

Expanding Horizons or Hitting Walls? Limits and Potentials of LLMs in Augmenting Lexical Knowledge Bases

Anonymous ACL submission

Abstract

This paper investigates the potential of Large Language Models (LLMs) to augment lexical knowledge bases (KBs) and to address their common limitations, such as static nature, limited coverage, and labor-intensive creation and maintenance. We propose a methodology that leverages LLMs to accurately reconstruct information from a source KB and generate new knowledge. Then, we evaluate this methodology using various LLMs and prompting techniques across three separate KBs. The results suggest that LLMs can accurately provide information when given ample contextual cues and when dealing with high-specificity concepts. However, they are prone to errors and inconsistencies when asked for rare or generic knowledge. The findings also indicate that LLMs can contribute to KB management by reducing the need for manual intervention. This study highlights the potential and limitations of LLMs in lexical semantics and emphasizes the importance of novel approaches to KB creation, maintenance, and integration.

1 Introduction

Lexical semantics represents a foundational aspect of Natural Language Processing (NLP), serving as the intersection where the meanings of words and their interrelationships converge. This discipline has always seen unstructured data become structured through the means of knowledge bases (KBs). These latter ones, however, face three common limitations: *i*) they exhibit a static nature, making it challenging to adapt to domains evolution; *ii*) they suffer from limited coverage, hindering their applicability across diverse domains; *iii*) their creation and maintenance is typically laborious, involving human-in-the-loop procedures.

The rise of Large Language Models (LLMs) within Generative AI highlights the importance of interpretable knowledge encapsulation, with KBs being crucial for both enhancing LLM training and

providing a means of error, inconsistency, and bias checking (Pan et al., 2024). This necessity becomes particularly pronounced given the expanding influence of Generative AI and its accompanying challenges, including issues such as hallucination (Ji et al., 2023).

This paper unveils an innovative methodology grounded in LLMs to tackle pivotal concerns within lexical semantics. In particular, our contribution is three-fold: *i*) we harness LLMs to reconstitute information encapsulated in a source KB to test their proficiency on this task; *ii*) subsequently, LLMs are deployed to create novel knowledge, proving their aptitude in crafting, and encoding KBs; *iii*) through a third-phase assessment of the newly generated content, we can finally evaluate the capability to expand upon the original KBs and, consequently, assess their completeness.

By conducting thorough experiments utilizing diverse LLMs and prompting techniques across three separate KBs, we elucidate the capacity of LLMs to furnish accurate information, particularly when supplied with substantial contextual cues and when dealing with concepts of high specificity. When confronted with requests for rare or generic knowledge, LLMs are instead prone to errors and inconsistencies.

2 Related Work

In the context of this work, it is essential to clarify that lexical semantic resources, KBs, ontologies and knowledge graphs represent facets or interpretations of the same underlying subject matter.

2.1 Knowledge Acquisition: KBs and KGs

Construction of KBs involves both manual and automatic methods, with famous KBs like WordNet (Fellbaum, 2020) and ConceptNet (Speer et al., 2017) initially depending on manual input. To reduce labor, automated IE techniques have been developed (Fader et al., 2011; Angeli et al., 2015;

Vo and Bagheri, 2017), extracting information from texts to update KBs. ML and NLP progress has also advanced in automatic Knowledge Graph (KG) construction, utilizing data to enhance traditional approaches (Chen et al., 2021).

2.2 Large Language Models

Recently, the advent of LLMs has opened new avenues for knowledge acquisition and representation. LLMs, such as GPT-3 (Brown et al., 2020) and LLama-2 (Touvron et al., 2023), have demonstrated remarkable capabilities in understanding and generating natural language text. Researchers have begun exploring the knowledge encoded within LLMs, probing their ability to serve as implicit KBs (Petroni et al., 2019; Razniewski and Weikum, 2021). This approach offers a novel means of accessing vast amounts of knowledge without explicit curation, although challenges remain in interpreting and validating the knowledge encoded in these models (Chang et al., 2023).

Despite the speed of breakthrough advancements in the field, LLMs still grapple with issues that fall into two main categories: architectural and data-related problems. Architectural problems are inherent to the model’s structure and necessitate a change in architecture for resolution. These include the prompt engineering problem, wherein models are non-deterministic and require the "perfect" prompt to elicit the correct response, as highlighted by Park et al. (2022). Conversely, data-related problems stem from the training methodologies and the datasets used, affecting the models’ mathematical and reasoning capabilities (Imani et al., 2023; Hendrycks et al., 2021), as well as their common sense understanding (West et al., 2022).

2.3 KBs/KGs and LLMs

Petersen and Potts (2023) demonstrate LLMs’ capability to interpret the word “break” and suggest that these models can advance lexical semantics. Their analysis reveals LLMs’ proficiency in identifying both known and novel meanings, as well as their superiority in semantic analysis. Kandpal et al. (2023) indicate that the knowledge representation in LLM training data affects content generation accuracy, particularly for uncommon concepts, challenging our understanding of LLMs’ semantic encoding. Cohen et al. (2023) propose a new approach for examining LLM knowledge using graph-based queries, which aligns with our emphasis on structured prompts to retrieve and leverage

knowledge from LLMs. In contrast to previous efforts, our work presents a scalable pipeline for standardized KB extension, leveraging LLMs and human-in-the-loop evaluation techniques.

3 A Methodology for KB Extension

Our proposed methodology encompasses different key modules, answering the following two main research questions: *RQ1) How well can LLMs mimic KB concepts and relationships?*; and *RQ2) Do LLMs possess the capability to produce novel information suitable for integration into existing KBs?* However, these inquiries serve as gateways to further exploration. Particularly in relation to *RQ1: what factors of both LLMs and KBs impact the quality of generated content?* This paper delves into the following considerations: LLM architecture (pre-trained, fine-tuned, and storytelling-oriented); prompting and extraction techniques (zero-shot versus one-shot); as well as the scale and intricacy of the KBs. Additionally, in relation to *RQ2: how does the quality of newly generated content compare to that of the original KB?* A manual assessment could provide insights into the completeness of the original resource, enhancing the proposed framework.

