41*	
Llama 3.1 8B SURP	0.19	0.1	
FlauBERT _{LARGE} SIMIL	0.06	0.07	
Llama 3.1 8B SIMIL	0.15	0.06	
Human acceptability rating	0.15	0.17	

Table 4: Correlation of word length and word frequency with model scoring and human evaluation.

6 Conclusions

In this paper, we investigated the effect of the graded regularity of polysemy patterns on the processing of novel metaphorical word senses by large language models. Using surprisal and semantic relatedness as proxies, we found evidence that models represent regularity of polysemy extensions in a human-like way. Especially surprisal proved to adequately model sense plausibility, showing a strong correlation with human judgment. Among models, RoBERTa delivered the best results. Furthermore, the distributions of model scores suggest sensitivity to different types of sense extensions and regularity degrees. Similarly to humans, LLMs could discriminate between attested polysemes, novel senses derived from regular polysemy patterns and non-sensical derivations. They were, however, less responsive to the gradations in regularity, only differentiating very regular and weakly regular patterns. These observations allow us to better understand how LLMs model lexical ambiguity and to what extent such factors as regularity, continuity and sense relatedness affect model representations.

References

- Wissam Antoun, Francis Kulumba, Rian Touchent, Éric de la Clergerie, Benoît Sagot, and Djamé Seddah. 2024. [Camembert 2.0: A smarter french language model aged to perfection](#). *Preprint*, arXiv:2411.08868.
- Jurij D. Apresjan. 1974. [Regular polysemy](#). *Linguistics*, 12(142):5–32.
- Lucie Barque and François-Régis Chaumartin. 2009. Regular polysemy in WordNet. *Journal for language technology and computational linguistics*, 24(2):5–18.
- Gemma Boleda, Sebastian Padó, and Jason Utt. 2012a. [Regular polysemy: A distributional model](#). In **SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 151–160, Montréal, Canada. Association for Computational Linguistics.
- Gemma Boleda, Sabine Schulte im Walde, and Toni Badia. 2012b. [Modeling regular polysemy: A study on the semantic classification of Catalan adjectives](#). *Computational Linguistics*, 38(3):575–616.
- Andreas Brocher, Stephani Foraker, and Jean-Pierre Koenig. 2016. [Processing of irregular polysemes in sentence reading](#). *Journal of experimental psychology. Learning, memory, and cognition*, 42 11:1798–1813.
- Andreas Brocher, Jean-Pierre Koenig, Gail Mauner, and Stephani Foraker. 2018. [About sharing and commitment: the retrieval of biased and balanced irregular polysemes](#). *Language, Cognition and Neuroscience*, 33:443 – 466.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Unsupervised cross-lingual representation learning at scale](#). *CoRR*, abs/1911.02116.
- Andrea de Varda and Marco Marelli. 2022. [The effects of surprisal across languages: Results from native and non-native reading](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2022*, pages 138–144, Online only. Association for Computational Linguistics.
- Marco Del Tredici and Núria Bel. 2015. [A word-embedding-based sense index for regular polysemy representation](#). In *Proceedings of the 1st Workshop on Vector Space Modeling for Natural Language Processing*, pages 70–78, Denver, Colorado. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Kevin Durkin and Jocelyn Manning. 1989. [Polysemy and the subjective lexicon: Semantic relatedness and the salience of intraword senses](#). *Journal of Psycholinguistic Research*, 18:577–612.
- S. Frank and Roel M. Willems. 2017. [Word predictability and semantic similarity show distinct patterns of brain activity during language comprehension](#). *Language, Cognition and Neuroscience*, 32:1192 – 1203.
- Abed Alhakim Freihat, Fausto Giunchiglia, and Biswanath Dutta. 2013. Regular polysemy in wordnet and pattern based approach. *International Journal On Advances in Intelligent Systems*, 6:199–212.
