

# Revisiting a Pain in the Neck: A Semantic Reasoning Benchmark for Language Models

Yang Liu<sup>1,2\*</sup>, Hongming Li<sup>1\*</sup>,  
Melissa Xiaohui Qin<sup>1</sup>, Qiankun Liu<sup>1</sup>, Chao Huang<sup>1†</sup>

<sup>1</sup>University of Science and Technology Beijing,

<sup>2</sup>State Key Laboratory of General Artificial Intelligence, BIGAI

liuyang@bigai.ai, hongmingli.lhm@gmail.com,

{qinxiaohui, liuqk3, chaohuang}@ustb.edu.cn

## Abstract

We present SEMANTICQA, an evaluation suite designed to assess language models (LMs) in semantic phrase processing tasks. The benchmark consolidates existing multiword expression (MWE) resources and reorganizes them into a unified testbed. It covers both general lexical phenomena, such as lexical collocations, and three fine-grained categories: idiomatic expressions, noun compounds, and verbal constructions. Through SEMANTICQA, we assess LMs of diverse architectures and scales in extraction, classification, and interpretation tasks, as well as sequential task compositions. We reveal substantial performance variation, particularly on tasks requiring semantic reasoning, highlighting differences in reasoning efficacy and semantic understanding of LMs, providing insights for pushing LMs with stronger comprehension on non-trivial semantic phrases. The evaluation harness and data of SEMANTICQA are available at <https://github.com/jacklanda/SemanticQA>.

## 1 Introduction

Semantic phrases (SP), also referred to as multiword expressions (MWE), are lexical combinations whose meanings or usages may not be fully derived from their individual components (Pasquer et al., 2020). They exhibit varying degrees of compositionality, idiomaticity, and fixedness (Sailer and Markantonatou, 2018; Ramisch, 2023). Despite extensive work in supervised and unsupervised paradigms, robust SP processing remains a fundamental challenge in NLP (Sag et al., 2002a; Constant et al., 2017a; Schwartz and Dagan, 2019; Ramisch et al., 2023a; Tanner and Hoffman, 2023).

Language models (LMs) are typically evaluated using benchmarks that emphasize mathematical reasoning (An et al., 2025; Balunovic et al., 2025),

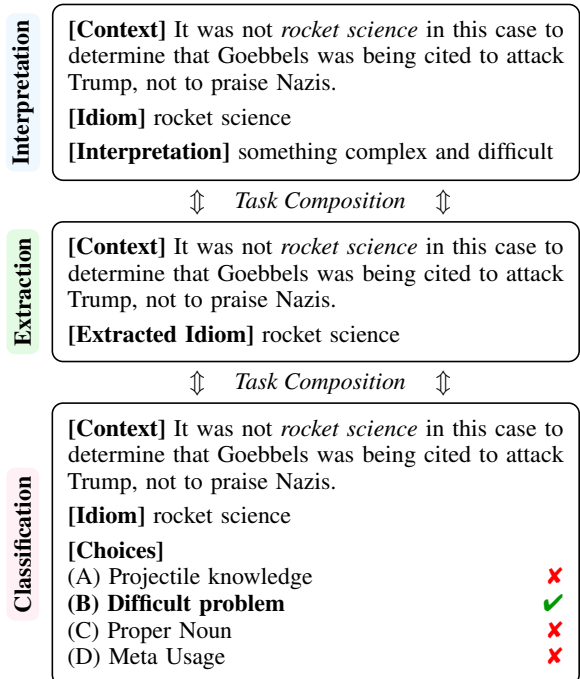


Figure 1: Atomic task exemplars of idiomatic expression in SEMANTICQA, grouped as task compositions.

code generation (Austin et al., 2021; Li et al., 2024), and logical reasoning (Li et al., 2025; Liu et al., 2026a). While these benchmarks effectively assess reasoning capacity and factual knowledge proficiency of LMs (Hu et al., 2024; Li et al., 2025; Yang et al., 2026), they largely overlook fine-grained semantic reasoning that operates over sub-sentential units. In particular, phrasal semantics, where meaning emerges from interactions between lexical constituents and context, remain under-explored and, when evaluated, are often assessed via isolated task formats that conflate multiple semantic operations (Liu et al., 2026b). As a result, it is difficult to determine whether strong performance reflects stable phrase-level semantic representations or task-specific heuristics. Therefore, recent work has called for diagnostic evaluations that disentangle semantic operations and investigate phrasal semantic

\*Equal contribution.

†Correspondence to: chaohuang@ustb.edu.cn

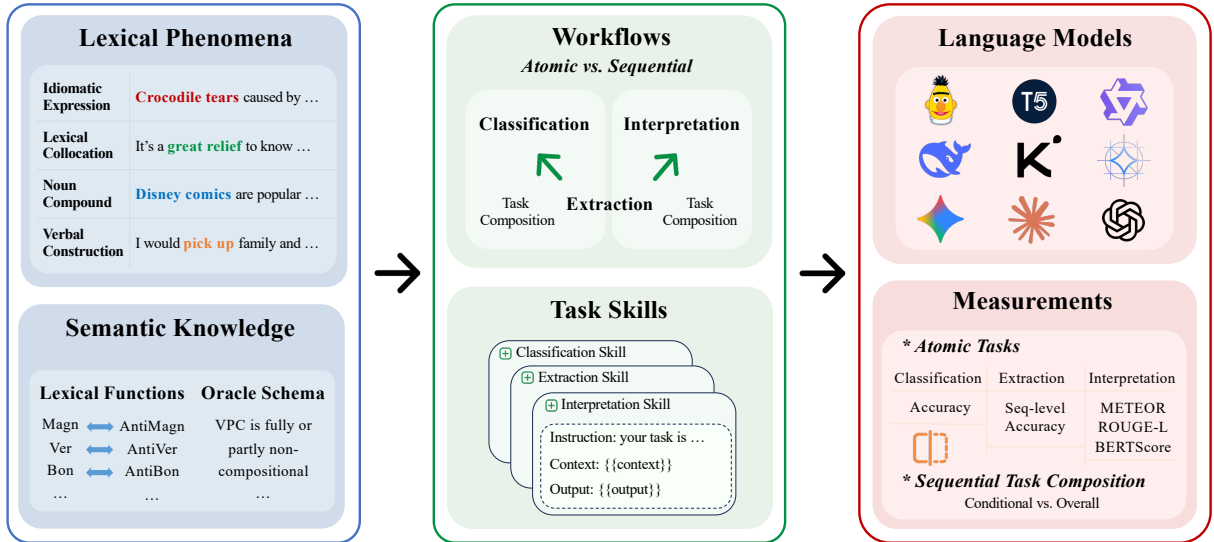


Figure 2: Overview of SEMANTICQA for benchmarking LMs on lexical phenomena.

behaviours beyond the understanding of superficial language (Miletić and Schulte im Walde, 2024).

We therefore ask: **How do language models behave when evaluated on phrasal semantics across distinct but structurally constrained task operations?** To answer this question, we introduce SEMANTICQA, an operation-aligned benchmark for semantic phrase processing. We adopt a deliberately operationalized view of semantic reasoning with respect to evaluation. Rather than requiring LMs to perform operations on the same instance, we examine whether phrasal semantic understanding generalizes across tasks that instantiate different operations. Specifically, we consider three atomic operations, including classification, extraction, and interpretation, which target the same underlying notion of phrasal sense while imposing structural constraints on LM generations. Under this formulation, semantic reasoning is assessed by the ability of LMs to exhibit compatible behavioral patterns across tasks, reflecting whether learnt phrasal semantics transfer across operations rather than overfitting to isolated task formats.

Under this definition, performance on a single task is insufficient. Instead, semantic reasoning is assessed via cross-operation consistency on SEMANTICQA, sensitivity to structural constraints, and robustness under compositional setups where applicable. Our contributions are threefold:

1. **Operation-aligned Semantic Evaluation.** SEMANTICQA does not introduce new semantic theories but evaluates phrasal competence through a set of controlled semantic opera-

tions with varying structural constraints. Its core contribution lies in **aligning existing SP tasks with the semantic operations they instantiate**, enabling systematic analyses of semantic behaviors across structural distinction yet related task families.

2. **Minimal and Controlled Design.** SEMANTICQA employs fixed prompt templates to reduce prompt-induced variance across LMs. By holding prompt structure constant while varying semantic operations, it supports fair comparison under shared conditions.
3. **Diagnostic Analyses of Cascade Sensitivity.** In explicitly designed sequential task setups, we show that strong LMs often fail to maintain semantic consistency across dependent operations, revealing phrase-level limitations that remain hidden in single-task evaluations.

## 2 Related Work

**Complex Reasoning.** Recent work evaluated LMs in a wide range of areas (Li et al., 2024; An et al., 2025; Balunovic et al., 2025). They focus on structured reasoning over explicit representations, such as compositional procedures in math or symbolic tasks. Although effective for formal reasoning, they overlook fine-grained semantic operations and omit their applications in context (Liu et al., 2024; An et al., 2025; Luong et al., 2025).

In contrast, semantic reasoning relies on the composition of phrasal meaning, contextual disambiguation, semantic-role inference, and paraphrase

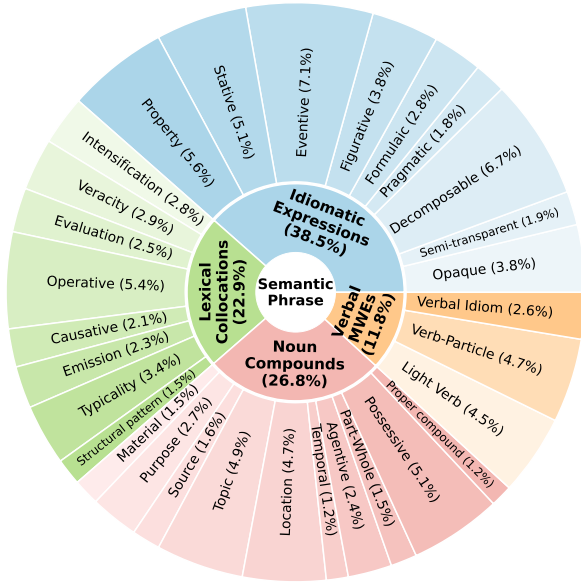


Figure 3: The coverage of coarse- and fine-grain semantic phrase categories of SEMANTICQA. The mapping of semantic phrase categories is listed in Table 14.

mapping. These aspects require the manipulation of latent semantic representations rather than symbolic rules. Prior work shows that even frontier LMs often depend on shallow heuristics, implying the need for assessment to examine semantic reasoning (Yu et al., 2024; Huang et al., 2025).

**Semantic Phrase Processing and Resources.** Semantic phrase processing has long been studied, with early work focusing on unsupervised phrase representations and compositional modeling (Vacareanu et al., 2020; Arase and Tsujii, 2020). Recent work examined idiom identification, contextual paraphrasing, and noun compound interpretation using transformers (Klubička et al., 2023; Wada et al., 2023). In parallel, a wide range of data resources was developed to evaluate phrasal semantics, covering idiomatic expressions (Tedeschi et al., 2022; Zhou et al., 2021), lexical collocations (Espinosa-Anke et al., 2019; Fisas et al., 2020; Espinosa-Anke et al., 2021), and verbal constructions (Savary et al., 2023; Ramisch et al., 2020).

Existing work typically isolates specific phrase types, task formats, and semantic phenomena, such as compositionality and idiomatic distinctions, without explicitly modeling the atomic semantic operations underlying phrase comprehension (Pham et al., 2023; Buijelaar and Pezzelle, 2023; Zeng et al., 2023). Consequently, evaluations are conducted in isolation, limiting cross-task and cross-phenomenon analyses and generalization. This

fragmentation motivates SEMANTICQA: a unified and operation-aligned benchmark for evaluating phrase-level semantic processing in LMs.

## 3 SEMANTICQA

### 3.1 Preliminaries

Semantic phrase exhibit diverse degrees of compositionality and idiomaticity. We consider four representative phrase types that capture popular sources of phrase variation. In addition, drawing on prior literature, we systematically consolidate and standardize the scope of SP considered in this work. We also employ LMs to label fine-grained categories of the phrases, with the resulting classification illustrated in Figure 3 (Nunberg et al., 1994; Sag et al., 2002b; Tratz and Hovy, 2010; Savary et al., 2017; Harish et al., 2021; Kolluru et al., 2022; Chakrabarty et al., 2022a; Mel’čuk, 2023).

**Lexical Collocations (LC).** LC forms a broad class of SPs with varying degrees of compositionality. They are characterized by conventionalized lexical relations between a *base* word and a *collocate* word, ranging from largely compositional combinations to idiom-style usages (Espinosa-Anke et al., 2021; Shvets and Wanner, 2022).