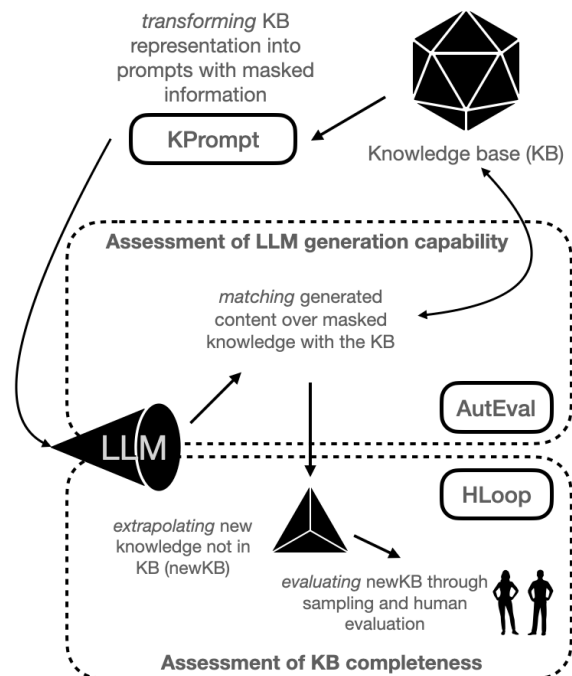


Figure 1: Architecture of the proposed framework for KBs extension. i) (*KPrompt*) encodes the source knowledge into masked prompts, that LLMs use for ii) (*AutEval*) re-generating existing knowledge and iii) (*HLoop*) generate new content to be manually-validated.

On this direction, the proposed framework includes three modules which are designed *i*) to systematically assess the proficiency of LLMs in delivering concepts aligned with existing KBs, *ii*) to probe their ability to generate novel concepts for potential integration into the KB, and *iii*) to ascertain the role of limited manual intervention in evaluating the completeness and coverage of the KB. An overview of the architecture is shown in Figure 1, while each module is presented in the following sections.

3.1 Knowledge Base-to-Prompt (*KPrompt*)

The first module involves the development of a Knowledge Base-to-Prompt strategy (*KPrompt*), which serves as the bridge between the lexical knowledge stored in the KB and the queries posed to LLMs. This strategy aims to convert the structured information within the KB into prompts that effectively capture the nuances and intricacies of the underlying semantic content. The objective is to enable LLMs to generate responses that align with the pre-existing knowledge stored in the KB, thus addressing the fundamental question of whether LLMs can proficiently deliver concepts consistent with the KB. For example, if a resource holds the information that x is connected with y through a semantic relation r (or, more formally, $r(a, b)$), then a generic template prompt for extracting b -candidates could be the following:

Given the relation r with the specific meaning $\langle r_description \rangle$, which concepts (like b) might be also connected to a through r ?

Depending on the kind of knowledge encoded in the target KB, this template may be adapted in different ways and through prompting strategies such as zero- and one-shot, which are defined later on in the experimental sections.

3.2 Automatic Evaluation (*AutEval*)

The second module focuses on an intrinsic evaluation of the LLM-based extension via knowledge masking (*AutEval*), assessing the capacity of LLMs to obtain both correct and novel knowledge by first masking existing semantic units in the source knowledge and then asking LLMs to generate possible candidates (see Section 3.1 example). By systematically matching the generated candidates with the original KBs, we assess the LLM’s capability to generate correct and existing information, as in the following example:

Here is a list of candidates to connect with $\langle a \rangle$ through $\langle r \rangle$: x, y, z

By checking the presence of x, y and z in the source KB (i.e. specifically of $r(a, x), r(a, y)$ and $r(a, z)$), it will be possible to give some answers related to the first research question *RQ1*.

3.3 Human-in-the-loop Strategy (*HLoop*)

The third module incorporates a human-in-the-loop strategy (*HLoop*) to evaluate the novel LLM-extracted knowledge not covered by the source KB, aiming to answer the nuanced question of whether LLMs can effectively extend the KB and, simultaneously, serve as a means to verify its completeness. In particular, human evaluators, through limited manual intervention, are asked to assess the relevance and accuracy of the novel LLMs-generated knowledge.

Continuing with the example of Section 3.2, if $r(a, y)$ and $r(a, z)$ are found not to be included in the source KB, a focused manual examination of such new content may be conducted to evaluate their accuracy.

4 Experimentation

In this section, we detail the experiment settings, i.e. the implementation of the three modules (*KPrompt*, *AutEval* and *HLoop*) on three knowledge bases: Semagram (Leone et al., 2020), MultiAlignNet (Grasso et al., 2022), and ConcepNet (Speer et al., 2017). The selection of these KBs has been carefully done by considering features such as scale and complexity of the encoded knowledge. By experimenting on this diversity, we aim to highlight insights and challenges under a reliable lens.

All the code for our experiments is openly available at <https://anonymous.4open.science/r/LLM-Semagram-2C44/>.