- Aina Garf Soler and Marianna Apidianaki. 2021. [Let’s play mono-poly: BERT can reveal words’ polysemy level and partitionability into senses](#). *Transactions of the Association for Computational Linguistics*, 9:825–844.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Alonsoius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der

750	Maaten, Lawrence Chen, Liang Tan, Liz Jenkins,	ban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun,	814
751	Louis Martin, Lovish Madaan, Lubo Malo, Lukas	Felix Kreuk, Feng Tian, Filippas Kokkinos, Firat	815
752	Blecher, Lukas Landzaat, Luke de Oliveira, Madeline	Ozgenel, Francesco Caggioni, Frank Kanayet, Frank	816
753	Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar	Seide, Gabriela Medina Florez, Gabriella Schwarz,	817
754	Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew	Gada Badeer, Georgia Swee, Gil Halpern, Grant	818
755	Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-	Herman, Grigory Sizov, Guangyi, Zhang, Guna	819
756	badur, Mike Lewis, Min Si, Mitesh Kumar Singh,	Lakshminarayanan, Hakan Inan, Hamid Shojanaz-	820
757	Mona Hassan, Naman Goyal, Narjes Torabi, Niko-	eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun	821
758	lay Bashlykov, Nikolay Bogoychev, Niladri Chatterji,	Habeeb, Harrison Rudolph, Helen Suk, Henry As-	822
759	Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick	pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim	823
760	Alrassy, Pengchuan Zhang, Pengwei Li, Petar Va-	Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis,	824
761	sic, Peter Weng, Prajjwal Bhargava, Pratik Dubal,	Irina-Elena Veliche, Itai Gat, Jake Weissman, James	825
762	Praveen Krishnan, Punit Singh Koura, Puxin Xu,	Geboski, James Kohli, Janice Lam, Japhet Asher,	826
763	Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj	Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-	827
764	Ganapathy, Ramon Calderer, Ricardo Silveira Cabral,	nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy	828
765	Robert Stojnic, Roberta Raileanu, Rohan Maheswari,	Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe	829
766	Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-	Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-	830
767	nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan	Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang,	831
768	Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-	Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khan-	832
769	hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-	delwal, Katayoun Zand, Kathy Matosich, Kaushik	833
770	hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sha-	Veeraraghavan, Kelly Michelen, Keqian Li, Ki-	834
771	ran Narang, Sharath Raparthy, Sheng Shen, Shengye	ran Jagadeesh, Kun Huang, Kunal Chawla, Kyle	835
772	Wan, Shruti Bhosale, Shun Zhang, Simon Van-	Huang, Lailin Chen, Lakshya Garg, Lavender A,	836
773	denhende, Soumya Batra, Spencer Whitman, Sten	Leandro Silva, Lee Bell, Lei Zhang, Liangpeng	837
774	Sootla, Stephane Collot, Suchin Gururangan, Syd-	Guo, Licheng Yu, Liron Moshkovich, Luca Wehrst-	838
775	ney Borodinsky, Tamar Herman, Tara Fowler, Tarek	edt, Madian Khabsa, Manav Avalani, Manish Bhatt,	839
776	Sheasha, Thomas Georgiou, Thomas Scialom, Tobias	Martynas Mankus, Matan Hasson, Matthew Lennie,	840
777	Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal	Matthias Reso, Maxim Groshev, Maxim Naumov,	841
778	Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh	Maya Lathi, Meghan Keneally, Miao Liu, Michael L.	842
779	Ramanathan, Viktor Kerkez, Vincent Gonguet, Vir-	Seltzer, Michal Valko, Michelle Restrepo, Mihir Pa-	843
780	ginie Do, Vish Vogeti, Vitor Albiero, Vladan Petro-	tel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark,	844
781	vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-	Mike Macey, Mike Wang, Miquel Jubert Hermoso,	845
782	ney Meers, Xavier Martinet, Xiaodong Wang, Xi-	Mo Metanat, Mohammad Rastegari, Munish Bansal,	846