**Idiomatic Expressions (IE).** IE are prototypical non-compositional phrases whose meanings can not be derived from their constituent words (*e.g.*, *kick the bucket*). Processing such expressions requires LMs to recover conventionalized meanings beyond literal composition. (Zhou et al., 2022; Zeng and Bhat, 2022; Haviv et al., 2023).

**Noun Compounds (NC).** NC are often compositional, but their interpretation frequently depends on implicit semantic relations, contextual cues, or world knowledge (*e.g.*, *baby oil vs. olive oil*) (Kolluru et al., 2022; Coil and Shwartz, 2023).

**Verbal Constructions (VC).** VC or verbal multiword expressions (VMWE), including light-verb constructions (LVC), verb-particle constructions (VPC), and verbal idioms (VID), are typically semi-compositional (Tanner and Hoffman, 2023; Savary et al., 2023; Ramisch et al., 2023b). Their meanings arise from an interaction between literal composition and conventional usage.

### 3.2 Benchmark Construction

SEMANTICQA is built upon prior resources (Harish et al., 2021; Espinosa-Anke et al., 2022; Garcia

Task	Data Source	Input ( $\mathcal{I}$ )	Output ( $\mathcal{O}$ )	Metrics	# Test Size	Phrase Type
IE Detection	Harish et al. (2021)	$\mathcal{P} \oplus \mathcal{S} \oplus \mathcal{IE}$	Choice from <i>Options</i>	ACC	273	IDIOMACITY
IE Extraction	Tedeschi et al. (2022)	$\mathcal{P} \oplus \mathcal{S}$	Extracted $\mathcal{IE}$	$ACC_s$	447	IDIOMACITY
IE Interpretation	Zhou et al. (2021); Chakrabarty et al. (2022b)	$\mathcal{P} \oplus \mathcal{S} \oplus \mathcal{IE}$	Interpretation of $\mathcal{IE}$	METEOR, ROUGE-L, BERTSCORE	818	IDIOMACITY
LC Categorization	Espinosa-Anke et al. (2021)	$\mathcal{P} \oplus \mathcal{T} \oplus \mathcal{S}$	Choice from <i>Options</i>	ACC	305	COLLOCATION
LC Extraction	Fisas et al. (2020)	$\mathcal{P} \oplus \mathcal{T} \oplus \mathcal{S}$	Extracted $\mathcal{LC}$	$ACC_s$	305	COLLOCATION
LC Interpretation	Espinosa-Anke et al. (2019, 2021)	$\mathcal{P} \oplus \mathcal{S} \oplus \mathcal{LC}$	Interpretation of $\mathcal{LC}$	METEOR, ROUGE-L, BERTSCORE	305	COLLOCATION
NC Compositionality	Garcia et al. (2021)	$\mathcal{P} \oplus \mathcal{S} \oplus \mathcal{NC}$	Choice from <i>Options</i>	ACC	242	NOUN COMPOUND
NC Extraction	Garcia et al. (2021); Kolluru et al. (2022)	$\mathcal{P} \oplus \mathcal{S}$	Extracted $\mathcal{NC}$	$ACC_s$	720	NOUN COMPOUND
NC Interpretation	Coil and Shwartz (2023)	$\mathcal{P} \oplus \mathcal{S} \oplus \mathcal{NC}$	Interpretation of $\mathcal{NC}$	METEOR, ROUGE-L, BERTSCORE	110	NOUN COMPOUND
VMWE Extraction	Savary et al. (2023)	$\mathcal{P} \oplus \mathcal{S}$	Extracted $\mathcal{VC}$	$ACC_s$	475	VERBAL MWE

Table 1: A summary of the data statistics in SEMANTICQA.  $\mathcal{P}$  refers to the prompt template,  $\mathcal{S}$  denotes the sentence context,  $\mathcal{T}$  represents the semantic taxonomy narrative,  $\mathcal{IE}$  denotes idiomatic expressions,  $\mathcal{LC}$  denotes lexical collocations, and  $\mathcal{NC}$  denotes noun compounds.

Lexical Function	Example	Semantic Relation
Magn	Magn( <i>rain</i> ) = <i>heavy</i>	“intense”, “strong”
AntiMagn	AntiMagn( <i>accent</i> ) = <i>slight</i>	“little”, “weak”
Ver	Ver( <i>message</i> ) = <i>clear</i>	“real”, “genuine”
AntiVer	AntiVer( <i>accusation</i> ) = <i>groundless</i>	“non-genuine”
Bon	Bon( <i>bread</i> ) = <i>fresh</i>	“positive”
AntiBon	AntiBon( <i>advantage</i> ) = <i>undue</i>	“negative”
Son	Son( <i>alarm clock</i> ) = <i>ring(s)</i>	“sound”, “voice”
Oper1	Oper1( <i>advice</i> ) = <i>give</i>	“perform”

Table 2: Partial semantic relations involved in this paper, with their exemplars. More relations in lexical functions (LFs) can be referred to the Table 9.

et al., 2021; Savary et al., 2023), which vary in annotation protocols, difficulty distributions, and semantic granularity. Rather than enforcing uniform difficulty or annotation consistency across sources, SEMANTICQA is designed to reflect the variation and is not intended for absolute comparisons on phrase types. We focus on within-task trends, as well as relative changes induced by semantic operations and sequential compositions. Semantic reasoning is grounded in performance patterns that are stable across multiple tasks and datasets, rather than in absolute scores. All experiments use the datasets described in Table 1 and §B.

### 3.3 Task Definitions

We organize tasks by both phrase types and atomic task operations, where each operation targets a distinct aspect. This allows tasks operating on the same underlying phrase meaning to differ in their

output structure and constraints.

For IE, we include detection (IED), extraction (IEE), and interpretation (IEI) tasks. Detection is formulated as a multiple-choice classification task, extraction requires exact span identification, and interpretation evaluates contextualized paraphrase generation. All datasets are adapted from existing annotated resources (Harish et al., 2021; Tedeschi et al., 2022; Zhou et al., 2021), with overlapping instances deduplicated and reformatted to ensure consistency across operations.

For LC, we design categorization (LCC), extraction (LCE), and interpretation (LCI) tasks. Categorization requires predicting the semantic relation of a collocation under a lexical function taxonomy (cf. Table 2 and Appendix §A) (Mel’čuk, 2023). Extraction identifies both the *base* word and *collocate* word in context, while interpretation (cf. Appendix §F) evaluates paraphrasing conditioned on context. Datasets are balanced across semantic relation categories to support controlled multi-class evaluation (Espinosa-Anke et al., 2021; Fisas et al., 2020; Espinosa-Anke et al., 2022, 2021).

For NC, we include compositionality classification (NCC), extraction (NCE), and interpretation (NCI) tasks, which evaluate compositionality judgement, structural identification, and literal meaning reconstruction in a given context, respectively (Garcia et al., 2021; Kolluru et al., 2022; Coil and Shwartz, 2023; Hendrickx et al., 2013).

For VMWE, we include VMWE extraction task, which requires identifying a single verbal construc-

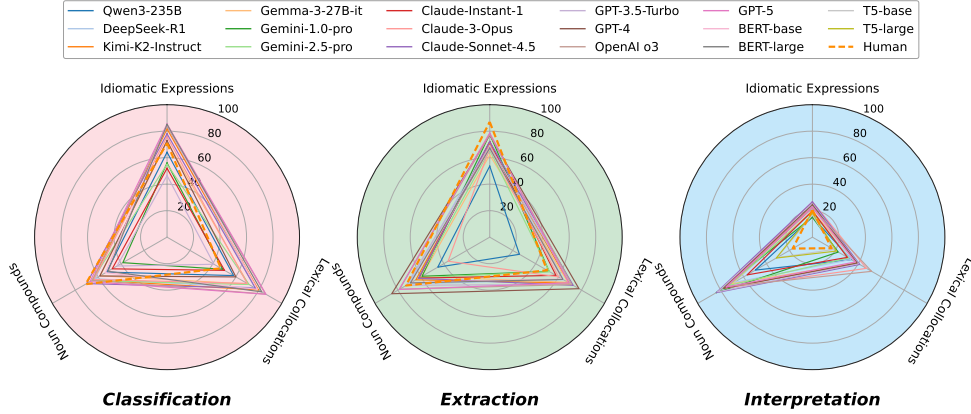


Figure 4: Overall the best performance (*i.e.*, capacity triangle  $\triangle$ ) of models on SEMANTICQA

tion in context, covering VPC (VPE), LVC (LVE), and VID (VIE) (Savary et al., 2023).

Finally, we formalize SP processing as a conditional generation problem under operation constraints. Given a prompt template  $\mathcal{P}$  (*cf.* Appendix §C) that specifies a target operation and an SP embedded in its context  $\mathcal{S}$ , a LM is required to generate an output  $\mathcal{O}$  that satisfies the instruction induced by  $\mathcal{P}$ . Concretely, the model input is constructed as  $\mathcal{I} := \mathcal{P} \oplus \mathcal{S}$ , where  $\oplus$  denotes a task-specific composition of instruction and contextualized phrase. The output  $\mathcal{O}$  varies according to the semantic operation being evaluated. For example, in extraction tasks,  $\mathcal{O}$  corresponds to the target phrase span identified from  $\mathcal{S}$  under the constraints specified by  $\mathcal{P}$ , whereas in classification or interpretation tasks,  $\mathcal{O}$  represents a semantic decision or reconstruction aligned with the given instruction.

For each task  $t$ , the configuration of the tuple  $(\mathcal{P}, \mathcal{S}, \mathcal{O})$  is instantiated according to a fixed template, as in Table 1. The dataset for task  $t$  is defined as  $\mathcal{D}^{(t)} := \{(p^{(t)}, s_i^{(t)}, o_i^{(t)})\}_{i=1}^N$ , where each example pairs a prompt, a contextualized SP, and a gold-standard output corresponding to the target semantic operation (*cf.* Figures 1 and 2).

### 3.4 Measurement

We adopt task-appropriate automatic metrics aligned with the output characteristics of each semantic operation. Classification tasks are evaluated using accuracy (ACC). Extraction tasks are evaluated using the accuracy of the exact match at the sequence-level ( $\text{ACC}_s$ ), which requires the exact recovery of the target phrase from the given context and avoids the inflation of the score from partial matches. Interpretation tasks are evaluated using METEOR (MTR) (Denkowski and Lavie, 2014) as

the primary metric, with ROUGE-L (R-L) (Lin, 2004) and BERTScore (B-S) (Zhang et al., 2020) reported for complementary analyses.

## 4 Results

### 4.1 Evaluation Setups

We evaluate a diverse set of LMs spanning different architectures, scales, and reasoning capabilities (*cf.* Appendix §E and Tables 8 and 11), including GPT-5 (Singh et al., 2025), OpenAI o3 (OpenAI, 2025), GPT-4 (Achiam et al., 2023), Claude-Sonnet-4.5 (Anthropic, 2025), Gemini-2.5-Pro (Comanici et al., 2025), Claude-3-Opus (Anthropic, 2024), DeepSeek-R1 (DeepSeek, 2025), Qwen3-235B (Yang et al., 2025), Gemma-3-27B-it (Kamath et al., 2025), and Kimi-K2-Instruct (Kimi, 2025), BERT-base/large (Devlin et al., 2019), and T5-base/large (Raffel et al., 2020), as summarized in Figure 4 and Table 12.

### 4.2 Benchmarking Results

**Overall Performance Patterns.** Table 3, Figures 5 and 4 show substantial variation across operations and phrase types (see Tables 12 and 10). Even within the same lexical phenomenon (*e.g.*, IE or LC), models behave differently in tasks, indicating that these operations impose distinct structural and semantic constraints. Interestingly, no model performs uniformly well across all setups, suggesting operation-specific strengths and weaknesses rather than a single transferable notion of phrase-level competence. Moreover, SEMANTICQA is neither saturated nor uniformly difficult: different tasks expose complementary failure modes, supporting its use as a diagnostic testbed rather than a leaderboard driven by aggregating scores.