4.1 Prompting Strategy

KBs encapsulate complex real-world information by codifying semantic relationships, presenting a challenge for LLMs, which are typically tailored to process natural language. No standard prompting method yet exists for repurposing KB data to align with LLMs’ textual processing. Our *KPrompt* methodology transfigures KB data into structured prompts for LLMs, however we do not propose *KPrompt* as a definitive standard but rather as an

253	innovative step towards bridging the gap between	4.2.2 MultiAlignNet	303
254	KB-based data and language model processing.	The MultiAlignNet KB, introduced by Grasso et al.	304
255	Within the plethora of prompt engineering	(2022), constitutes a recently-developed lexical-	305
256	methodologies present in literature, we selected	semantic resource constructed using plain textual	306
257	those that do not require an interactive dialogue	information gathered from several corpora in mul-	307
258	with a LLM and are considered state of the art:	iple languages. It encompasses knowledge across	308
259	Zero-Shot prompt (Kojima et al., 2022), and Few	1,047 noun concepts called <i>heads</i> and it results in	309
260	Shot prompt (Min et al., 2022; Touvron et al., 2023).	21,514 interconnected concepts. It is also linked	310
261	In our prompts, we instruct LLMs to return 10 con-	to WordNet (Fellbaum, 2020) and BabelNet (Nav-	311
262	cepts in order to align them with the automatic	igli and Ponzetto, 2010) synsets. In a simplified	312
263	evaluation in Section 5.1.	depiction, its internal framework resembles a KG	313
264	The output generated by a LLM typically con-	comprising three primary node types — <i>noun</i> , <i>verb</i>	314
265	sists of plain text that enumerates various concepts.	and <i>adjective</i> nodes, alongside two distinct relation-	315
266	To isolate these concepts (or entities), we employ	ship types — paradigmatic and syntagmatic. Our	316
267	regular expressions, which serve as a necessary	experiment centered on the latter category, formu-	317
268	step due to the model’s potential to “hallucinate”	lating prompts such as the following:	318
269	- that is to append extraneous descriptions to the	Provide a list of 10 English nouns	319
270	actual list of concepts. To address this, we crafted	related to the concept “shape, form,	320
271	the following regular expression “\b\w+\b”. We	configuration” in the form of a	321
272	also experimented with a simpler one, \w+, but it	comma-separated list of lowercase lemmas.	322
273	yielded sub-optimal results across different KBs.	Examples: solubility, mean, packing,	323
274	4.2 Knowledge Bases	weight, load, color, size, style, art	324
275	In this section, we overview the KBs chosen for	4.2.3 ConceptNet	325
276	experimentation.	ConceptNet, introduced by Speer et al. (2017),	326
277	4.2.1 Semagram	serves as a multilingual KB that captures the con-	327
278	The Semagram KB, introduced by Leone et al.	nections and common-sense relationships among	328
279	(2020), boasts a versatile structure that captures the	words. The inclusion of words and relationships	329
280	semantics of a given concept through a slot-filler	stems from diverse sources, ranging from crowd-	330
281	representation. The current version encompasses	sourced inputs to expert-generated content. The	331
282	over 300 concepts and 26 slots (i.e., semantic re-	dataset boasts more than 21M edges and over 8M	332
283	lationships). Each concept is also interconnected	nodes, with the English vocabulary alone compris-	333
284	with other resources, e.g. BabelNet (Navigli and	ing around 1,5M nodes. ConceptNet is character-	334
285	Ponzetto, 2010). Following (Ventrice and Siragusa,	ized by two fundamental types of relations: sym-	335
286	2023), we observed that these descriptions adhered	metric relations and asymmetric relations. In par-	336
287	to straightforward ontology relations; for example,	ticular, we focused on <i>UsedFor</i> (symmetric) and	337
288	the <i>material</i> slot could be translated as “ <i>can be</i>	<i>RelatedTo</i> (asymmetric).	338
289	<i>made of</i> ”. Consequently, we opted to craft simple	We then designed a straightforward prompt that	339
290	sentences, each posing a criterion to the LLM. Each	receives a concept as input and instructs the LLM to	340
291	criterion was then coupled with all its associated	identify 10 concepts that possess either a “ <i>related</i>	341
292	fillers. Subsequently, we devised a concise prompt:	<i>to</i> ” or “ <i>used for</i> ” relation with the given concept.	342
293	Provide a list of 10 words that satisfy	An example of prompt (“ <i>used for</i> ”) is as follow:	343
294	the condition.	Given the concept ‘car’, list 10	344
295	Desired output: comma-separated list of	concepts for which ‘car’ is used for, in	345
296	words	the form of a comma-separated list.	346
297	Condition: can be made of wood	4.3 LLMs Selection	347
298	Here, <i>condition</i> encompasses the textual interpre-	Among all the different types of openly available	348
299	tation of the corresponding slot. This prompt struc-	LLMs, one way to select the optimal model is	349
300	ture serves as a streamlined and effective means to	through Open LLM Leaderboard , a widely recog-	350
301	elicit targeted responses from the LLM based on	nised LLM competition list. At the time of the	351
302	the semantics encoded in Semagram.	selection of the model, the LLama-2 architecture	352

(Touvron et al., 2023) and the models fine-tuned from its associated weights were the highest ranked. The models are filtered by: pre-trained, fine-tuned on domain-specific datasets and chat models. The second and third categories are both derived from a fine-tuning on the first one. Another sub-distinction that we argue being ever so important in the current LLM panorama are storytelling fine-tuned models (Xie et al., 2023). One goal of this paper is also to discover if these models can enhance the capabilities to carry out the task under study.

For our purpose, we used three principles for LLMs selection: *i)* State-of-the-art for their respective categories at the time of selection. These were selected via an average score over different benchmarks for language capabilities of LLMs: ARC (Chollet, 2019), HellaSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2020), TruthfulQA (Lin et al., 2022), and WinoGrande (Sakaguchi et al., 2021). *ii)* With more than 30B parameters. This is justified by empirical evidence suggesting that larger language models tend to outperform smaller ones across various language tasks. *iii)* Pertaining to the three categories illustrated in Section 4.3.

For the aforementioned reasons, our choices fell on: *i)* **Yi-34B**: a model trained from scratch with the Llama-2 architecture; *ii)* **Tulu-2-70B**: a model that combines instruction and RLHF tuned chat models on a mix of publicly available, synthetic and human-created datasets; *iii)* **Aetheria-L2-70B**: a model specifically tailored for storytelling, that combines Euryale v1.3 base with the DPO training of the Tulu v2 model, and the GOAT Storytelling model. The LimaRP v3 QLoRA was then added (RoyalLab, 2023).