783	aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xin-	Nandhini Santhanam, Natascha Parks, Natasha	847
784	feng Xie, Xuchao Jia, Xuewei Wang, Yaelle Gold-	White, Navyata Bawa, Nayan Singhal, Nick Egebo,	848
785	schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen,	Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich	849
786	Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao,	Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz,	850
787	Zacharie Delpierre Coudert, Zheng Yan, Zhengxing	Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin	851
788	Chen, Zoe Papakipos, Aaditya Singh, Aayushi Sri-	Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pe-	852
789	vastava, Abha Jain, Adam Kelsey, Adam Shajnfeld,	dro Rittner, Philip Bontrager, Pierre Roux, Piotr	853
790	Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,	Dollar, Polina Zvyagina, Prashant Ratanchandani,	854
791	Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei	Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel	855
792	Baevski, Allie Feinstein, Amanda Kallet, Amit San-	Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu	856
793	gani, Amos Teo, Anam Yunus, Andrei Lupu, An-	Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,	857
794	dres Alvarado, Andrew Caples, Andrew Gu, Andrew	Raymond Li, Rebekkah Hogan, Robin Battey, Rocky	858
795	Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-	Wang, Russ Howes, Ruty Rinott, Sachin Mehta,	859
796	dani, Annie Dong, Annie Franco, Anuj Goyal, Apar-	Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara	860
797	jita Saraf, Arkabandhu Chowdhury, Ashley Gabriel,	Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov,	861
798	Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-	Satadru Pan, Saurabh Mahajan, Saurabh Verma,	862
799	dan, Beau James, Ben Maurer, Benjamin Leonhardi,	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-	863
800	Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,	864
801	Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-	Shengxin Cindy Zha, Shishir Patil, Shiva Shankar,	865
802	cock, Bram Wasti, Brandon Spence, Brani Stojkovic,	Shuqiang Zhang, Shuqiang Zhang, Sinong Wang,	866
803	Brian Gamido, Britt Montalvo, Carl Parker, Carly	Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala,	867
804	Burton, Catalina Mejia, Ce Liu, Changhan Wang,	Stephanie Max, Stephen Chen, Steve Kehoe, Steve	868
805	Changkyu Kim, Chao Zhou, Chester Hu, Ching-	Satterfield, Sudarshan Govindaprasad, Sumit Gupta,	869
806	Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-	Summer Deng, Sungmin Cho, Sunny Virk, Suraj	870
807	ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty,	Subramanian, Sy Choudhury, Sydney Goldman, Tal	871
808	Daniel Kreymer, Daniel Li, David Adkins, David	Remez, Tamar Glaser, Tamara Best, Thilo Koehler,	872
809	Xu, Davide Testuggine, Delia David, Devi Parikh,	Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim	873
810	Diana Liskovich, Didem Foss, Dingkan Wang, Duc	Matthews, Timothy Chou, Tzook Shaked, Varun	874
811	Le, Dustin Holland, Edward Dowling, Eissa Jamil,	Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai	875
812	Elaine Montgomery, Eleonora Presani, Emily Hahn,	Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad	876
813	Emily Wood, Eric-Tuan Le, Erik Brinkman, Este-	Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu,	877

878	Vladimir Ivanov, Wei Li, Wenchen Wang, Wen-	reading data. In <i>Text, Speech, and Dialogue</i> , pages	935
879	wen Jiang, Wes Bouaziz, Will Constable, Xiaocheng	35–47, Cham. Springer International Publishing.	936
880	Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo		
881	Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,	Roger Levy. 2008. Expectation-based syntactic compre-	937
882	Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,	hension . <i>Cognition</i> , 106(3):1126–1177.	938
883	Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao,		
884	Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary	Jiangtian Li. 2024. Semantic minimalism and the con-	939