MODEL	IDIOM			COLLOCATION			NOUN COMPOUND			VMWE		
	IED	IEE	IEI	LCC	LCE	LCI	NCC	NCE	NCI	VPE	LVE	VIE
METRIC (%)	ACC	ACC <sub>s</sub>	MTR	ACC	ACC <sub>s</sub>	MTR	ACC	ACC <sub>s</sub>	MTR	ACC <sub>s</sub>	ACC <sub>s</sub>	ACC <sub>s</sub>
HUMAN	71.0	<b>87.0</b>	20.5	47.0	50.0	16.7	<b>71.0</b>	73.0	17.2	<b>85.0</b>	<b>55.0</b>	<b>78.0</b>
DEEPSEEK-R1: <i>zero-shot</i>	71.1	69.4	<b>12.4</b>	66.6	<b>31.5</b>	<b>31.8</b>	60.2	51.3	<b>31.4</b>	76.8	26.7	50.5
↪ + <i>three-shot</i>	79.1	70.6	19.4	76.4	55.6	33.6	62.7	66.3	68.3	74.7	26.7	<b>59.1</b>
↪ + <i>five-shot</i>	84.3	72.3	19.2	76.1	64.3	32.9	60.6	70.7	68.7	81.6	35.8	57.1
KIMI-K2-INSTRUCT: <i>zero-shot</i>	68.5	63.1	13.9	68.5	34.4	33.7	60.6	45.4	65.4	<b>55.8</b>	28.9	46.7
↪ + <i>three-shot</i>	77.7	68.9	23.5	79.0	67.9	39.1	59.3	64.4	71.4	79.5	39.4	43.8
↪ + <i>five-shot</i>	81.7	69.6	21.7	79.7	69.2	36.9	64.7	63.6	76.7	81.1	<b>43.3</b>	46.7
GEMMA-3-27B-IT: <i>zero-shot</i>	<b>55.0</b>	<b>57.3</b>	13.5	<b>58.0</b>	38.4	35.0	58.3	<b>39.9</b>	43.8	66.8	19.4	38.1
↪ + <i>three-shot</i>	69.6	62.0	19.9	70.1	63.7	37.3	56.7	57.2	68.3	74.1	28.3	45.7
↪ + <i>five-shot</i>	72.1	61.6	19.2	70.8	68.2	38.7	56.2	59.2	70.5	70.5	35.0	52.4
CLAUDE-SONNET-4.5: <i>zero-shot</i>	72.5	68.5	17.0	67.5	40.1	34.8	<b>51.0</b>	45.1	77.2	69.8	<b>16.1</b>	<b>41.9</b>
↪ + <i>three-shot</i>	77.7	72.0	25.8	77.1	70.5	41.2	61.4	59.3	81.2	76.8	30.6	42.9
↪ + <i>five-shot</i>	78.0	72.0	<b>26.7</b>	76.1	<b>72.7</b>	40.8	<b>70.1</b>	62.1	<b>83.8</b>	<b>82.0</b>	37.2	47.6
OPENAI o3: <i>zero-shot</i>	57.1	65.1	12.6	72.1	37.7	35.9	65.2	62.9	45.7	67.9	25.6	51.4
↪ + <i>three-shot</i>	79.5	77.4	21.3	85.9	65.3	<b>41.6</b>	58.9	77.5	68.2	76.3	29.1	52.4
↪ + <i>five-shot</i>	83.5	74.7	21.9	83.6	71.5	35.9	63.5	78.6	74.5	77.3	36.9	50.0
GPT-5: <i>zero-shot</i>	82.8	67.6	13.9	75.4	36.7	33.7	66.8	64.3	57.3	74.2	28.9	56.2
↪ + <i>three-shot</i>	82.1	78.3	22.6	<b>86.2</b>	67.2	35.4	61.8	77.1	70.1	74.7	33.3	51.4
↪ + <i>five-shot</i>	<b>85.4</b>	<b>78.7</b>	22.5	84.3	68.9	37.4	67.2	<b>79.0</b>	75.3	74.7	38.3	50.5

Table 3: Major experimental results on SEMANTICQA. **Digits** highlight cases in which human scores are higher than those of all evaluated models, serving as a coarse reference. Light Pink indicates the human baseline. Light Green and Light Blue present open-source models and proprietary models, respectively. Green indicates the highest performance across all models (zero-shot and few-shot) within each task category, while red indicates the lowest.

Model	IEI		LCI		NCI	
	R-L	B-S	R-L	B-S	R-L	B-S
DEEPSEEK-R1	<b>14.7</b>	<b>85.1</b>	42.0	90.2	<b>37.6</b>	<b>91.3</b>
↪ 3-shot	25.2	88.1	44.9	91.6	73.0	96.3
↪ 5-shot	25.0	88.0	44.6	91.8	75.5	96.6
KIMI-K2-INST.	18.8	86.7	<b>40.2</b>	90.3	68.9	95.6
↪ 3-shot	<b>27.9</b>	88.4	<b>52.2</b>	92.8	77.8	96.7
↪ 5-shot	26.4	88.2	48.3	<b>97.2</b>	<b>83.7</b>	<b>97.2</b>
OPENAI o3	17.3	86.5	41.5	<b>89.8</b>	49.9	93.8
↪ 3-shot	26.2	88.5	51.2	92.6	71.2	96.0
↪ 5-shot	26.8	88.6	44.8	91.6	76.5	96.5
GPT-5	19.2	86.6	40.6	89.9	56.4	93.3
↪ 3-shot	27.5	<b>88.7</b>	46.9	92.3	70.9	96.4
↪ 5-shot	27.1	88.6	47.7	92.3	77.7	96.8

Table 4: Interpretation task results on IEI, LCI, and NCI. We report ROUGE-L (R-L) and BERTSCORE (B-S) scores. **Green** indicates the best performance and **red** indicates the worst performance within each column.

**Effect of In-Context Learning (ICL).** The impact of ICL varies by task type (*cf.* Tables 3 and 4). Interpretation tasks benefit most consistently from few-shot prompting. Across IEI, LCI, and NCI, three- or five-shot demonstrations yield clear gains in MTR. However, as shown in Table 4, complementary metrics reveal that these improvements primarily reflect exemplar-guided reconstruction rather than strict semantic grounding, as embedding-based similarity can be high even when lexical overlap remains limited. Few-shot ICL improves both R-L and B-S, but gains vary by

phrase type, reflecting the inherently open-ended nature of interpretation outputs.

For classification tasks, ICL exhibits hybrid effects. Models with weaker zero-shot performance often improve, whereas others plateau or regress, such as OpenAI o3 on LCC and NCC, indicating sensitivity to exemplar selection and task formulation. Extraction tasks are the most unstable under ICL. While demonstrations can substantially improve task performance when span structure is clearly illustrated, performance may degrade when test instances diverge from the demonstrated patterns. Overall, ICL is consistently beneficial for interpretation, variably effective for classification, and highly task-dependent for extraction.

### 4.3 Human Performance

We estimate human performance using annotations from three linguistics graduate students, each labeling 100 randomly sampled examples per task in SEMANTICQA, following a two-stage protocol inspired by SuperGLUE (Sarlin et al., 2020). Human scores are reported as a contextual reference to situate task difficulty, rather than an upper bound on performance (*cf.* Table 3). Differences between human and model results may arise from metric properties, task ambiguity, and response normalization effects, especially for interpretation tasks.

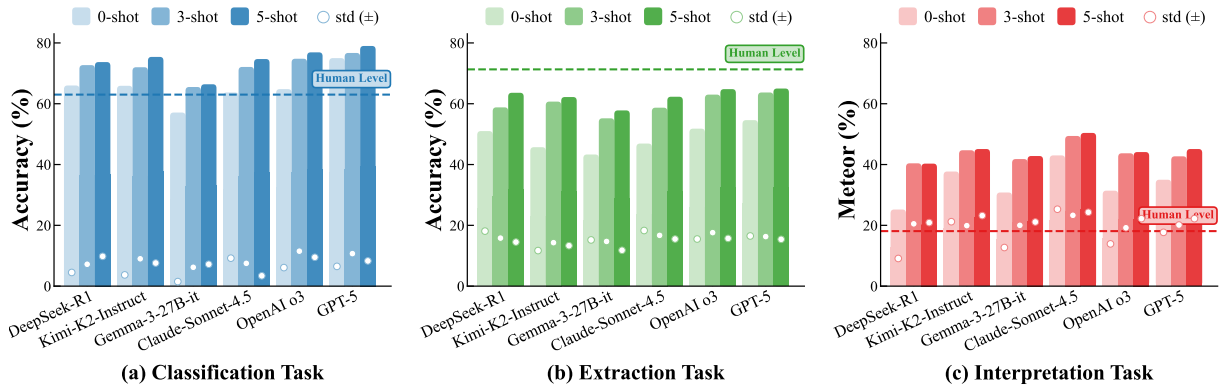


Figure 5: Grouped bars represent the mean performance of each model, while circular markers denote the population standard deviation (std), computed across tasks within the same category.

Accordingly, we avoid strong claims based on the absolute human vs. model comparisons. Instead, human performance is used to contextualize task difficulty and to illustrate evaluation challenges under varying output constraints of SEMANTICQA.

#### 4.4 Semantic Category Scaling with ICL

To examine how LMs encode semantic distinctions among lexical relations, we further investigate the LCC task under an increasing number of target categories. We construct a controlled scaling setup with varying the category size from 1 to 16 by log scale, and evaluate supervised models and four representative LMs under zero- and few-shot settings. Overall results are shown in Figures 6 and 20. Additionally, Figure 21 compares proprietary GPT-family models across all tasks, while Figure 22 presents the confusion matrix of GPT-5 on eight lexical function-based relations.

Across all settings, models consistently outperform random and majority baselines, indicating non-trivial semantic reasoning even without demonstrations. Accuracy decreases as the number of categories grows, but the degradation rate varies substantially across model families. Supervised baselines remain relatively stable, whereas frontier LMs exhibit sharper drops at larger scales. For example, DeepSeek-R1 decreases from 81.7% to 35.4% as category size increases, suggesting that in-context semantic reasoning alone does not fully substitute for explicit supervised signals when fine-grained relational distinctions are required.

#### 4.5 Sequential Task Compositions

To approximate realistic semantic phrase processing workflows, we evaluate **Sequential Task Compositions**, where models must perform multiple

*dependent* semantic operations in sequence, such as extraction followed by interpretation or categorization. Tables 5 and 6 report results for sequential interpretation and classification compositions (Ram et al., 2024; Alazraki et al., 2025).

For interpretation, conditional performance (Cond. MTR) on correctly extracted phrases is consistently higher than overall scores (Overall MTR) and shows only limited gains from few-shot prompting across both IE and LC settings. This gap indicates that accurate extraction remains a primary bottleneck for downstream interpretation, and that fluent semantic reconstruction does not reliably compensate for upstream structural errors. Compositional classification degrades more sharply as task complexity increases. While leading models perform well in four-class LC settings, accuracy drops substantially in eight- and sixteen-class scenarios, with similar trends observed for IE and NC. Few-shot prompting partially mitigates this degradation but does not remove the strong dependence on extraction quality. Overall, performance drops in compositional settings should be viewed as a diagnostic signal rather than evidence of complex error propagation. They indicate that current models struggle to robustly integrate intermediate semantic outputs, even when individual operations perform well in isolation. By separating atomic semantic operations from their compositions, SEMANTICQA exposes a persistent gap between performance on isolated atomic tasks and the stability of end-to-end semantic pipelines.

#### 4.6 VMWE Extraction with Oracle Schema

We analyze prompting strategies for VMWE extraction under zero-shot and few-shot ICL settings and introduce ORACLE SCHEMA, which augments

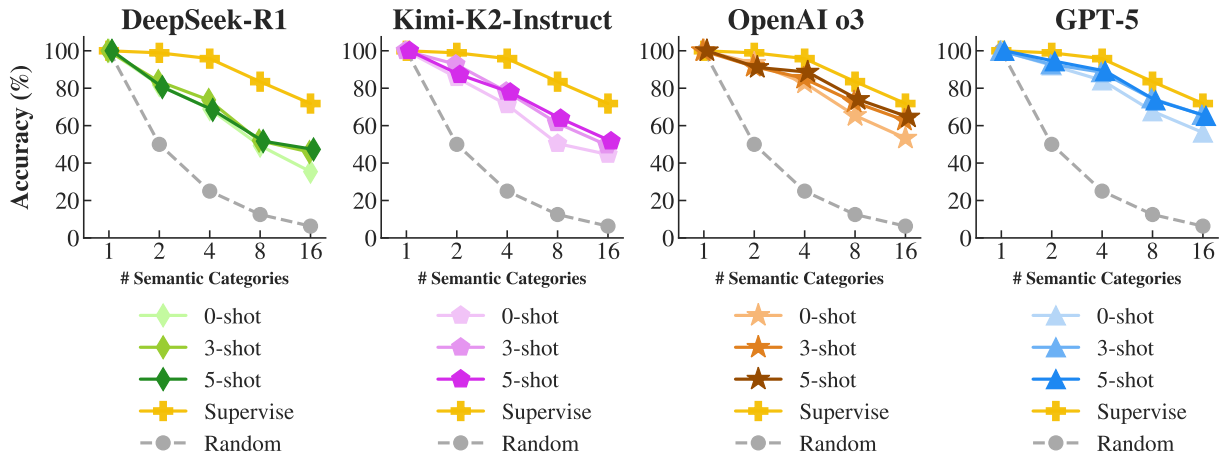


Figure 6: The ability of semantic relation categorization of  $\mathcal{LC}$  with different numbers of in-context exemplars and semantic category scale. The number  $n$  of classes is chosen from  $N := \{1, 2, 4, 8, 16\}$ . Each model is prompted with the  $k$ -shot settings, where  $k \in \{0, 3, 5\}$ , respectively. Accuracy scores are calculated by the mean values based on 30 examples sampled per class from the test split of Espinosa-Anke et al. (2021), partial categories ( $n \leq 8$ ) are run with three-class combinations in random selection, finally result in the mean value as the average.