Each model was tasked to process the full set of prompts for each KB, with the sampling parameters configured to a top-p (nucleus) sampling value of 0.95, temperature of 0.4, PagedAttention enabled (Kwon et al., 2023), and a maximum token limit of 100. We used AWQ 4-bit quantization (Lin et al., 2023) to reduce memory utilization without losing language capability (Yao et al., 2023).

5 Evaluation

Within our framework, we assess two key aspects: *i)* the proficiency of the LLMs in accurately extracting verified knowledge from the KB (*AutEval*), and *ii)* the extent to which novel knowledge is extracted that was not originally encoded in the KB (*HLoop*), through manual annotation.

5.1 The *AutEval* process

We adopted standard evaluation metrics to assess the performance of LLMs. Formally, let p represent a prompt from the set P , $C_p = \{c_1, c_2, \dots, c_n\}$ denote the list of concepts returned by the LLM, $C_k = \{c_1, \dots, c_k\}$ the set of the first k concepts of C_p , and $K_p = \{k_1, k_2, \dots, k_m\}$ denote the list of concepts existing in the KB and related to the prompt. The metrics are defined as follows: **Precision@K**: proportion of the returned items in the top- k (C_k) that are actually relevant; **Recall@K**: proportion of relevant items found in the top- k recommendations (K_q); and **F-Measure@K**: the harmonic mean between Precision and Recall. We also provide an asymptotic Recall value based on truncating the concept list to 10 items in the prompt.

5.1.1 Semagram

Table 1 presents the performance scores for both zero-shot and one-shot prompts across the three LLMs. The scores in the table reveal that Tulu outperforms the other models. Additionally, the performance disparity between zero-shot and one-shot prompts is marginal (0.41 percentage points for F1@10).

Prompt		P@10	R@10	F1@10
Zero-shot	Aetheria	8.73	44.00	14.57
	Tulu	9.59	46.36	15.89
	Yi	5.82	25.30	9.46
One-shot	Aetheria	7.52	36.21	12.45
	Tulu	9.42	43.36	15.48
	Yi	5.9	28.11	9.75

Table 1: Results on Semagram for Precision, Recall and F1-Measure at 10.

Figure 2 illustrates the Precision and Recall scores for varying values of k : 1, 2, 5, and 10. A notable trend is that at lower values of k , one-shot prompts outperform the zero-shot in both metrics. This trend reverses when $k > 5$.

5.1.2 MultiAlignNet

We constructed distinct prompts for each syntactic category—*noun*, *verb*, and *adjective*—to evaluate potential performance discrepancies across these types. Their outcomes are detailed in Table 2. Our findings suggest that model efficacy varies significantly with syntactic category; specifically, Aetheria and Tulu demonstrate superior Precision and Recall for *noun* nodes, outstripping *verb* and

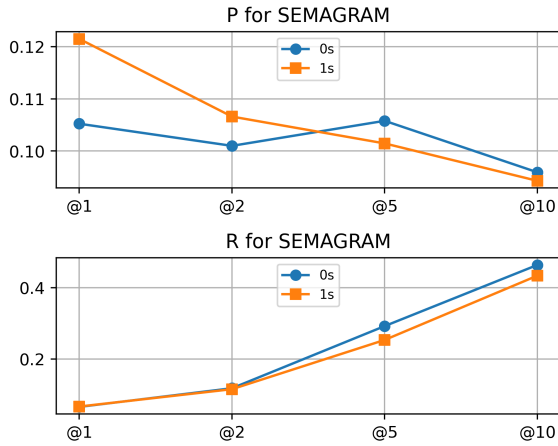


Figure 2: Precision and Recall on Semagram as k increases (best performing model). The asymptotic Recall value (at 10) is 0.98.

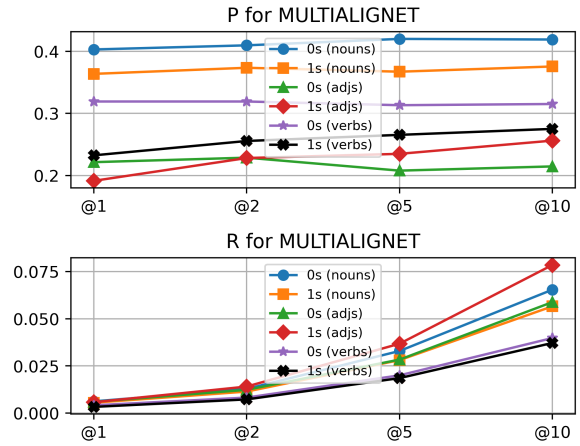


Figure 3: Precision and Recall on MultiAlignNet as k increases (best performing models). The asymptotic Recall value (at 10) for *nouns* is 0.17, for *adjectives* is 0.32, for *verbs* is 0.14.

441 *adjective* nodes. Conversely, for *verb* and *adjective*
 442 nodes, Aetheria and Yi lead in zero-shot and
 443 one-shot settings. Overall, the zero-shot prompting
 444 strategy on Aetheria demonstrates the most efficacy,
 445 except for *adjective* nodes.

Prompt		P@10	R@10	F1@10
Zero-shot (Nouns)	Aetheria	41.89	6.54	11.30
	Tulu	39.47	6.00	10.40
	Yi	12.21	1.78	3.11
Zero-shot (Adjs)	Aetheria	21.45	5.88	9.23
	Tulu	16.16	4.36	6.87
	Yi	4.74	1.24	1.97
Zero-shot (Verbs)	Aetheria	31.5	3.98	7.07
	Tulu	30.20	3.50	6.27
	Yi	8.58	1.02	1.83
One-shot (Nouns)	Aetheria	36.83	5.52	9.59
	Tulu	37.55	5.66	9.84
	Yi	26.23	4.25	7.33
One-shot (Adjs)	Aetheria	20.54	4.94	7.97
	Tulu	22.32	5.77	9.18
	Yi	25.61	7.84	12.00
One-shot (Verbs)	Aetheria	27.29	3.20	5.73
	Tulu	26.86	3.29	5.86
	Yi	27.5	3.72	6.55

Table 2: Results on MultiAlignNet for Precision, Recall and F-Measure at 10.