885	DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang,	tinuous nature of polysemy. <i>Mind & Language</i> ,	940
886	Zhiwei Zhao, and Zhiyu Ma. 2024. The Llama 3	39(5):680–705.	941
887	Herd of Models . <i>Preprint</i> , arXiv:2407.21783.		
888	Janosch Haber and Massimo Poesio. 2020. Assessing	Jiangtian Li and Blair C Armstrong. 2024. Probing the	942
889	polyseme sense similarity through co-predication ac-	representational structure of regular polysemy via	943
890	ceptability and contextualised embedding distance .	sense analogy questions: Insights from contextual	944
891	In <i>Proceedings of the Ninth Joint Conference on Lex-</i>	word vectors. <i>Cognitive Science</i> , 48(3):e13416.	945
892	<i>ical and Computational Semantics</i> , pages 114–124,		
893	Barcelona, Spain (Online). Association for Computa-	Jiangtian Li and Marc Joanisse. 2021. Word senses as	946
894	tional Linguistics.	clusters of meaning modulations: A computational	947
		model of polysemy . <i>Cognitive Science</i> , 45.	948
895	Janosch Haber and Massimo Poesio. 2021. Patterns of	Tong Liu, Iza Škrjanec, and Vera Demberg. 2024.	949
896	polysemy and homonymy in contextualised language	Temperature-scaling surprisal estimates improve fit	950
897	models . In <i>Findings of the Association for Computa-</i>	to human reading times – but does it do so for the	951
898	<i>tional Linguistics: EMNLP 2021</i> , pages 2663–2676,	“right reasons”? In Proceedings of the 62nd Annual	952
899	Punta Cana, Dominican Republic. Association for	Meeting of the Association for Computational Lin-	953
900	Computational Linguistics.	guistics (Volume 1: Long Papers) , pages 9598–9619,	954
		Bangkok, Thailand. Association for Computational	955
901	Janosch Haber and Massimo Poesio. 2024. Polysemy—	Linguistics.	956
902	Evidence from linguistics, behavioral science, and	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	957
903	contextualized language models . <i>Computational Lin-</i>	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	958
904	<i>guistics</i> , 50(1):351–417.	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	959
905	John Hale. 2001. A probabilistic Earley parser as a psy-	RoBERTa: A Robustly Optimized BERT Pretrain-	960
906	cholinguiistic model . In <i>Second Meeting of the North</i>	ing Approach . <i>arXiv preprint arXiv:1907.11692</i> .	961
907	<i>American Chapter of the Association for Computa-</i>		
908	<i>tional Linguistics</i> .	Alizée Lombard, Richard Huyghe, Lucie Barque, and	962
		Doriane Gras. 2023. Regular polysemy and novel	963
909	Carina Kauf and Anna Ivanova. 2023. A better way to	word-sense identification. <i>The Mental Lexicon</i> ,	964
910	do masked language model scoring . In <i>Proceedings</i>	18(1):94–119.	965
911	<i>of the 61st Annual Meeting of the Association for</i>		
912	<i>Computational Linguistics (Volume 2: Short Papers)</i> ,	Alizée Lombard, Anastasia Ulicheva, Maria Korochk-	966
913	pages 925–935, Toronto, Canada. Association for	ina, and Kathy Rastle. 2024. The regularity of poly-	967
914	Computational Linguistics.	semy patterns in the mind: Computational and exper-	968
		imental data . <i>Glossa Psycholinguistics</i> , 3(1).	969
915	Ekaterini Klepousniotou, G. Bruce Pike, Karsten Stein-	Anastasiya Lopukhina and Konstantin Lopukhin. 2016.	970
916	hauer, and Vincent L. Gracco. 2012. Not all ambigu-	Regular polysemy: from sense vectors to sense pat-	971
917	ous words are created equal: An eeg investigation	terns. In <i>Proceedings of the 5th Workshop on Cog-</i>	972
918	of homonymy and polysemy . <i>Brain and Language</i> ,	<i>nitive Aspects of the Lexicon (CogALex-V)</i> , pages	973
919	123:11–21.	19–23.	974
920	Sun Kun, Qiuying Wang, and Xiaofei Lu. 2023. An	Qing Lyu, Marianna Apidianaki, and Chris Callison-	975
921	interpretable measure of semantic similarity for pre-	burch. 2023. Representation of lexical stylistic fea-	976
922	dicting eye movements in reading . <i>Psychonomic</i>	tures in language models’ embedding space . In <i>Pro-</i>	977
923	<i>Bulletin & Review</i> , 30:1227 – 1242.	<i>ceedings of the 12th Joint Conference on Lexical and</i>	978
924	Hang Le, Loïc Vial, Jibril Frej, Vincent Segonne, Max-	<i>Computational Semantics (*SEM 2023)</i> , pages 370–	979
925	imin Coavoux, Benjamin Lecouteux, Alexandre Al-	387, Toronto, Canada. Association for Computational	980
926	lauzen, Benoît Crabbé, Laurent Besacier, and Didier	Linguistics.	981
927	Schwab. 2020. Flaubert: Unsupervised language	Hector Martinez Alonso. 2013. <i>Annotation of Regular</i>	982
928	model pre-training for french . In <i>Proceedings of The</i>	<i>Polysemy: An empirical assessment of the underspec-</i>	983
929	<i>12th Language Resources and Evaluation Confer-</i>	<i>ified sense</i> . Ph.D. thesis, Universitat Pompeu Fabra	984
930	<i>ence</i> , pages 2479–2490, Marseille, France. European	and University of Copenhagen, Barcelona, Spain.	985
931	Language Resources Association.	James A. Michaelov, Megan D. Bardolph, Cyma K. Van	986
932	Sidney Leal, Edresson Casanova, Gustavo Paetzold, and	Petten, Benjamin K. Bergen, and Seana Coulson.	987
933	Sandra Aluísio. 2021. Evaluating semantic similarity	2023. Strong prediction: Language model surprisal	988
934	methods to build semantic predictability norms of		

989	explains multiple n400 effects. <i>Neurobiology of Lan-</i>		
990	<i>guage</i> , 5:107 – 135.		