Type	Model	# Shot	Ext. (ACC)	Cond. (MTR)	Overall (MTR)
LC	DEEPSEEK-R1	0-shot	27.9	<b>35.8</b>	10.0
		3-shot	34.4	38.8	13.4
		5-shot	33.8	<b>42.3</b>	14.3
	GPT-5	0-shot	26.2	37.6	<b>9.9</b>
		3-shot	39.7	40.1	15.9
		5-shot	41.3	41.8	<b>17.3</b>
IE	DEEPSEEK-R1	0-shot	51.3	<b>12.0</b>	<b>6.2</b>
		3-shot	57.3	13.0	7.4
		5-shot	57.0	13.4	7.6
	GPT-5	0-shot	48.3	<b>17.4</b>	8.4
		3-shot	55.7	17.2	9.6
		5-shot	59.3	17.1	<b>10.1</b>

Table 5: Performance comparison on sequential extraction-interpretation tasks. Ext. Acc denotes phrase extraction accuracy. Cond. MTR evaluates interpretation of correctly extracted phrases; Overall MTR reflects end-to-end performance. **Green** indicates the highest performance across all models within each task category, while **red** indicates the lowest.

task instructions with the target type and its definition (*cf.* Appendix §D). Table 7 shows that this strategy consistently improves performance across models. For example, DeepSeek-R1 increases from 51.6% to 64.1%, demonstrating that providing explicit semantic descriptions of the target expression substantially enhances VMWE extraction.

## 5 Discussion and Takeaways

Rather than restating performance trends, we distill what operation-aligned evaluation reveals about the assessment and modeling of semantics.

**Phrasal Semantics Requires Multi-dimensional Evaluations.** We show that phrase-level semantic competence can not be captured by any single

Type	Model	0-shot		3-shot		5-shot	
		Cond.	Overall	Cond.	Overall	Cond.	Overall
LC	DEEPSEEK-R1						
	↪ 4-class	73.4	36.4	74.9	44.2	80.5	44.4
	↪ 8-class	56.1	26.7	79.7	39.2	71.7	38.8
	↪ 16-class	<b>34.7</b>	<b>16.0</b>	51.0	25.6	54.5	27.7
	GPT-5						
	↪ 4-class	<b>91.3</b>	45.7	89.9	<b>58.1</b>	89.9	55.0
↪ 8-class	76.2	38.8	83.3	50.0	80.3	49.2	
↪ 16-class	63.4	33.1	69.4	43.4	73.4	44.8	
IE	DEEPSEEK-R1						
	↪ 4-class	63.8	46.5	65.0	46.9	<b>61.9</b>	<b>45.8</b>
	GPT-5						
↪ 4-class	79.3	65.9	77.7	<b>66.3</b>	<b>79.7</b>	65.9	
NC	DEEPSEEK-R1						
	↪ 4-class	<b>63.5</b>	<b>33.2</b>	<b>71.2</b>	36.9	71.1	37.8
	GPT-5						
↪ 4-class	68.8	36.5	64.7	37.3	66.9	<b>38.6</b>	

Table 6: Classification performance comparison. Cond.: accuracy given correct extraction; Overall: end-to-end accuracy. **Green** indicates the highest performance across all models (zero-shot and few-shot) within each task category, while **red** indicates the lowest.

task or metric. Interpretation, extraction, and categorization probe distinct aspects of semantic phrase processing and differ substantially in structural constraint. While extraction and categorization require explicit grounding in linguistic structure or semantic relations, interpretation operates in a weakly constrained output space. Consequently, performance on open-ended interpretation alone risks conflating fluent semantic generation with structurally grounded understanding.

**Metric Sensitivity Shapes Apparent Model Strengths.** The contrast between strong inter-

Model	# Shot	w/ ORACLE		w/o ORACLE	
		Acc( $\Delta$ ) $\uparrow$		Acc $\uparrow$	
DEEPSEEK-R1	0-shot	64.1 (+12.5)	vs.	51.6	
	3-shot	<b>72.3 (+8.9)</b>	vs.	63.4	
	5-shot	70.5 (+1.2)	vs.	69.3	
KIMI-K2-INST.	0-shot	53.3 (+9.1)	vs.	44.2	
	3-shot	67.6 (+1.1)	vs.	66.5	
	5-shot	<b>69.5 (+3.8)</b>	vs.	65.7	
OPENAI O3	0-shot	54.8 (+6.9)	vs.	47.9	
	3-shot	67.3 (+5.1)	vs.	62.2	
	5-shot	<b>70.7 (+3.6)</b>	vs.	67.1	
GPT-5	0-shot	59.6 (+7.6)	vs.	52.0	
	3-shot	66.8 (+5.1)	vs.	61.7	
	5-shot	<b>72.6 (+6.9)</b>	vs.	65.7	

Table 7: We report the accuracy of selected frontier LMs on the VMWE extraction task under different ICL setups, both with and without ORACLE SCHEMA.

pretation scores (B-S; cf. Tables 3, 4, and 12) and weaker extraction performance highlights how evaluation metrics shape perceived model capabilities. Flexible similarity-based metrics used for interpretation primarily reward paraphrasing ability and instruction-following behavior, whereas strict span-based evaluations expose brittleness in structural identification. As a result, high interpretation scores should be interpreted as evidence of improved exemplar-guided semantic reconstruction rather than conclusive semantic correctness. This discrepancy suggests that current evaluation practices may overestimate semantic robustness when structural constraints are not explicitly enforced.

**Workflow Robustness Remains Limited.** Sequential evaluations further reveal that semantic workflows are highly sensitive to upstream errors. Interpretation does not reliably compensate for failures in extraction or categorization; instead, structural errors propagate and often remain undetected under flexible metrics. This lack of robustness under error accumulation remains a key challenge for structured semantic evaluation settings and is largely obscured by atomic benchmarks.

## 6 Conclusion

We introduce SEMANTICQA, a benchmark for evaluating semantic phrase processing of LMs. We perform evaluations on a wide range of models with the introduced measurement, complemented by targeted human comparisons across ten tasks. The results show that, despite strong performance in general benchmarks, LMs continue to face substantial challenges in SEMANTICQA, revealing persistent limitations in semantic phrase understanding. Our

analyses further characterize model behavior across task types and highlights directions for future research on more robust and structurally grounded semantic processing.

## Limitations

This work has several limitations that suggest directions for future research. First, although SEMANTICQA covers four common phrase phenomena, it is restricted to English and does not capture the long tail of SP types, such as multiword named entities or complex function words (Constant et al., 2017b; Miletić and Schulte im Walde, 2024). Second, while multiple task formats are included, future benchmarks should incorporate more complex sequential task compositions and additional evaluation paradigms, such as semantic retrieval (Espinosa-Anke et al., 2021; Pham et al., 2023). Finally, although we evaluate many representative models, rapid progress in LM architectures calls for continual updates and broader coverage. We encourage future work to extend SEMANTICQA toward more comprehensive and multilingual resources (Espinosa-Anke et al., 2019).

## Ethical Considerations

This research uses publicly available datasets in accordance with their original licenses and does not include any private, sensitive, or personally identifiable information. The benchmark is intended solely for research and diagnostic purposes, and known limitations are explicitly documented to avoid overgeneralization. Computational resources were used responsibly, and potential risks related to data misuse and model evaluation were considered. Where human annotations were involved, annotators were recruited under fair labor practices and received appropriate compensation.

## Acknowledgement

This work was supported in part by the National Natural Science Foundation of China under Grant 62372039 and Grant 62002016, and by the Fundamental Research Funds for the Central Universities (FRF-BRA-25-012).

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman,

- Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Lisa Alazraki, Lihu Chen, Ana Brassard, Joe Stacey, Hossein A. Rahmani, and Marek Rei. 2025. **Agentcoma: A compositional benchmark mixing common-sense and mathematical reasoning in real-world scenarios**. *Preprint*, arXiv:2508.19988.
- Shengnan An, Xunliang Cai, Xuezhi Cao, Xiaoyu Li, Yehao Lin, Junlin Liu, Xinxuan Lv, Dan Ma, Xuanlin Wang, Ziwen Wang, and Shuang Zhou. 2025. **Amo-bench: Large language models still struggle in high school math competitions**. *Preprint*, arXiv:2510.26768.
- Anthropic. 2024. **The claude 3 model family: Opus, sonnet, haiku**. In *Anthropic Blog*.
- Anthropic. 2025. Anthropic. <https://www.anthropic.com/news/claude-sonnet-4-5>. September 30, 2025.
- Yuki Arase and Jun'ichi Tsujii. 2020. **Compositional phrase alignment and beyond**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1611–1623, Online. Association for Computational Linguistics.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. 2021. **Program synthesis with large language models**. *ArXiv*, abs/2108.07732.
- Mislav Balunovic, Jasper Dekoninck, Ivo Petrov, Nikola Jovanović, and Martin Vechev. 2025. **Matharena: Evaluating LLMs on uncontaminated math competitions**. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Ekaba Bisong. 2019. Google colab. *Building machine learning and deep learning models on google cloud platform: a comprehensive guide for beginners*, pages 59–64.
- Lars Buijelaar and Sandro Pezzelle. 2023. **A psycholinguistic analysis of BERT's representations of compounds**. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2230–2241, Dubrovnik, Croatia. Association for Computational Linguistics.
- Tuhin Chakrabarty, Yejin Choi, and Vered Shwartz. 2022a. **It's not rocket science: Interpreting figurative language in narratives**. *Transactions of the Association for Computational Linguistics*, 10:589–606.
- Tuhin Chakrabarty, Yejin Choi, and Vered Shwartz. 2022b. **It's not rocket science: Interpreting figurative language in narratives**. *Transactions of the Association for Computational Linguistics*, 10:589–606.
- I-Hsuan Chen, Yunfei Long, Qin Lu, and Chu-Ren Huang. 2017. **Leveraging eventive information for better metaphor detection and classification**. In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, pages 36–46, Vancouver, Canada. Association for Computational Linguistics.
- Albert Coil and Vered Shwartz. 2023. **From chocolate bunny to chocolate crocodile: Do language models understand noun compounds?** In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 2698–2710, Toronto, Canada. Association for Computational Linguistics.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. **Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities**. *arXiv preprint arXiv:2507.06261*.
- Mathieu Constant, Gülşen Eryiğit, Johanna Monti, Lonneke Van Der Plas, Carlos Ramisch, Michael Rosner, and Amalia Todirascu. 2017a. **Multiword expression processing: A survey**. *Computational Linguistics*, 43(4):837–892.
- Mathieu Constant, Gülşen Eryiğit, Johanna Monti, Lonneke van der Plas, Carlos Ramisch, Michael Rosner, and Amalia Todirascu. 2017b. **Survey: Multiword expression processing: A Survey**. *Computational Linguistics*, 43(4):837–892.
- DeepSeek. 2025. **Deepseek-r1 incentivizes reasoning in llms through reinforcement learning**. *Nature*, 645:633–638.
- Michael Denkowski and Alon Lavie. 2014. **Meteor universal: Language specific translation evaluation for any target language**. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 376–380, Baltimore, Maryland, USA. 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.
- Luis Espinosa-Anke, Joan Codina-Filba, and Leo Wanner. 2021. **Evaluating language models for the retrieval and categorization of lexical collocations**. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1406–1417, Online. Association for Computational Linguistics.
- Luis Espinosa-Anke, Steven Schockaert, and Leo Wanner. 2019. **Collocation classification with unsupervised relation vectors**. In *Proceedings of the 57th*