446 Figure 3 presents the performance metrics across
 447 varying values of k . The observed trends align with
 448 those reported in Section 5.1.1, albeit with notable
 449 distinctions. For *verb* and *adjective* nodes, there is
 450 a minor but consistent enhancement in Precision.

5.1.3 ConceptNet

451 As discussed in Section 4.2.3, we focused on *RelatedTo*
 452 and *UsedFor* relationships. Table 3 presents
 453 the results, where Tulu, employing a one-shot
 454 prompt strategy, achieves the highest scores. The
 455 performance gap between zero-shot and one-shot
 456 is about 1.03 percentage points. Differently, Yi
 457 demonstrated a notably poorer performance. Fi-
 458 nally, the one-shot strategy aided the models in
 459 comprehending the task, thereby enhancing their
 460 ability to retrieve more accurate concepts.

461 In Figure 4, Tulu consistently shows the high-
 462 est Precision and Recall values for the *RelatedTo*
 463 relationship across different values of k . The Re-
 464 call scores remain comparable up to $k = 2$, with a
 465 slight improvement for the one-shot prompt beyond
 466 this point. Differently, the *UsedFor* relationship ex-
 467 hibits lower performance overall, with a peak on
 468 Precision at $k = 2$.
 469

5.2 The HLoop process

470 Our comprehensive automatic evaluation indicated
 471 a consistent pattern of moderate to low Precision
 472 and Recall scores across both fine-tuned (Aetheria,
 473 Tulu) and non-fine-tuned (Yi) models. Such results
 474 pointed out the necessity for an additional layer
 475 of scrutiny. Therefore, we structured a manual
 476 evaluation to understand whether the unsatisfactory
 477 scores stemmed from model errors or the genera-
 478 tion of novel data absent from the KBs.
 479

480 To this end, we selected a sample of 300 prompt
 481 outputs, from the best performing models, refer-
 482 ring respectively to each of the three employed

Prompt		P@10	R@10	F1@10
Zero-shot (RelTo)	Aetheria	19.62	10.68	13.83
	Tulu	20.25	10.38	13.72
	Yi	9.37	4.60	6.17
Zero-shot (UsedFor)	Aetheria	6.79	3.29	4.43
	Tulu	8.7	3.98	5.46
	Yi	1.92	0.8	1.13
One-shot (RelTo)	Aetheria	20.8	11.43	14.75
	Tulu	22.18	11.96	15.54
	Yi	8.16	4.07	5.43
One-shot (UsedFor)	Aetheria	7.06	3.43	4.62
	Tulu	8.46	4.09	5.51
	Yi	4.20	1.80	2.50

Table 3: Results on ConceptNet for Precision, Recall and F-Measure at 10.

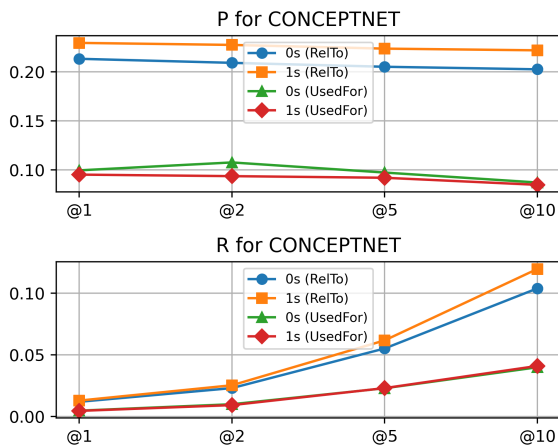


Figure 4: Precision and Recall on ConceptNet as k increases (best performing model). The asymptotic Recall value (at 10) for *RelatedTo* is 0.67, for *UsedFor* is 0.56.

KBs. We then asked three annotators to examine the 900 prompt outputs and to determine if a generated concept is related to the target one within the predefined relationship. The potential verdicts for each entry were categorized as correct, incorrect, or misspelled, with the latter specifically denoting any grammatical inaccuracies introduced by the models.

We used Fleiss (1971)’s kappa ($F-\kappa$) and Randolph (2005)’s multirated kappa ($R-\kappa$) to evaluate the Inter Annotator Agreement (IAA). Both metrics provide a lower and upper bound on the IAA. For Semagram, we obtained an $F-\kappa$ of 0.43 – moderate agreement – and an $R-\kappa$ of 0.60 – substantial agreement. MultiAlignNet has an $F-\kappa$ of 0.51 – moderate agreement – and an $R-\kappa$ of 0.64 – sub-

stantial agreement. Finally, ConceptNet obtained the lowest agreement scores, having an $F-\kappa$ of 0.33 – fair agreement – and an $R-\kappa$ of 0.5 – moderate agreement.

Figure 5 shows the results of the annotation. Both Semagram and MultiAlignNet have almost the same amount of correct and incorrect concepts, with a small difference on the misspelled (8 vs 6). Having such a small pool of misspelled concepts demonstrates that our prompting methodology does not elicit word hallucinations. ConceptNet, instead, has a large amount of incorrect concepts (179 incorrect vs 115 correct). These results, discussed in Section 5.3, show that LLMs “prefer” more fine-grained and specific relationships, whereas they hallucinate on more generic and abstract relationships.

