

On the Role of Entity and Event Level Conceptualization in Generalizable Reasoning: A Survey of Tasks, Methods, Applications, and Future Directions

Anonymous ACL submission

Abstract

Conceptualization, a fundamental element of human cognition, plays a pivotal role in human generalizable reasoning. Generally speaking, it refers to the process of sequentially abstracting specific instances into higher-level concepts and then forming abstract knowledge that can be applied in unfamiliar or novel situations. This enhances models' inferential capabilities and supports the effective transfer of knowledge across various domains. Despite its significance, the broad nature of this term has led to inconsistencies in understanding conceptualization across various works, as there exists different types of instances that can be abstracted in a wide variety of ways. There is also a lack of a systematic overview that comprehensively examines existing works on the definition, execution, and application of conceptualization to enhance reasoning tasks. In this paper, we address these gaps by first proposing a categorization of different types of conceptualizations into four levels based on the types of instances being conceptualized, in order to clarify the term and define the scope of our work. Then, we present the first comprehensive survey of over 150 papers, surveying various definitions, resources, methods, and downstream applications related to conceptualization into a unified taxonomy, with a focus on the entity and event levels. Furthermore, we shed light on potential future directions in this field and hope to garner more attention from the community.

1 Introduction

Conceptualization has been widely recognized as a fundamental component of human intelligence, spanning fields from psychology (Kahneman, 2011; Evans, 2003; Bransford and Franks, 1971) to computational linguistics (Bengio et al., 2021; Tenenbaum et al., 2011; Lachmy et al., 2022). In the era of deep learning, numerous studies have emerged focusing on conceptualization as a means to achieve generalizable reasoning with (Large)

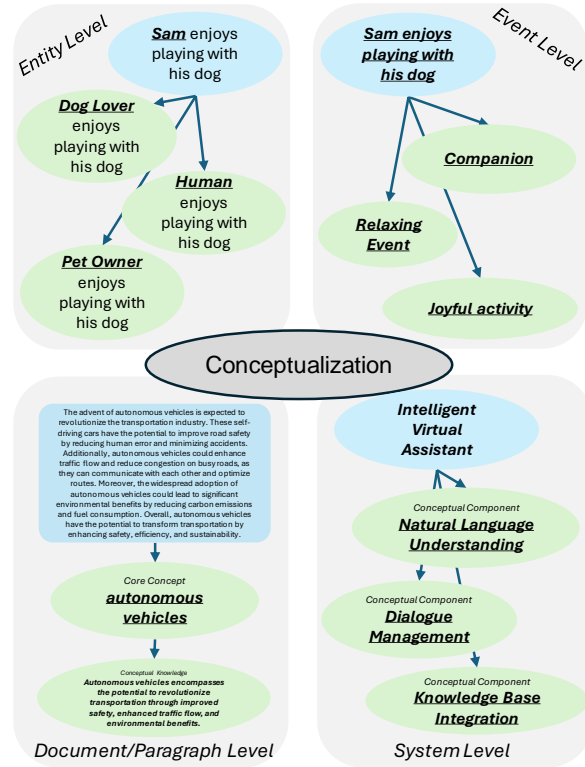


Figure 1: Examples of performing conceptualization at different semantic levels.

Language Models (LLMs; OpenAI, 2022, 2023; Touvron et al., 2023a,b; Mesnard et al., 2024; Reid et al., 2024) in areas such as commonsense reasoning (Wang et al., 2023b,a, 2024a), causal reasoning (Feder et al., 2021; Kunda et al., 1990), physical reasoning (Bisk et al., 2020; Wang et al., 2023c; Hong et al., 2021), and more.

In general terms, conceptualization refers to the process of consolidating specific instances with shared properties or characteristics into a cohesive concept that represents a vast collection of instances. It is a sub-type of abstraction (Giunchiglia and Walsh, 1992), but specifically requires the presence of a concept as the base for such abstraction. With proper conceptualization, abstract knowledge can be subsequently derived by associating original knowledge at the instance level with that concept.

When encountering unfamiliar or novel scenarios, concepts in abstract knowledge can be instantiated to new instances to support downstream reasoning (Tenenbaum et al., 2011). This process can occur at various levels, including entity (Wu et al., 2012; Liang et al., 2017; Alukaev et al., 2023; Liu et al., 2023c), event (He et al., 2024; Wang et al., 2024a,c), paragraph/document (Falke and Gurevych, 2019; Falke et al., 2017), and system levels (Subramonian et al., 2023; Kadioglu and Kleynhans, 2024), ultimately forming a hierarchy that contribute to a comprehensive understanding and representation of knowledge.

Despite its significance, the field lacks a comprehensive and unified taxonomy to categorize existing research on conceptualization. On the one hand, the term “conceptualization” is inherently broad, encompassing various types of conceptualizations across different instances and performed in various ways, all included under a single term. As illustrated in Figure 1, the conceptualization of entities and documents requires two distinct paradigms; however, the current terminology fails to adequately address these differences. This has led to confusion and miscommunication among works that apply conceptualization in their methodologies. On the other hand, the methods for conceptualizing different types of instances in a scalable and accurate manner remain unclear. Finally, it is essential to summarize the benefits that conceptualization can bring to downstream tasks to gather insights for future applications and new research directions.