- Annual Meeting of the Association for Computational Linguistics*, pages 5765–5772, Florence, Italy. Association for Computational Linguistics.
- Luis Espinosa-Anke, Alexander Shvets, Alireza Mohammadshahi, James Henderson, and Leo Wanner. 2022. [Multilingual extraction and categorization of lexical collocations with graph-aware transformers](#). In *Proceedings of the 11th Joint Conference on Lexical and Computational Semantics*, pages 89–100, Seattle, Washington. Association for Computational Linguistics.
- Afsaneh Fazly, Paul Cook, and Suzanne Stevenson. 2009. Unsupervised type and token identification of idiomatic expressions. *Computational Linguistics*, 35(1):61–103.
- Beatriz Fisas, Luis Espinosa-Anke, Joan Codina-Filbá, and Leo Wanner. 2020. [CollFrEn: Rich bilingual English–French collocation resource](#). In *Proceedings of the Joint Workshop on Multiword Expressions and Electronic Lexicons*, pages 1–12, online. Association for Computational Linguistics.
- Marcos Garcia, Tiago Kramer Vieira, Carolina Scarton, Marco Idiart, and Aline Villavicencio. 2021. [Assessing the representations of idiomaticity in vector models with a noun compound dataset labeled at type and token levels](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2730–2741, Online. Association for Computational Linguistics.
- Alexander Gelbukh and 1 others. 2012. *Semantic analysis of verbal collocations with lexical functions*, volume 414. Springer.
- Tayyar Madabushi Harish, Gow-Smith Edward, Scarton Carolina, and Villavicencio Aline. 2021. [AStitchInLanguageModels: Dataset and methods for the exploration of idiomaticity in pre-trained language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3464–3477, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Adi Haviv, Ido Cohen, Jacob Gidron, Roei Schuster, Yoav Goldberg, and Mor Geva. 2023. [Understanding transformer memorization recall through idioms](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 248–264, Dubrovnik, Croatia. Association for Computational Linguistics.
- Iris Hendrickx, Zornitsa Kozareva, Preslav Nakov, Diarmuid Ó Séaghdha, Stan Szpakowicz, and Tony Veale. 2013. [SemEval-2013 task 4: Free paraphrases of noun compounds](#). In *Second Joint Conference on Lexical and Computational Semantics (\*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*, pages 138–143, Atlanta, Georgia, USA. Association for Computational Linguistics.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. In *International Conference on Learning Representations*.
- Xuming Hu, Junzhe Chen, Xiaochuan Li, Yufei Guo, Lijie Wen, Philip S. Yu, and Zhijiang Guo. 2024. [Towards understanding factual knowledge of large language models](#). In *The Twelfth International Conference on Learning Representations*.
- Sirui Huang, Yanggan Gu, Zhonghao Li, Xuming Hu, Li Qing, and Guandong Xu. 2025. [StructFact: Reasoning factual knowledge from structured data with large language models](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 7521–7552, Vienna, Austria. Association for Computational Linguistics.
- Gemma Team Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean-Bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, Gael Liu, and 191 others. 2025. [Gemma 3 technical report](#). *ArXiv*, abs/2503.19786.
- Kimi. 2025. [Kimi k2: Open agentic intelligence](#). *Preprint*, arXiv:2507.20534.
- Filip Klubička, Vasudevan Nedumpozhimana, and John Kelleher. 2023. [Idioms, probing and dangerous things: Towards structural probing for idiomaticity in vector space](#). In *Proceedings of the 19th Workshop on Multiword Expressions (MWE 2023)*, pages 45–57, Dubrovnik, Croatia. Association for Computational Linguistics.
- Olga Kolesnikova. 2020. Automatic detection of lexical functions in context. *Computación y sistemas*, 24(3):1337–1352.
- Keshav Kolluru, Gabriel Stanovsky, and Mausam. 2022. [“covid vaccine is against covid but Oxford vaccine is made at Oxford!” semantic interpretation of proper noun compounds](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10407–10420, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Jia Li, Ge Li, Xuanming Zhang, Yunfei Zhao, Yihong Dong, Zhi Jin, Binhua Li, Fei Huang, and Yongbin Li. 2024. [Evocodebench: An evolving code generation benchmark with domain-specific evaluations](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 57619–57641. Curran Associates, Inc.

- Jiaqi Li, Xinyi Dong, Yang Liu, Zhizhuo Yang, Quansen Wang, Xiaobo Wang, Song-Chun Zhu, Zixia Jia, and Zilong Zheng. 2025. **ReflectEvo: Improving meta introspection of small LLMs by learning self-reflection**. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 16948–16966, Vienna, Austria. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. **ROUGE: A package for automatic evaluation of summaries**. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Hongwei Liu, Zilong Zheng, Yuxuan Qiao, Haodong Duan, Zhiwei Fei, Fengzhe Zhou, Wenwei Zhang, Songyang Zhang, Dahua Lin, and Kai Chen. 2024. **MathBench: Evaluating the theory and application proficiency of LLMs with a hierarchical mathematics benchmark**. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 6884–6915, Bangkok, Thailand. Association for Computational Linguistics.
- Yang Liu, Jiaqi Li, and Zilong Zheng. 2026a. **Rulereasoner: Reinforced rule-based reasoning via domain-aware dynamic sampling**. In *The Fourteenth International Conference on Learning Representations*.
- Yang Liu, Jiaye Yang, Weikang Li, Jiahui Liang, Yang Li, and Lingyong Yan. 2026b. **LM-lexicon: Improving definition modeling via harmonizing semantic experts**. In *Proceedings of the 19th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1–22, Rabat, Morocco. Association for Computational Linguistics.
- Thang Luong, Dawsen Hwang, Hoang H Nguyen, Golnaz Ghiasi, Yuri Chervonyi, Insuk Seo, Junsu Kim, Garrett Bingham, Jonathan Lee, Swaroop Mishra, Alex Zhai, Huiyi Hu, Henryk Michalewski, Jimin Kim, Jeonghyun Ahn, Junhwi Bae, Xingyou Song, Trieu Hoang Trinh, Quoc V Le, and Junehyuk Jung. 2025. **Towards robust mathematical reasoning**. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 35406–35430, Suzhou, China. Association for Computational Linguistics.
- Igor A. Mel'čuk. 1998. Collocations and lexical functions. *Phraseology. Theory, analysis, and applications*, pages 23–53.
- Igor A. Mel'čuk. 2023. *General phraseology: Theory and practice*. John Benjamins.
- Filip Miletić and Sabine Schulte im Walde. 2024. **Semantics of multiword expressions in transformer-based models: A survey**. *Transactions of the Association for Computational Linguistics*, 12:593–612.
- Geoffrey Nunberg, Ivan A. Sag, and Thomas Wasow. 1994. **Idioms**. *Language*, 70(3):491–538.
- OpenAI. 2025. Openai. <https://openai.com/index/introducing-o3-and-o4-mini>. April 16, 2025.
- Caroline Pasquer, Agata Savary, Carlos Ramisch, and Jean-Yves Antoine. 2020. **Verbal multiword expression identification: Do we need a sledgehammer to crack a nut?** In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3333–3345, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, and 1 others. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.
- Thang Pham, Seunghyun Yoon, Trung Bui, and Anh Nguyen. 2023. **PiC: A phrase-in-context dataset for phrase understanding and semantic search**. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1–26, Dubrovnik, Croatia. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. **Language models are unsupervised multitask learners**. *OpenAI*. Accessed: 2024-11-15.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.
- Parikshit Ram, Tim Klinger, and Alexander G. Gray. 2024. **What makes models compositional? a theoretical view**. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI '24*.
- Carlos Ramisch. 2023. *Multiword expressions in computational linguistics*. Habilitation à diriger des recherches, Aix Marseille Université (AMU).
- Carlos Ramisch, Agata Savary, Bruno Guillaume, Jakub Waszczuk, Marie Candito, Ashwini Vaidya, Verginica Barbu Mititelu, Archana Bhatia, Uxoa Iñurieta, Voula Giouli, Tunga Güngör, Menghan Jiang, Timm Lichte, Chaya Liebeskind, Johanna Monti, Renata Ramisch, Sara Stymne, Abigail Walsh, and Hongzhi Xu. 2020. **Edition 1.2 of the PARSEME shared task on semi-supervised identification of verbal multiword expressions**. In *Proceedings of the Joint Workshop on Multiword Expressions and Electronic Lexicons*, pages 107–118, online. Association for Computational Linguistics.
- Carlos Ramisch, Abigail Walsh, Thomas Blanchard, and Shiva Taslimipoor. 2023a. **A survey of mwe identification experiments: The devil is in the details**. In *Proceedings of the 19th Workshop on Multiword Expressions (MWE 2023)*, pages 106–120.

- Carlos Ramisch, Abigail Walsh, Thomas Blanchard, and Shiva Taslimipour. 2023b. [A survey of MWE identification experiments: The devil is in the details](#). In *Proceedings of the 19th Workshop on Multiword Expressions (MWE 2023)*, pages 106–120, Dubrovnik, Croatia. Association for Computational Linguistics.
- María A Barrios Rodríguez. 2003. The domain of the lexical functions fact0, causfact0 and reall1. *learning*, page 64.
- Ivan A Sag, Timothy Baldwin, Francis Bond, Ann Copestake, and Dan Flickinger. 2002a. Multiword expressions: A pain in the neck for nlp. In *Computational Linguistics and Intelligent Text Processing: Third International Conference, CICLing 2002 Mexico City, Mexico, February 17–23, 2002 Proceedings 3*, pages 1–15. Springer.
- Ivan A. Sag, Timothy Baldwin, Francis Bond, Ann Copestake, and Dan Flickinger. 2002b. Multiword expressions: A pain in the neck for nlp. In *Computational Linguistics and Intelligent Text Processing*, pages 1–15, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Manfred Sailer and Stella Markantonatou. 2018. *Multiword expressions: Insights from a multi-lingual perspective*. Language Science Press.
- Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. 2020. Superglue: Learning feature matching with graph neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4938–4947.
- Agata Savary, Cherifa Ben Khelil, Carlos Ramisch, Voula Giouli, Verginica Barbu Mititelu, Najet Hadj Mohamed, Cvetana Krstev, Chaya Liebeskind, Hongzhi Xu, Sara Stymne, Tunga Güngör, Thomas Pickard, Bruno Guillaume, Eduard Bejček, Archana Bhatia, Marie Candito, Polona Gantar, Uxoa Iñurrieta, Albert Gatt, and 9 others. 2023. [PARSEME corpus release 1.3](#). In *Proceedings of the 19th Workshop on Multiword Expressions (MWE 2023)*, pages 24–35, Dubrovnik, Croatia. Association for Computational Linguistics.
- Agata Savary, Carlos Ramisch, Silvio Cordeiro, Federico Sangati, Veronika Vincze, Behrang Qasemizadeh, Marie Candito, Fabienne Cap, Voula Giouli, Ivelina Stoyanova, and Antoine Doucet. 2017. [The PARSEME shared task on automatic identification of verbal multiword expressions](#). In *Proceedings of the 13th Workshop on Multiword Expressions (MWE 2017)*, pages 31–47, Valencia, Spain. Association for Computational Linguistics.
- Alexander Shvets and Leo Wanner. 2022. [The relation dimension in the identification and classification of lexically restricted word co-occurrences in text corpora](#). *Mathematics*, 10(20).
- Vered Shwartz and Ido Dagan. 2019. [Still a pain in the neck: Evaluating text representations on lexical composition](#). *Transactions of the Association for Computational Linguistics*, 7:403–419.
- Aaditya Singh, Adam Fry, Adam Perelman, Adam Tart, Adi Ganesh, Ahmed El-Kishky, Aidan McLaughlin, Aiden Low, AJ Ostrow, Akhila Ananthram, and 1 others. 2025. Openai gpt-5 system card. *arXiv preprint arXiv:2601.03267*.
- Giorgos Spathas and Dimitris Michelioudakis. 2021. [States in the decomposition of verbal predicates](#). *Natural Language & Linguistic Theory*, 39(4):1253–1306.
- Joshua Tanner and Jacob Hoffman. 2023. [MWE as WSD: Solving multiword expression identification with word sense disambiguation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 181–193, Singapore. Association for Computational Linguistics.
- Simone Tedeschi, Federico Martelli, and Roberto Navigli. 2022. [ID10M: Idiom identification in 10 languages](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 2715–2726, Seattle, United States. Association for Computational Linguistics.
- Christine Thielen. 1999. Review of "turning a bilingual dictionary into a lexical-semantic database" by thierry fontenelle. max niemeyer verlag 1997. *Comput. Linguist.*, 25(3):447–449.
- Stephen Tratz and Eduard Hovy. 2010. [A taxonomy, dataset, and classifier for automatic noun compound interpretation](#). In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 678–687, Uppsala, Sweden. Association for Computational Linguistics.
- Robert Vacareanu, Marco A. Valenzuela-Escárcega, Rebecca Sharp, and Mihai Surdeanu. 2020. [An unsupervised method for learning representations of multi-word expressions for semantic classification](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3346–3356, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Takashi Wada, Yuji Matsumoto, Timothy Baldwin, and Jey Han Lau. 2023. [Unsupervised paraphrasing of multiword expressions](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 4732–4746, Toronto, Canada. Association for Computational Linguistics.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, and 3 others. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

- Alison Wray. 2002. *Formulaic language and the lexicon*, volume 20. Cambridge University Press Cambridge.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.
- Qianyu Yang, Yang Liu, Jiaqi Li, Jun Bai, Hao Chen, Kaiyuan Chen, Tiliang Duan, Jiayun Dong, Xiaobo Hu, Zixia Jia, and 1 others. 2026. \$onemillion-bench: How far are language agents from human experts? *arXiv preprint arXiv:2603.07980*.
- Fei Yu, Hongbo Zhang, Prayag Tiwari, and Benyou Wang. 2024. [Natural language reasoning, a survey](#). *ACM Comput. Surv.*, 56(12).
- Ziheng Zeng and Suma Bhat. 2022. [Getting BART to ride the idiomatic train: Learning to represent idiomatic expressions](#). *Transactions of the Association for Computational Linguistics*, 10:1120–1137.
- Ziheng Zeng, Kellen Cheng, Srihari Nanniyur, Jianing Zhou, and Suma Bhat. 2023. [IEKG: A common-sense knowledge graph for idiomatic expressions](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14243–14264, Singapore. Association for Computational Linguistics.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). In *International Conference on Learning Representations*.
- Jianing Zhou, Hongyu Gong, and Suma Bhat. 2021. [PIE: A parallel idiomatic expression corpus for idiomatic sentence generation and paraphrasing](#). In *Proceedings of the 17th Workshop on Multiword Expressions (MWE 2021)*, pages 33–48, Online. Association for Computational Linguistics.
- Jianing Zhou, Ziheng Zeng, Hongyu Gong, and Suma Bhat. 2022. Idiomatic expression paraphrasing without strong supervision. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 11774–11782.