991	James A. Michaelov and Benjamin K. Bergen. 2022. Do		
992	language models make human-like predictions about		
993	the coreferents of Italian anaphoric zero pronouns?		
994	In <i>Proceedings of the 29th International Conference</i>		
995	<i>on Computational Linguistics</i> , pages 1–14, Gyeongju,		
996	Republic of Korea. International Committee on Com-		
997	putational Linguistics.		
998	Kanishka Misra. 2022. minicons: Enabling flexible be-		
999	havioral and representational analyses of transformer		
1000	language models. <i>arXiv preprint arXiv:2203.13112</i> .		
1001	Sathvik Nair and Philip Resnik. 2023. Words, subwords,		
1002	and morphemes: What really matters in the surprisal-		
1003	reading time relationship? In <i>Findings of the As-</i>		
1004	<i>sociation for Computational Linguistics: EMNLP</i>		
1005	2023, pages 11251–11260, Singapore. Association		
1006	for Computational Linguistics.		
1007	Sathvik Nair, Mahesh Srinivasan, and Stephan Mey-		
1008	lan. 2020. Contextualized word embeddings encode		
1009	aspects of human-like word sense knowledge. In <i>Pro-</i>		
1010	<i>ceedings of the Workshop on the Cognitive Aspects</i>		
1011	<i>of the Lexicon</i> , pages 129–141, Online. Association		
1012	for Computational Linguistics.		
1013	Sanni Nimb and Bolette Sanford Pedersen. 2000. Treat-		
1014	ing metaphoric senses in a danish computational		
1015	lexicon –different cases of regular polysemy. In		
1016	<i>Proceedings of the 9th EURALEX International</i>		
1017	<i>Congress</i> , pages 679–691, Stuttgart, Germany. In-		
1018	stitut für Maschinelle Sprachverarbeitung.		
1019	Kangsang Noh, Eunjeong Oh, and Sanghoun Song. 2024.		
1020	Testing language models’ syntactic sensitivity to		
1021	grammatical constraints: a case study of wanna con-		
1022	traction. <i>Frontiers in Communication</i> , 9.		
1023	Byung-Doh Oh and William Schuler. 2022. Why does		
1024	surprisal from larger transformer-based language		
1025	models provide a poorer fit to human reading times?		
1026	<i>Transactions of the Association for Computational</i>		
1027	<i>Linguistics</i> , 11:336–350.		
1028	Marina Ortega-Andrés and Agustín Vicente. 2019. Pol-		
1029	ysemy and co-predication. <i>Glossa: a journal of gen-</i>		
1030	<i>eral linguistics</i> , 4(1).		
1031	Wim Peters and Ivonne Peters. 2000. Lexicalised sys-		
1032	tematic polysemy in WordNet. In <i>Proceedings of the</i>		
1033	<i>Second International Conference on Language Re-</i>		
1034	<i>sources and Evaluation (LREC’00)</i> , Athens, Greece.		
1035	European Language Resources Association (ELRA).		
1036	Alexander Pollatsek, Barbara Jean Juhasz, Erik D. Re-		
1037	ichle, Debra Machacek, and Keith Rayner. 2008. Im-		
1038	mediate and delayed effects of word frequency and		
1039	word length on eye movements in reading: a reversed		
1040	delayed effect of word length. <i>Journal of experimen-</i>		
1041	<i>tal psychology. Human perception and performance</i> ,		
1042	34 3:726–50.		
1043	James Pustejovsky. 1991. The Generative Lexicon.		
1044	<i>Computational Linguistics</i> , 17(4):409–441.		