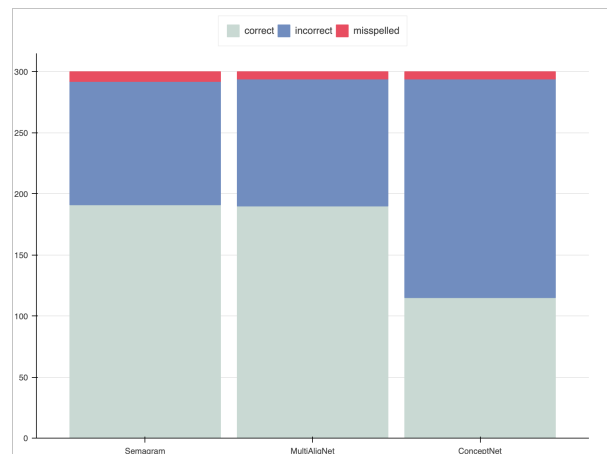


Figure 5: The ratio of correct, incorrect, and misspelled concepts on the three KBs.

5.3 Discussion

Throughout the evaluation phases, we noticed further interesting patterns that serve as additional contributions of the paper.

5.3.1 KB Size and Recall Relationship

From the observations gathered across varying sizes of KBs, we noticed an increase of Recall rates over the KB sizes. This phenomenon can be attributed to bigger data coverage within the KB (i.e., the greater the volume of entries within a KB, the higher the probability that it encompasses information pertinent to a broader spectrum of queries (Kandpal et al., 2023)).

5.3.2 Fine-tuned Models are Better

Another confirmed hypothesis regards the relationship between LLMs fine-tuning and performance in

generating new and correct data for the KBs. Both Tulu and Aetheria performed 2 and 2.1 times better than the standard pre-trained Yi, given their Recall scores on the three KBs. Finally, the storytelling model, Aetheria, performed slightly better.

5.3.3 LLMs Cannot Improvise

LLMs could only recall up to 40% of the information within the tested resources, often falling short of this benchmark, verging on null accuracy. This underscores a noteworthy constraint in their capabilities. Even with guidance (i.e. using a one-shot prompt), their low recall capability does pose three possible hypothesis: *i*) LLMs cannot infer semantic relationships; *ii*) the given prompts are not accurate enough to give an explanation of the task to the model, and so contextual information may be missing. LLMs can struggle with understanding and maintaining context, especially if the KB is complex and the semantic relationships are too general; *iii*) LLMs might misunderstand the meaning of terms within the KB and thus fail to recall relevant information that depends on a different interpretation of these terms.

All of the above can be in fact true knowing that *i*) LLMs do not "understand" semantics because they have no formal grounding or theory of mind (Pavlick, 2023; Ullman, 2023); *ii*) LLMs are heavily context dependent (Shi et al., 2023); *iii*) when unguided, they fail to resolve ambiguity in language most of the time (Zhao et al., 2021; Liu et al., 2021; Zhang et al., 2022).

5.3.4 Results from Manual Analysis

Diving deeper into the results, it is peculiar to see some divergency. On one hand, on KBs that can be considered "fine-grained" (e.g. Semagram, Multi-AligNet), LLMs seem to perform better and generate new and usable knowledge. On the other, common-sense KBs seem to heavily challenge LLMs because of the aforementioned missing context and semantic ambiguity.

This brought to a vast amount of confabulated content, mostly dependent to directionality problems in semantic relationships. This is due to the "Reversal Curse", discovered in (Berglund et al., 2023), that states "if a model has been trained on a sentence of the form "A is B", it will not automatically generalize to the reverse direction "B is A" whatsoever".

6 Conclusions

This study embarked on an exploration of Large Language Models to reconfigure lexical-semantic information, leveraging three existing resources. Our objectives encompassed assessing their efficacy and precision in this endeavor, alongside examining their ability to autonomously enrich these resources, as verified through human evaluation. Implicitly, our methodology also scrutinized the comprehensiveness of the original resources.

In synthesis, LLMs, in their current state of training, evince notable limitations in the domain of lexical semantics, extending beyond the realm of prompt variability and output alignment challenges (Kim et al., 2023). In resources teeming with rich contextual nuances (thus not solely reliant on decontextualized *x-rel-y* relations), LLMs manifest a pronounced capacity for generating novel knowledge. Finally, the moderate concordance among human evaluators in assessing LLM outputs underscores significant inadequacies within the resources themselves. The encoded information, or its attempted encoding, appears markedly unstable, subjective, and frequently incomplete, thereby signifying a pressing imperative for further refinement and augmentation.

6.1 Future Work

Although our paper exposes the limits of LLMs in generating data for existing KBs, recent work suggests that, while these models may not be suitable for truth-telling, they excel at revising incorrect data and identifying mistakes (Gou et al., 2023; Tyen et al., 2023). One potential future research direction is to modify our prompt engineering methodology from a purely zero/one-shot approach to a pipeline that incorporates both prompting and revision. Another research direction could be the use of RAG (Retrieval Augmented Generation) and KEG (Knowledge Enriched Generation) (Lewis et al., 2020; Yu et al., 2021; Gao et al., 2023) to enhance the context capabilities of these models. Retrieving current KB data instead of explaining through examples might be the key to unlock their capabilities. Finally, another direction involves the creation of a dataset of probing questions to assess the ability of LLMs in generating accurate and coherent data for KBs, serving as a benchmark to compare LLM performance and advance the development of more sophisticated models for this task.

630 Limitations

631 Our study introduces a novel approach to enhancing
632 lexical semantic KBs using LLMs, yet it comes
633 with several limitations that warrant attention. The
634 methodology is deeply entwined with the capabilities
635 of LLMs, meaning that any intrinsic limitations,
636 such as biases in training data or a lack of
637 deep context understanding, are directly reflected
638 in the quality of our generated knowledge. The
639 complexity of semantic relationships within the
640 KB also significantly influences LLM performance,
641 with fine-grained KBs yielding better results compared
642 to those with more abstract, common-sense relationships.
643

644 Furthermore, our focus on English-language resources
645 limits the applicability of our findings to KBs
646 in other languages, particularly those with complex
647 morphology or unique syntax. The validation process
648 also revealed a moderate level of agreement among
649 human annotators, highlighting the subjective nature
650 of interpreting LLM outputs, which could introduce
651 inconsistencies in the assessment of knowledge completeness
652 and validity.