## Appendix

<b>A Semantic Gloss for Lexical Functions</b>	<b>15</b>
<b>B Additional Details of Datasets</b>	<b>15</b>
B.1 Idiomacity Detection . . . . .	15
B.2 Idiom Extraction . . . . .	15
B.3 Idiom Interpretation . . . . .	16
B.4 Lexical Collocation Categorization	16
B.5 Lexical Collocation Extraction . .	16
B.6 Lexical Collocation Interpretation	16
B.7 Noun Compound Compositionality	16
B.8 Noun Compound Extraction . . .	17
B.9 Noun Compound Interpretation . .	17
B.10 VMWE Extraction . . . . .	17
<b>C Example Prompt</b>	<b>17</b>
<b>D Oracle Prompt</b>	<b>18</b>
<b>E Additional Experiment Details</b>	<b>18</b>
<b>F Annotation Guideline</b>	<b>19</b>

## A Semantic Gloss for Lexical Functions

In recent years, there has been an increasing interest in assigning lexical functions as labels to annotated MWE in the sense of the meaning-text theory (Mel’čuk, 2023). The lexical function is a multi-valued function, which  $f$  associates a lexical unit  $L$  with a set  $f(L)$  of lexical expressions.

As seen in Table 9, we constructed a collection of the representative lexical functions with their semantic glosses from the existing work. We compiled the prompts with the task descriptions.

## B Additional Details of Datasets

### B.1 Idiomacity Detection

In the initial dataset<sup>1</sup> proposed by (Harish et al., 2021), there exists three or four possible meanings (*i.e.*, interpretation) for each instance. For instances with only three interpretations, we add the option “None of the above” to keep consistency to the four-choices form. We deduplicate according to the unique (idiom, choice) pair for all instances. As a result, we collate 273 examples (*cf.* Table 1). Figure 7 shows an example of data.

[Context] There is also a covered pavilion. It is located next to <i>Silver Lining</i> Tire Recycling. The hours are 6:00 am to 10:00 pm, year round.	
[Choices]	
(A) grey lining	✗
(B) unexpected advantage	✗
<b>(C) Proper Noun</b>	✓
(D) Meta Usage	✗

Figure 7: A data example of idiomacity detection (IED).

### B.2 Idiom Extraction

The original dataset<sup>2</sup> consists of instances with or without idiom  $\mathcal{IE}$ . Since the inference-only experiments comprise most of our work, we filter out all the examples without the  $\mathcal{IE}$  existing to increase the coverage diversity of idioms; then, we deduplicate according to the unique item of the occurred  $\mathcal{IE}$ . The final prepared test set consists of 447 examples with a unique item of  $\mathcal{IE}$  existing in each. Figure 8 shows an example of data.

<sup>1</sup><https://github.com/H-TayyarMadabushi/AStitchInLanguageModels>

<sup>2</sup><https://github.com/Babelscape/ID10M>

[Context] In the screenplay by Lorenzo Semple Jr. , and David Rayfiel , Turner very early on stumbles upon the existence of a kind of super - C.I.A. within the C.I.A. , after which his life is *not worth a plug nickel* .


[Idiom] “not worth a plug nickel” 

Figure 8: A data example of idiom extraction (IEE).

### B.3 Idiom Interpretation

We collected 916 instances in total from the PIE (Zhou et al., 2021)<sup>3</sup> and (Chakrabarty et al., 2022b)<sup>4</sup>, after deduplication by occurred items of idiom  $\mathcal{IE}$ . Figure 9 shows an example of data.

[Context] The remission at this stage of having cancer was truly the *turning point* of her life .

[Idiom] “turning point”


[Interpretation] “the time of significant change (mostly positive) in situation” 

Figure 9: A data example of idiom interpretation (IEI).

### B.4 Lexical Collocation Categorization

We collect the collocation data with the annotated labels from the expanded LEXFUNC<sup>5</sup> (Espinosa-Anke et al., 2021). We inherited the training and validation sets of the initial data and sampled 50 examples per semantic category from the test set in classification concerning the computation efficiency. Figure 10 shows an example of data.

[Context] In genoa, the *violent storm* knocked down power lines, blacking out the homes of 5,000 residents.


[Category] **Magn** (strong semantic). 

Figure 10: A data example is the lexical collocation categorization (LCC) by semantic relations. Note that the taxonomy included in the prompt is omitted here.

<sup>3</sup>[https://github.com/zhjjn/MWE\\_PIE](https://github.com/zhjjn/MWE_PIE)

<sup>4</sup><https://github.com/tuhinjbcse/FigurativeNarrativeBenchmark>

<sup>5</sup><https://github.com/luisespinoaanke/lexicalcollocations>

### B.5 Lexical Collocation Extraction

The initial dataset is collected from (Fisas et al., 2020)<sup>6</sup>. We select the English part of the data and perform deduplication to filter out overlap collocations. We downsample 50 instances randomly for each semantic category to form our test set and reuse the training and validation sets of the original data. Figure 11 shows an example of data. We conduct  $\mathcal{LC}$  extraction but not identification task, and not query models to distinguish the base and the collocate to simplify the task in this work.

[Context] He still gets up the moment the *alarm clock rings* .

[Semantic relation] Strong or intense degree in the lexical semantic relation.


[Collocation] “alarm clock rings” 

Figure 11: A data example of lexical collocation extraction (LCE).

### B.6 Lexical Collocation Interpretation

The data<sup>7</sup> we used is proposed in (Espinosa-Anke et al., 2021). We perform random sampling from the original data and get the 400 examples (50 per class) as our test set. We manually annotated and revised the test examples, and get the Cohen’s kappa coefficient  $\kappa = 0.718$ , to confirm the quality. An example of data is shown in Figure 12.

[Context] Through robert bennett, his lawyer, the president continued friday to call mrs. jones’ *baseless accusation*.

[Collocation] “baseless accusation”


[Interpretation] “Groundless claim made without substantiation” 

Figure 12: A data example of lexical collocation interpretation (LCI).

### B.7 Noun Compound Compositionality

The annotated noun compound data is collected from the NCTTI<sup>8</sup> (Garcia et al., 2021). After data

<sup>6</sup><https://github.com/TaInUPF/CollFrEn>

<sup>7</sup><https://github.com/luisespinoaanke/lexicalcollocations>

<sup>8</sup><https://github.com/marcospln/nctti>

processing, we filtered out the compound without reference context, collated 237 examples, and split them into training, validation, and test sets. Figure 13 shows an example of data.

[Context] <i>Fair play</i> incorporates the concepts of friendship, respect for others and always playing in the right spirit.	
[Noun compound] “Fair play”	
[Choices]	
(A) Compositional	✗
<b>(B) Partly compositional</b>	✓
(C) None of the above	✗
(D) Non-compositional	✗

Figure 13: A data example of noun compound compositionality (NCC).

## B.8 Noun Compound Extraction

As our beginning, we sampled the test set from the PRONCI<sup>9</sup> (Kolluru et al., 2022). We used the training and validation sets to leverage the compositional part of noun compounds in the original dataset. We randomly sampled from the test set to form the new test set with 720 examples. We demonstrate a data example in Figure 14.

[Context] The rhombus shape of the patches arose by adaptation to the <i>Paris fashion</i> of the 17th century by Biancolelli.	
[Noun compound] “Paris fashion”	↖

Figure 14: A data example of noun compound extraction (NCE).

## B.9 Noun Compound Interpretation

We leverage the initial training, validation, and test data splits from (Coil and Shwartz, 2023)<sup>10</sup>. To provide a context for each noun compound, we use ChatGPT to generate a reference sentence. To verify the quality of synthetic data, we performed a manual inspection, which resulted in  $acc > 98\%$ . A data example is shown in the figure 15.

<sup>9</sup><https://github.com/dair-iitd/pronci>

<sup>10</sup><https://github.com/jordancoil/noun-compound-interpretation>

[Context] She used a straightedge to draw a <i>ruler line</i> across the paper, ensuring her graph was perfectly aligned.	
[Noun compound] “ruler line”	
[Interpretation] “line drawn with a ruler” “	

Figure 15: A data example of noun compound interpretation (NCI).

## B.10 VMWE Extraction

We used the English corpus of PARSEME v1.3<sup>11</sup> (Savary et al., 2023), the existing largest annotated corpora of VMWE. The initial data is used to conduct extraction instead of identification tasks. Figure 16 shows an example of the data.

[Context] Harry tore back across the room as the landing light <i>clicked on</i> .	
[VMWE] “clicked on”	↖

Figure 16: A data example of VMWE Extraction.

## C Example Prompt

We manually create a unified prompt template for all tasks that can be adapted to each task with specific filling arguments. The prompt format is shown in the Figure 17. The detailed prompt for each task can be accessed in our code base<sup>12</sup>.

<sup>11</sup>[https://gitlab.com/parseme/parseme\\_corpus\\_en](https://gitlab.com/parseme/parseme_corpus_en)

<sup>12</sup><https://github.com/lexbench/LexBench/tree/main/lexbench/prompts>

Unified Prompt Template	
Assume that you are a linguist who researches <i>{{semantic phrases}}</i> .	
You will be given a sentence that contains only an item of <i>{{semantic phrase}}</i> .	
Your task is to ...	
Please make sure you read and understand these instructions carefully.	
Few-shot Examples:	
Phrase: <i>{{an example of the phrase}}</i>	
Context: <i>{{a context of the example}}</i>	
Output: <i>{{an output of the example}}</i>	
...	
Phrase: <i>{{phrase}}</i>	
Context: <i>{{context}}</i>	
Output:	

Figure 17: Unified prompt template used in the work.

## D Oracle Prompt

Oracle Prompt Template	
Assume that you are a linguist who conducts research on <i>{{verbal multiword expressions (VMwEs)}}</i> .	
You will be given a context that includes only one <i>{{verbal multiword expression}}</i> .	
Your task is to ...	
Please make sure you read and understand these instructions carefully.	
Few-shot Examples:	
Context: <i>{{a context of the example}}</i>	
Output: <i>{{an output of the example}}</i>	
...	
VMwE Definition: <i>{{Definition of VMwE}}</i>	
VMwE Definition Example: Definition of VMwE: <i>{{"Verb-particle construction (VPC) is sometimes called phrasal or phrasal-prepositional verb. The meaning of the VPC is fully or partly non-compositional."}}</i>	
Context: <i>{{context}}</i>	
Output:	

Figure 18: Oracle prompt template used in the work.