	Hugh Rabagliati and Jesse Snedeker. 2013. The truth	1045	
	about chickens and bats: Ambiguity avoidance dis-	1046	
	tinguishes types of polysemy. <i>Psychological science</i> ,	1047	
	24(7):1354–1360.	1048	
	Julian Salazar, Davis Liang, Toan Q. Nguyen, and Ka-	1049	
	trrin Kirchhoff. 2020. Masked language model scor-	1050	
	ing. In <i>Proceedings of the 58th Annual Meeting of</i>	1051	
	<i>the Association for Computational Linguistics</i> , pages	1052	
	2699–2712, Online. Association for Computational	1053	
	Linguistics.	1054	
	Lavinia Salicchi, Emmanuele Chersoni, and Alessandro	1055	
	Lenci. 2023. A study on surprisal and semantic relat-	1056	
	edness for eye-tracking data prediction. <i>Frontiers in</i>	1057	
	<i>Psychology</i> , 14.	1058	
	Lavinia Salicchi and Yu-Yin Hsu. 2025. Not every met-	1059	
	ric is equal: Cognitive models for predicting n400	1060	
	and p600 components during reading comprehen-	1061	
	sion. In <i>Proceedings of the 31st International Con-</i>	1062	
	<i>ference on Computational Linguistics</i> , pages 3648–	1063	
	3654, Abu Dhabi, UAE. Association for Computa-	1064	
	tional Linguistics.	1065	
	Lavinia Salicchi, Alessandro Lenci, and Emmanuele	1066	
	Chersoni. 2021. Looking for a role for word embed-	1067	
	dings in eye-tracking features prediction: Does se-	1068	
	mantic similarity help? In <i>Proceedings of the 14th In-</i>	1069	
	<i>ternational Conference on Computational Semantics</i>	1070	
	<i>(IWCS)</i> , pages 87–92, Groningen, The Netherlands	1071	
	(online). Association for Computational Linguistics.	1072	
	Cory Shain, Clara Meister, Tiago Pimentel, Ryan Cot-	1073	
	terrell, and Roger Levy. 2024. Large-scale evidence	1074	
	for logarithmic effects of word predictability on read-	1075	
	ing time. <i>Proceedings of the National Academy of</i>	1076	
	<i>Sciences of the United States of America</i> , 121.	1077	
	Nathalie Sørensen, Sanni Nimb, and Bolette Sandford	1078	
	Pedersen. 2023. How do we treat systematic poly-	1079	
	semy in wordnets and similar resources?—using hu-	1080	
	man intuition and contextualized embeddings as guid-	1081	
	ance. In <i>Proceedings of the 12th Global Wordnet</i>	1082	
	<i>Conference</i> , pages 117–126.	1083	
	Mahesh Srinivasan and Hugh Rabagliati. 2015. How	1084	
	concepts and conventions structure the lexicon:	1085	
	Cross-linguistic evidence from polysemy. <i>Lingua</i> ,	1086	
	157:124–152.	1087	
	William Timkey and Marten van Schijndel. 2021. All	1088	
	bark and no bite: Rogue dimensions in transformer	1089	
	language models obscure representational quality.	1090	
	In <i>Proceedings of the 2021 Conference on Empiri-</i>	1091	
	<i>cal Methods in Natural Language Processing</i> , pages	1092	
	4527–4546, Online and Punta Cana, Dominican Re-	1093	
	public. Association for Computational Linguistics.	1094	
	Sean Trott and Benjamin Bergen. 2023. Word meaning	1095	
	is both categorical and continuous. <i>Psychological</i>	1096	
	<i>Review</i> , 130(5):1239.	1097	
	Agustín Vicente. 2024. Polysemies and the one repre-	1098	
	sentation hypothesis. <i>The Mental Lexicon</i> .	1099	

- Gregor Wiedemann, Steffen Remus, Avi Chawla, and Chris Biemann. 2019. Does BERT Make Any Sense? Interpretable Word Sense Disambiguation with Contextualized Embeddings. In *Proceedings of the 15th Conference on Natural Language Processing (KONVENS 2019): Long Papers*, pages 161–170, Erlangen, Germany. German Society for Computational Linguistics & Language Technology.
- Ethan G. Wilcox, Tiago Pimentel, Clara Meister, Ryan Cotterell, and Roger P. Levy. 2023. [Testing the Predictions of Surprisal Theory in 11 Languages](#). *Transactions of the Association for Computational Linguistics*, 11:1451–1470.
- Weijie Xu, Jason Chon, Tianran Liu, and Richard Futrell. 2023. [The linearity of the effect of surprisal on reading times across languages](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 15711–15721, Singapore. Association for Computational Linguistics.
- Rebecca Zhu. 2021. [Preschoolers’ acquisition of producer-product metonymy](#). *Cognitive Development*, 59:101075.
- George Kingsley Zipf. 1945. [The meaning-frequency relationship of words](#). *The Journal of general psychology*, 33:251–6.