653 Additionally, while human evaluation is critical
654 for ensuring the quality of LLM outputs, it is not
655 scalable and requires substantial manual effort,
656 posing a challenge for larger KBs or ongoing
657 maintenance. The issue of LLMs potentially generating
658 plausible but incorrect information, known as
659 "hallucinations" or using a better word, "confabulations"¹,
660 persists despite our efforts to minimize it through
661 strategic prompting. Finally, our study's success
662 hinged on the meticulous crafting of prompts, a
663 process lacking standardized best practices, and
664 remains a significant challenge in eliciting consistently
665 accurate and relevant responses from LLMs.
666

667 Ethics Statement

668 The experiment was designed with the idea of providing
669 beneficial knowledge and not harm any individual
670 or group. Our primary goal was to develop a methodology
671 for expanding the coverage of existing lexical KBs
672 using Large Language Models (LLMs). We recognize
673 that the use of LLMs carries potential risks and
674 ethical considerations, and we have taken steps
675 to mitigate these risks throughout our research.
676 For example, using open weights LLMs and open
677 sourcing our software helps the

¹<https://www.beren.io/2023-03-19-LLMs-confabulate-not-hallucinate/>

678 community understand the concept of risk mitigation
679 and reproducibility in experiments.

680 We recognize that the use of LLMs can have
681 environmental and social impacts. We have made
682 efforts to minimize the environmental impact of our
683 research by optimizing our code and using energy-
684 efficient hardware. By using AWQ quantization,
685 we allowed the models to run on one A100 GPU.
686 The estimated working hours of the single GPU
687 was of 20 hours, for a CO2 emission of 2.8 kg. In
688 our commitment to offset these emissions, we have
689 initiated the establishment of a forest through
690 [Tree-dome](https://www.treedom.net/)². As an initial endeavor, we have
691 planted a tree, uniquely identified with the code
692 [YMZ-6K66](#).

692 References

- 693 Gabor Angeli, Michael Johnson Premkumar, and
694 Christopher D Manning. 2015. Leveraging linguistic
695 structure for open domain information extraction. In
696 *Proceedings of the 53rd Annual Meeting of the Association
697 for Computational Linguistics and the 7th International
698 Joint Conference on Natural Language Processing (Volume 1: Long
699 Papers)*, pages 344–354.
- 700 Lukas Berglund, Meg Tong, Max Kaufmann, Mikita
701 Balesni, Asa Cooper Stickland, Tomasz Korbak, and
702 Owain Evans. 2023. [The reversal curse: LLMs trained on "a is b" fail to learn "b is a"](#). In *NeurIPS
703 Workshop on Attributing Model Behavior at Scale*.
704
- 705 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
706 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
707 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
708 Askell, et al. 2020. Language models are few-shot
709 learners. *Advances in neural information processing
710 systems*, 33:1877–1901.
- 711 Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,
712 Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi,
713 Cunxiang Wang, Yidong Wang, et al. 2023. A survey
714 on evaluation of large language models. *ACM
715 Transactions on Intelligent Systems and Technology*.
- 716 Muhao Chen, Yaxin Tian, Mohan Yang, Puneet Mathur,
717 and Yaxin Chen. 2021. Knowledge graph embedding:
718 A survey of approaches and applications. *ACM
719 Computing Surveys (CSUR)*, 54(5):1–35.
- 720 François Chollet. 2019. On the measure of intelligence.
721 *arXiv preprint arXiv:1911.01547*.
- 722 Roi Cohen, Mor Geva, Jonathan Berant, and Amir
723 Globerson. 2023. [Crawling the internal knowledge-base of language models](#). In *Findings of the Association
724 for Computational Linguistics: EACL 2023*,
725 pages 1856–1869, Dubrovnik, Croatia. Association
726 for Computational Linguistics.
727