## E Additional Experiment Details

For models accessed via API endpoints, the evaluation probes both zero-shot and few-shot (three- and five-shot) performance. Throughout all experiments, we set the sampling temperature to  $\tau = 0$  and employ top-p decoding (Holtzman

Model	# Params	Arch.	Creator	Public	Post Training
BERT base <sup>†</sup>	110M	Enc.	Google	✓	FT
BERT large <sup>†</sup>	340M	Enc.	Google	✓	FT
T5 base <sup>†</sup>	220M	Enc.+Dec.	Google	✓	FT
T5 large <sup>†</sup>	770M	Enc.+Dec.	Google	✓	FT
Qwen3-235B <sup>‡</sup>	235B	Dec.(MoE)	Qwen Team	✓	SFT
DeepSeek-R1 <sup>‡</sup>	685B	Dec.(MoE)	DeepSeek-AI	✓	SFT + RL
Kimi-K2-Instruct <sup>‡</sup>	1T	Dec.(MoE)	Kimi Team	✓	SFT
Gemma-3-27B-it <sup>‡</sup>	27B	Dec.	Gemma Team	✓	SFT
Gemini-1.0-pro <sup>‡</sup>	*	*	Google	✗	SFT + RL
Gemini-2.5-pro <sup>‡</sup>	*	*	Google	✗	SFT + RL
Claude-Instant-1 <sup>‡</sup>	*	*	Anthropic	✗	SFT + RL
Claude-3-Opus <sup>‡</sup>	*	*	Anthropic	✗	SFT + RL
Claude-Sonnet-4.5 <sup>‡</sup>	*	*	Anthropic	✗	SFT + RL
GPT-3.5-Turbo <sup>‡</sup>	*	*	OpenAI	✗	SFT + RL
GPT-4 <sup>‡</sup>	*	*	OpenAI	✗	SFT + RL
OpenAI o3 <sup>‡</sup>	*	*	OpenAI	✗	SFT + RL
GPT-5 <sup>‡</sup>	*	*	OpenAI	✗	SFT + RL

Table 8: A list of LMs tested in this paper: “Public” indicates whether the model weights are open. In detail, Light Pink text delineates the supervised fine-tuned models. Light Green and Light Blue parts present open-source models and proprietary models, respectively. “Post Training” indicates whether the model is trained further in some ways after pre-training. <sup>†</sup>We perform trivial full-set fine-tuning for the models. <sup>‡</sup>We use the official API for the model inference.

et al., 2019) with  $p = 1.0$ . Inference is accelerated and deployed using vLLM (Kwon et al., 2023). For non-API-based models, we apply the following configuration. For the sequence classification tasks such as LCC, we employ bert-base/large-uncased as our tuning initiation. Similarly, we construct primary baselines for extraction tasks that leverage the B-I-O scheme to conduct sequence labeling. The training is run with an NVIDIA A100-40GB on Google Colab (Bisong, 2019). For interpretation tasks, we use t5-base/large model to conduct vanilla fine-tuning. Additionally, We train all models for a specific number of epochs shown in Table 11 and perform early stopping over the validation set. Model checkpoints used in our experiment are implemented by PyTorch (Paszke et al., 2019), and Hugging Face Transformers (Wolf et al., 2020). The input format of the prompt and the few-shot demonstration settings we used during the experiment are shown in Figure 17. Since each model has different generation styles, we conduct a pre-run before each test. Then, we develop ad hoc heuristics based on the response generated by models to parse predictions accurately. The perplexity computing in the interpretation tasks is to feed the phrase and its interpretation into the template “*The meaning of phrase {{phrase}} in context is {{interpretation}}*”, and then we compute the token-level perplexity by GPT-2-XL (Radford et al., 2019).

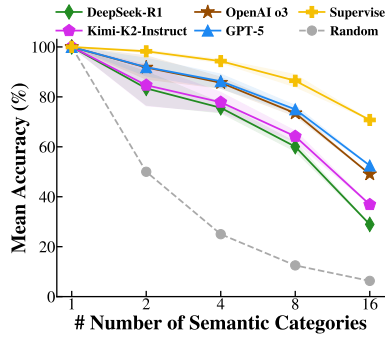


Figure (20) Each model is run with zero-shot prompting in the semantic relation classification with category scaling. Mean accuracy (%) of different models are average over runs in three sampled sets. For comparative reasons, we also plotted the level of random baseline.

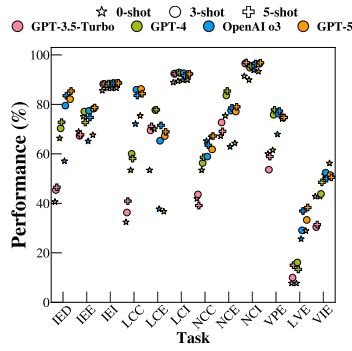


Figure (21) Model performance (GPT-3.5-Turbo, GPT-4, OpenAI o3, GPT-5) across all twelve tasks. Note that the y-axis denotes task-specific metrics, and thus absolute values should not be compared across different tasks.

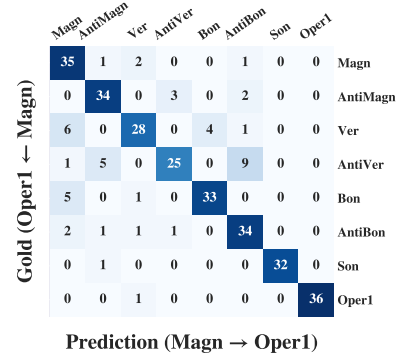


Figure (22) Confusion matrix for the best-performing model with ICL (GPT-5 in 5-shot setting) in categorizing eight semantic relations described by lexical functions (*cf.* Table 2). The x-axis denotes the prediction results, and the y-axis represents the gold standards.

## F Annotation Guideline

We established the following criteria for compiling the dataset of collocation interpretation (§3.3).

1. **Objective:** Interpret each lexical collocation in five distinct narratives for comprehensive understanding according to the given context.
2. **Dataset Overview:** Contains context and collocations paired with base and collocate.
3. **Annotation Format:** Include collocation, five narratives (N1-N5), and rationale.
4. **Consistency and Accuracy:** Maintain consistent and accurate interpretations across the five narratives in the same semantic meaning.

Lexical Function	Semantic Gloss	Complete Description
Magn (Mel'čuk, 1998)	Intense, strong degree, an intensifier of semantic relation for base lexeme.	Intensify the base lexeme to a high level, strengthening its semantic relation with the associated concept via the collocate lexeme.
AntiMagn (Mel'čuk, 1998)	Slight and weak degree, a de-intensifier	Weaken meaning intensity, diminishing the semantic relationship between the base lexeme and its associated concept.
Ver (Gelbukh et al., 2012)	Lat. verus, real, genuine	"As it should be", "Meet the intended requirements of <i>K</i> ".
AntiVer (Mel'čuk, 1998)	Non-genuine	Characterize something as non-genuine, not authentic, not in its intended or proper state, and not meeting the required standards or expectations.
Bon (Espinosa-Anke et al., 2021)	Positive	Something is good or in a positive situation.
AntiBon (Espinosa-Anke et al., 2021)	Negative	Something is bad or in a negative situation.
IncepPredPlus (Thielen, 1999)	Start to increase.	Denote initiating a process or action that leads to an increase or enhancement of something.
FinFunc0 (Kolesnikova, 2020)	End.existence	The value means "the <i>K</i> of FinFunc0 ceases to be experienced".
Fact0 (Mel'čuk, 1998)	Lat. factum, fact. To fulfil the requirement of <i>K</i> , and the argument of this function fulfills its own requirement.	Fulfill the base requirement, do something with the base, and do what you are supposed to do with the base.
CausFunc0 (Gelbukh et al., 2012)	The agent does something so that the event denoted by the noun occurs	Do something so that <i>K</i> begins occurring.
Caus1Func0 (Espinosa-Anke et al., 2021)	Cause the existence. 1st argument.	Bring about something's presence or creation, with the first argument indicating the responsible agent or entity.
CausFact0 (Rodríguez, 2003)	To cause something to function according to its destination.	Denote causing something to function according to its intended purpose or destination.
CausPredMinus (Thielen, 1999)	Cause to decrease.	Describe the act of causing a decrease or reduction in something.
CausFunc1 (Gelbukh et al., 2012)	The non-agentive participant does something such that the event denoted by the noun occurs.	A person/object, different from the agent of <i>K</i> , does something so that <i>K</i> occurs and has an effect on the agent of <i>K</i> .
LiquFunc0 (Espinosa-Anke et al., 2021)	Cause termination of the existence	Cause termination of the existence.
Son (Kolesnikova, 2020)	Lat. <i>sonare</i> : sound.	The <i>K</i> is usually a noun, and the value means "emit a characteristic sound".
Oper1 (Kolesnikova, 2020)	Lat. <i>operari</i> : perform, do, act something. The subject is as the 1st argument.	Represent a light verb linking the event's first participant (subject) with the event's name (direct object).
Oper2 (Espinosa-Anke et al., 2021)	Lat. <i>operari</i> : perform, do, act something. The subject is as the 2nd argument.	Represent a light verb linking the event's first participant (subject) with the event's name (indirect object).
IncepOper1 (Gelbukh et al., 2012)	Incep is from Lat. <i>incipere</i> : begin. Begin to do, perform, experience, carry out <i>K</i> .	Signify the start of an action or event, linking the event's subject with its name using a light verb.
FinOper1 (Kolesnikova, 2020)	Fin is from Lat. <i>finire</i> : cease.	Terminate doing something.
Real1 (Rodríguez, 2003)	Fulfill a requirement imposed by the noun or performing an action typical for the noun.	To fulfill the requirement of <i>K</i> , to act according to <i>K</i> .
Real2 (Kolesnikova, 2020)	Acting as expected. Something be realized as expected	<i>K</i> that is normally expected of the second participant
AntiReal2 (Kolesnikova, 2020)	Not acting as expected. Something not be realized as expected.	The <i>V</i> is the negation of an internal element of the argument of this function.

Table 9: All lexical functions with their semantic gloss in this paper. The column "semantic gloss" provides the definition for each LF, and we use a sentence to describe the complete meaning of LF in column "Complete Description". *K* denotes the keyword/base word of a LF, and *V* denotes the value/collocate word of a LF.

VMWE	BERT-base			BERT-large			# Support
	P	R	F1	P	R	F1	
IAV	60.7 <sub>5.6</sub>	38.0 <sub>4.3</sub>	46.5 <sub>3.3</sub>	46.5 <sub>3.6</sub>	38.9 <sub>5.6</sub>	42.3 <sub>4.8</sub>	36.0
LVC.cause	46.4 <sub>12.2</sub>	18.4 <sub>4.0</sub>	26.2 <sub>5.6</sub>	26.4 <sub>12.4</sub>	20.7 <sub>9.1</sub>	23.2 <sub>10.5</sub>	29.0
LVC.full	52.1 <sub>4.4</sub>	61.1 <sub>2.0</sub>	56.1 <sub>2.2</sub>	55.2 <sub>2.5</sub>	56.8 <sub>8.4</sub>	55.9 <sub>5.4</sub>	172.0
MVC	95.9 <sub>4.0</sub>	80.5 <sub>2.0</sub>	87.5 <sub>2.8</sub>	100.0 <sub>0.0</sub>	80.5 <sub>2.0</sub>	89.2 <sub>1.2</sub>	29.0
VID	52.4 <sub>5.3</sub>	36.1 <sub>0.9</sub>	42.7 <sub>1.8</sub>	63.8 <sub>5.0</sub>	36.1 <sub>1.9</sub>	46.1 <sub>2.0</sub>	108.0
VPC.full	64.3 <sub>3.3</sub>	78.4 <sub>1.6</sub>	70.6 <sub>1.5</sub>	64.4 <sub>0.1</sub>	79.4 <sub>0.5</sub>	71.1 <sub>0.2</sub>	194.0
VPC.semi	55.9 <sub>38.7</sub>	8.9 <sub>6.9</sub>	12.9 <sub>7.3</sub>	38.8 <sub>6.4</sub>	35.6 <sub>3.9</sub>	37.1 <sub>5.0</sub>	30.0
Micro Avg.	63.2 <sub>1.9</sub>	61.2 <sub>0.6</sub>	62.3 <sub>1.3</sub>	64.2 <sub>0.5</sub>	62.4 <sub>2.7</sub>	63.3 <sub>1.6</sub>	85.4

Table 10: We report the full results of VMWE extraction reproduced on MTLB-STRUCT. The performance of all categories are defined in the corpora PARSEME 1.3. The corresponding standard deviation is calculated by the results of three runnings with the selected seeds {21, 42, 84}.

Computing Infrastructure			
1 × A100 40GB GPU (Google Colab)			
Hyperparameter	Assignment	Hyperparameter	Assignment
architecture	BERT-{base, large}	architecture	T5-{base, large}
tokens per sample	150	tokens per sample	128
batch size	4,800	batch size	2,048
number of workers	8	number of workers	4
learning rate	$3e^{-5}$	learning rate	$5e^{-5}$
number of epochs	10	number of epochs	5
save interval (epoch)	1	save interval (epoch)	1
validation interval (epoch)	1	validation interval (epoch)	1
ratio of warmup steps	3%	ratio of warmup steps	3%
learning rate scheduler	Polynomial decay	learning rate scheduler	Cosine decay
learning rate optimizer	Adam	learning rate optimizer	Adam
Adam beta weights	(0.9, 0.99)	Adam beta weights	(0.9, 0.99)
Adam epsilon	$1e^{-6}$	Adam epsilon	$1e^{-6}$
weight decay	0	weight decay	0
random seed	21, 42, 84	random seeds	21, 42, 84

Table 11: Hyperparameters for finetuning BERT-Taggers and T5 Generators.

MODEL	IDIOM			COLLOCATION			NOUN COMPOUND			VMWE		
	IED	IEE	IEI	LCC	LCE	LCI	NCC	NCE	NCI	VPE	LVE	VIE
METRIC (%)	Acc	ACC <sub>s</sub>	B-S	Acc	ACC <sub>s</sub>	B-S	Acc	ACC <sub>s</sub>	B-S	ACC <sub>s</sub>	ACC <sub>s</sub>	ACC <sub>s</sub>
<b>HUMAN</b>	71.0	<b>87.0</b>	87.6	47.0	50.0	86.8	<b>71.0</b>	73.0	80.3	<b>85.0</b>	<b>55.0</b>	<b>78.0</b>
<b>SUPERVISED METHODS</b>												
BERT <sub>B</sub> : <i>fine-tuned</i>	85.0	66.8	-	78.8	63.1	-	53.6	68.5	-	68.7	<b>52.2</b>	36.1
BERT <sub>L</sub> : <i>fine-tuned</i>	85.1	67.2	-	82.6	63.8	-	51.5	69.1	-	74.1	41.7	34.2
T5 <sub>B</sub> : <i>fine-tuned</i>	-	-	86.8	-	-	87.2	-	-	<b>89.7</b>	-	-	-
T5 <sub>L</sub> : <i>fine-tuned</i>	-	-	87.1	-	-	87.7	-	-	89.8	-	-	-
<b>PROMPT-BASED METHODS</b>												
Qwen3-235B: <i>zero-shot</i>	64.1	53.6	86.7	58.0	<b>25.9</b>	90.3	52.7	45.3	93.5	56.8	19.4	39.1
DeepSeek-R1: <i>zero-shot</i>	71.1	69.4	<b>85.1</b>	66.6	31.5	90.2	60.2	51.3	91.3	76.8	26.7	50.5
↪ + <i>three-shot</i>	79.1	70.6	88.1	76.4	55.6	91.6	62.7	66.3	96.3	74.7	26.7	<b>59.1</b>
↪ + <i>five-shot</i>	84.3	72.3	88.0	76.1	64.3	91.8	60.6	70.7	96.6	81.6	35.8	57.1
Kimi-K2-Instruct: <i>zero-shot</i>	68.5	63.1	86.7	68.5	34.4	90.3	60.6	45.4	95.6	55.8	28.9	46.7
↪ + <i>three-shot</i>	77.7	68.9	88.4	79.0	67.9	92.8	59.3	64.4	96.7	79.5	39.4	43.8
↪ + <i>five-shot</i>	81.7	69.6	88.2	79.7	69.2	92.3	64.7	63.6	97.2	81.1	43.3	46.7
Gemini-3-27B-it: <i>zero-shot</i>	55.0	57.3	86.4	58.0	38.4	89.5	58.3	39.9	92.1	66.8	19.4	38.1
↪ + <i>three-shot</i>	69.6	62.0	88.1	70.1	63.7	91.1	56.7	57.2	95.3	74.1	28.3	45.7
↪ + <i>five-shot</i>	72.1	61.6	87.9	70.8	68.2	90.7	56.2	59.2	95.9	70.5	35.0	52.4
Gemini-1.0-pro: <i>zero-shot</i>	56.0	77.8	86.9	48.5	51.8	89.5	<b>38.5</b>	59.0	91.8	43.8	<b>6.7</b>	43.8
Gemini-2.5-pro: <i>zero-shot</i>	55.0	65.6	87.4	71.5	52.1	89.4	65.6	61.2	93.7	<b>42.6</b>	27.4	42.9
Claude-Instant-1: <i>zero-shot</i>	51.2	72.2	85.7	40.5	42.6	89.7	43.2	50.9	91.9	59.2	11.6	39.0
↪ + <i>three-shot</i>	47.9	60.8	86.5	49.8	54.7	<b>87.0</b>	47.8	59.1	94.1	48.9	18.8	35.5
↪ + <i>five-shot</i>	52.0	<b>47.4</b>	87.0	50.1	57.7	87.1	44.9	61.8	94.5	53.1	15.0	38.4
Claude-3-Opus: <i>zero-shot</i>	66.3	62.8	87.1	61.3	34.7	88.5	50.4	36.3	91.7	67.3	28.3	42.8
↪ + <i>three-shot</i>	75.8	64.8	88.1	69.5	56.7	92.8	56.7	33.6	93.1	74.7	37.2	47.6
↪ + <i>five-shot</i>	72.8	67.1	88.2	69.8	60.0	92.8	63.9	<b>30.9</b>	96.0	75.7	35.5	43.2
Claude-Sonnet-4.5: <i>zero-shot</i>	72.5	68.5	87.2	67.5	40.1	88.9	51.0	45.1	94.4	69.8	16.1	41.9
↪ + <i>three-shot</i>	77.7	72.0	88.4	77.1	70.5	91.8	61.4	59.3	96.8	76.8	30.6	42.9
↪ + <i>five-shot</i>	78.0	72.0	88.5	76.1	72.7	90.9	<b>70.1</b>	62.1	<b>97.6</b>	<b>82.0</b>	37.2	47.6
GPT-3.5-Turbo: <i>zero-shot</i>	<b>40.6</b>	68.9	85.6	<b>32.4</b>	53.4	88.9	41.9	67.2	91.4	60.0	7.7	42.8
↪ + <i>three-shot</i>	45.4	67.3	88.2	36.3	69.5	92.4	43.6	72.7	96.5	53.6	10.0	<b>30.4</b>
↪ + <i>five-shot</i>	46.5	67.7	88.3	40.9	71.1	92.4	39.1	69.1	96.9	58.9	15.0	31.4
GPT-4: <i>zero-shot</i>	66.3	75.1	86.5	53.4	70.1	89.4	53.4	75.4	89.9	61.5	7.7	42.8
↪ + <i>three-shot</i>	70.3	77.1	88.1	60.0	<b>77.7</b>	<b>92.9</b>	56.3	83.6	94.8	75.8	16.1	43.8
↪ + <i>five-shot</i>	72.8	72.7	88.4	58.1	<b>77.8</b>	92.7	58.6	<b>85.4</b>	95.5	77.8	13.3	48.5
OpenAI o3 : <i>zero-shot</i>	57.1	65.1	86.5	72.1	37.7	89.8	65.2	62.9	93.8	67.9	25.6	51.4
↪ + <i>three-shot</i>	79.5	77.4	88.5	85.9	65.3	92.6	58.9	77.5	96.0	76.3	29.1	52.4
↪ + <i>five-shot</i>	83.5	74.7	88.6	83.6	71.5	91.6	63.5	78.6	96.5	77.3	36.9	50.0
GPT-5: <i>zero-shot</i>	82.8	67.6	86.6	75.4	36.7	89.9	66.8	64.3	93.3	74.2	28.9	56.2
↪ + <i>three-shot</i>	82.1	<b>78.3</b>	<b>88.7</b>	<b>86.2</b>	67.2	92.3	61.8	77.1	96.4	74.7	33.3	51.4
↪ + <i>five-shot</i>	<b>85.4</b>	<b>78.7</b>	88.6	84.3	68.9	92.3	67.2	79.0	96.8	74.7	38.3	50.5

Table 12: Complete Experimental Results in SEMANTICQA. “-” denotes the model that is unavailable or inappropriate for the task. **Digits** highlight cases in which human scores are higher than those of all evaluated models, serving as a coarse reference. Light Pink text delineates the baselines with supervised fine-tuning. Light Green and Light Blue parts present open-source models and proprietary models.

System	Acc@1	Acc@2	Acc@4	Acc@8	Acc@16
<b><i>Baselines</i></b>					
Random	100.0	50.0	25.0	12.5	6.3
Majority	100.0	50.0	25.0	12.5	6.3
<b><i>Small language models</i></b>					
BERT <sub>B</sub>	100.0 <sub>0,0</sub>	98.9 <sub>1,9</sub>	89.4 <sub>4,9</sub>	79.9 <sub>6,4</sub>	69.9 <sub>0,0</sub>
BERT <sub>L</sub>	100.0 <sub>0,0</sub>	<b>98.9<sub>1,0</sub></b>	<b>95.8<sub>1,4</sub></b>	<b>83.5<sub>5,2</sub></b>	<b>71.8<sub>0,0</sub></b>
<b><i>Large language models</i></b>					
DeepSeek-R1					
↔ + 0-shot	100.0 <sub>0,0</sub>	81.7 <sub>11,9</sub>	68.9 <sub>4,2</sub>	49.3 <sub>4,0</sub>	35.4 <sub>0,0</sub>
↔ + 3-shot	100.0 <sub>0,0</sub>	83.8 <sub>7,5</sub>	73.7 <sub>2,4</sub>	52.0 <sub>6,2</sub>	45.9 <sub>0,0</sub>
↔ + 5-shot	100.0 <sub>0,0</sub>	80.6 <sub>11,8</sub>	68.4 <sub>7,8</sub>	51.7 <sub>7,7</sub>	47.3 <sub>0,0</sub>
Kimi-K2-Instruct					
↔ + 0-shot	100.0 <sub>0,0</sub>	85.6 <sub>14,6</sub>	71.1 <sub>8,2</sub>	50.4 <sub>4,6</sub>	44.6 <sub>0,0</sub>
↔ + 3-shot	100.0 <sub>0,0</sub>	92.8 <sub>6,7</sub>	78.3 <sub>3,8</sub>	61.5 <sub>7,9</sub>	49.6 <sub>0,0</sub>
↔ + 5-shot	100.0 <sub>0,0</sub>	87.2 <sub>4,2</sub>	77.8 <sub>6,9</sub>	63.8 <sub>5,2</sub>	51.7 <sub>0,0</sub>
OpenAI o3					
↔ + 0-shot	100.0 <sub>0,0</sub>	93.3 <sub>8,3</sub>	82.8 <sub>4,6</sub>	65.4 <sub>2,5</sub>	53.3 <sub>0,0</sub>
↔ + 3-shot	100.0 <sub>0,0</sub>	91.7 <sub>5,9</sub>	85.3 <sub>3,1</sub>	72.6 <sub>3,4</sub>	62.9 <sub>0,0</sub>
↔ + 5-shot	100.0 <sub>0,0</sub>	91.1 <sub>5,5</sub>	88.6 <sub>4,5</sub>	74.0 <sub>1,6</sub>	64.6 <sub>0,0</sub>
GPT-5					
↔ + 0-shot	100.0 <sub>0,0</sub>	92.2 <sub>6,3</sub>	84.2 <sub>5,9</sub>	67.7 <sub>3,4</sub>	56.3 <sub>0,0</sub>
↔ + 3-shot	100.0 <sub>0,0</sub>	92.8 <sub>9,1</sub>	88.6 <sub>3,8</sub>	<u>74.6<sub>2,8</sub></u>	<u>65.8<sub>0,0</sub></u>
↔ + 5-shot	100.0 <sub>0,0</sub>	<u>94.4<sub>6,7</sub></u>	<u>89.2<sub>4,1</sub></u>	73.5 <sub>3,6</sub>	65.2 <sub>0,0</sub>

Table 13: Our best experimental results (avg<sub>std</sub>). The mean accuracy scores with their standard deviation are computed by averaging the results of three independent runs with different random seeds. Results of baselines are also provided including random choice as well as the majority of class instances over each sub categorization tasks. The **Bold** and underlined texts denote the best and second-best performance in the specific category, respectively.

<b>SP Category</b>	<b>SP Subtype</b>	<b>Representative Literature</b>
<b>Idiomatic Expressions</b>	Opaque	Nunberg et al. (1994)
	Semi-transparent	Nunberg et al. (1994)
	Decomposable	Nunberg et al. (1994)
	Pragmatic	Wray (2002)
	Formulaic	Wray (2002)
	Figurative	Sag et al. (2002b)
	Eventive	Chen et al. (2017)
	Stative	Spathas and Michelioudakis (2021)
Property	Fazly et al. (2009)	
<b>Lexical Collocations</b>	Intensification	Mel'čuk (1998)
	Veracity	Mel'čuk (1998)
	Evaluation	Mel'čuk (1998)
	Operative	Mel'čuk (1998)
	Causative	Mel'čuk (1998)
	Emission	Mel'čuk (1998)
	Typicality	Mel'čuk (1998)
	Structural pattern	Mel'čuk (1998)
<b>Noun Compounds</b>	Material	Tratz and Hovy (2010)
	Purpose	Tratz and Hovy (2010)
	Source	Tratz and Hovy (2010)
	Topic	Tratz and Hovy (2010)
	Location	Tratz and Hovy (2010)
	Temporal	Tratz and Hovy (2010)
	Agentive	Tratz and Hovy (2010)
	Part-Whole	Tratz and Hovy (2010)
	Possessive	Tratz and Hovy (2010)
	Proper compound	Tratz and Hovy (2010)
<b>Verbal MWEs</b>	Light Verb	Savary et al. (2017)
	Verb-Particle	Savary et al. (2017)
	Verbal Idiom	Savary et al. (2017)

Table 14: Correspondence between the proposed semantic phrase subtypes and related categories in established MWE literature with representative references.