²<https://www.treedom.net/>

728	Anthony Fader, Stephen Soderland, and Oren Etzioni. 2011. Identifying relations for open information extraction. In <i>Proceedings of the Conference on Empirical Methods in Natural Language Processing</i> , pages 1535–1545. Association for Computational Linguistics.	783
729		784
730		785
731		786
732		
733		
734	Christiane Fellbaum. 2020. <i>WordNet</i> . Princeton University Press.	787
735		788
736	Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. <i>Psychological bulletin</i> , 76(5):378.	789
737		790
738		791
739	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. <i>arXiv preprint arXiv:2312.10997</i> .	792
740		793
741		
742		794
743		795
744	Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujia Yang, Nan Duan, and Weizhu Chen. 2023. CritiC: Large language models can self-correct with tool-interactive critiquing .	796
745		797
746		798
747		799
748	Francesca Grasso, Vladimiro Lopera Rulfi, and Luigi Di Caro. 2022. Multialignet: Cross-lingual knowledge bridges between words and senses. In <i>International Conference on Knowledge Engineering and Knowledge Management</i> , pages 36–50. Springer.	
749		
750		
751		
752		
753	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. In <i>International Conference on Learning Representations</i> .	
754		
755		
756		
757		
758	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset . In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks</i> , volume 1. Curran.	
759		
760		
761		
762		
763		
764	Shima Imani, Liang Du, and Harsh Shrivastava. 2023. Mathprompter: Mathematical reasoning using large language models. In <i>ICLR 2023 Workshop on Trustworthy and Reliable Large-Scale Machine Learning Models</i> .	
765		
766		
767		
768		
769	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. <i>ACM Computing Surveys</i> , 55(12):1–38.	
770		
771		
772		
773		
774	Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In <i>Proceedings of the 40th International Conference on Machine Learning, ICML’23</i> . JMLR.org.	
775		
776		
777		
778		
779	Sungdong Kim, Sanghwan Bae, Jamin Shin, Soyoun Kang, Donghyun Kwak, Kang Min Yoo, and Minjoon Seo. 2023. Aligning large language models through synthetic feedback. In <i>EMNLP</i> .	
780		
781		
782		
	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. <i>ArXiv</i> , abs/2205.11916.	
	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the 29th Symposium on Operating Systems Principles</i> , pages 611–626.	
	Valentina Leone, Giovanni Siragusa, Luigi Di Caro, and Roberto Navigli. 2020. Building semantic grams of human knowledge . In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 2991–3000, Marseille, France. European Language Resources Association.	
	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in Neural Information Processing Systems</i> , 33:9459–9474.	
	Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Xingyu Dang, and Song Han. 2023. Awq: Activation-aware weight quantization for llm compression and acceleration. <i>arXiv preprint arXiv:2306.00978</i> .	
	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252.	
	Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What makes good in-context examples for gpt-3? <i>arXiv preprint arXiv:2101.06804</i> .	
	Dabin Min, Kaeun Kim, Jong Hyuk Lee, Yisak Kim, and Chang Min Park. 2022. RRED : A radiology report error detector based on deep learning framework . In <i>Proceedings of the 4th Clinical Natural Language Processing Workshop</i> , pages 41–52, Seattle, WA. Association for Computational Linguistics.	
	Roberto Navigli and Simone Paolo Ponzetto. 2010. Babelnet: Building a very large multilingual semantic network. In <i>Proceedings of the 48th annual meeting of the association for computational linguistics</i> , pages 216–225.	
	Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. 2024. Unifying large language models and knowledge graphs: A roadmap. <i>IEEE Transactions on Knowledge and Data Engineering</i> .	
	Joon Sung Park, Lindsay Popowski, Carrie Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2022. Social simulacra: Creating populated	

838	prototypes for social computing systems. In <i>Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology</i> , pages 1–18.	
839		
840		
841	Ellie Pavlick. 2023. Symbols and grounding in large language models. <i>Philosophical Transactions of the Royal Society A</i> , 381(2251):20220041.	
842		
843		
844	Erika Petersen and Christopher Potts. 2023. Lexical semantics with large language models: A case study of English “break” . In <i>Findings of the Association for Computational Linguistics: EACL 2023</i> , pages 490–511, Dubrovnik, Croatia. Association for Computational Linguistics.	
845		
846		
847		
848		
849		
850	Fabio Petroni, Tim Rocktäschel, Mike Lewis, and Sebastian Riedel. 2019. Language models as knowledge bases? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> .	
851		
852		
853		
854	Justus J Randolph. 2005. Free-marginal multirater kappa (multirater k [free]): An alternative to fleiss’ fixed-marginal multirater kappa. <i>Online submission</i> .	
855		
856		
857	Simon Razniewski and Gerhard Weikum. 2021. Zero-shot open knowledge base question answering. In <i>Proceedings of the 30th ACM International Conference on Information and Knowledge Management</i> , pages 1095–1104.	
858		
859		
860		
861		
862	RoyalLab. 2023. royallab/Aetheria-L2-70B · Hugging Face — huggingface.co. https://huggingface.co/royallab/Aetheria-L2-70B . [Accessed 01-02-2024].	
863		
864		
865		
866	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. <i>Communications of the ACM</i> , 64(9):99–106.	
867		
868		
869		
870	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In <i>International Conference on Machine Learning</i> , pages 31210–31227. PMLR.	
871		
872		
873		
874		
875		
876	Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 31.	
877		
878		
879		
880	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	
881		
882		
883		
884		
885		
886	Gladys Tyen, Hassan Mansoor, Peter Chen, Tony Mak, and Victor Cărbune. 2023. Llms cannot find reasoning errors, but can correct them! <i>arXiv preprint arXiv:2311.08516</i> .	
887		
888		
889		
890	Tomer Ullman. 2023. Large language models fail on trivial alterations to theory-of-mind tasks. <i>arXiv preprint arXiv:2302.08399</i> .	
891		
892		
	Laura Ventrice and Giovanni Siragusa. 2023. Enhancing semantic resources via large language models. <i>GENERAL’23: GENerative, Explainable and Reasonable Artificial Learning Workshop 2023</i> .	893 894 895 896
	Duyu Vo and Ebrahim Bagheri. 2017. Extracting structured data from natural language with semi-structured knowledge. <i>arXiv preprint arXiv:1710.10723</i> .	897 898 899
	Peter West, Chandra Bhagavatula, Jack Hessel, Jena Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. 2022. Symbolic knowledge distillation: from general language models to commonsense models . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4602–4625, Seattle, United States. Association for Computational Linguistics.	900 901 902 903 904 905 906 907 908 909
	Zhuohan Xie, Trevor Cohn, and Jey Han Lau. 2023. The next chapter: A study of large language models in storytelling . In <i>International Conference on Natural Language Generation</i> .	910 911 912 913
	Zhewei Yao, Cheng Li, Xiaoxia Wu, Stephen Youn, and Yuxiong He. 2023. A comprehensive study on post-training quantization for large language models. <i>arXiv preprint arXiv:2303.08302</i> .	914 915 916 917
	Wenhao Yu, Meng Jiang, Zhiting Hu, Qingyun Wang, Heng Ji, and Nazneen Rajani. 2021. Knowledge-enriched natural language generation. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts</i> , pages 11–16.	918 919 920 921 922 923
	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4791–4800.	924 925 926 927 928
	Yiming Zhang, Shi Feng, and Chenhao Tan. 2022. Active example selection for in-context learning . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9134–9148, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	929 930 931 932 933 934
	Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In <i>International Conference on Machine Learning</i> , pages 12697–12706. PMLR.	935 936 937 938 939