To address these issues, we present the first-ever survey that systematically taxonomizes conceptualization. First, in Section 2, we define four types of conceptualization based on different semantic levels of the instances being conceptualized: entity, event, document, and system. In later sections, we focus on two main types of conceptualization based on the entity and event levels, as they are largely uncharted in existing literature and play a key role in human reasoning. We then propose a set of objectives to select and survey papers that feature conceptualization as their core idea, review more than 150 papers, and organize them into three main categories, as shown in Figure 2. We summarize the main representative tasks and datasets available for these types of conceptualization in Section 3. Subsequently, in Section 4, we categorize conceptualization acquisition methods into extraction, retrieval, and generative-based methods.

The downstream benefits of conceptualization are discussed in Section 5, with a specific focus on several reasoning tasks. Finally, in Section 6, we propose two future directions that can benefit from conceptualization. We hope our work can serve as a practical handbook for researchers and pave the way for further advancements in the field of conceptualization.

2 Four Levels of Conceptualization

We first define four levels of conceptualization according to the type of instances being conceptualized. They are categorized into four levels: entity level, event level, document level, and system level. Running examples are shown in Figure 1.

Entity Level: Entity-level conceptualization involves grouping multiple entities under a shared concept (Yang et al., 2021; Peng et al., 2022). It is the most common form of conceptualization in human cognition and is frequently applied for knowledge acquisition (Carey, 1991; Murphy, 2004). For instance, entities like “apple,” “pear,” and “grape,” can be categorized together under the broader concept of “fruit.” By doing so, abstract knowledge can be derived by reintegrating the concept into the context of specific instances, such as the assertion “fruit is delicious,” with “apple is delicious” serving as the specific source. When someone encounters an unknown fruit, they can quickly understand its properties by associating it with the abstract knowledge of fruit, such as its possible taste or nutrition.

Event Level: While a concept can capture the semantic meaning of a group of entities, it can also represent events at a higher level of conceptualization. Event-level conceptualization aims to broaden the scope from entities to include events as well (He et al., 2024; Wang et al., 2024c). It seeks to associate different events under a shared concept that preserves the original semantic meaning to the maximum extent possible. For instance, activities like “Sam playing with his dog,” “Alex dancing in the club,” and “Bob doing yoga” can all be conceptualized as “relaxing events.” Abstract knowledge can then follow, stating that “If someone engages in relaxing events, they feel happy and relaxed.” When someone encounters an unknown or unfamiliar event, such as “Charlie likes painting the sunset,” they can infer that painting the sunset is a relaxing event and that Charlie feels happy and relaxed when doing so.

Document Level: Document-level conceptualiza-

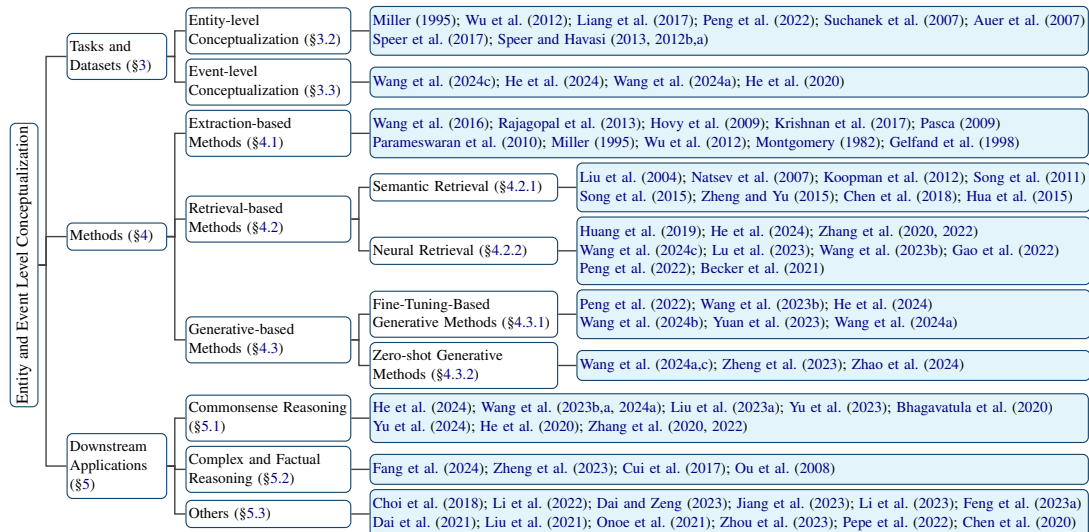


Figure 2: Taxonomy of representative works in entity and event level **conceptualization** categorized by tasks and datasets (§3), methods in performing conceptualization (§4), and downstream applications (§5).

tion further extends the scope of the instance from entities and events to paragraphs or even entire documents. It aims to generate a summary that captures the main ideas and essential information while maintaining the overall semantic and context of the original text. Previous works on abstractive summarization (Ladhak et al., 2022; Wang et al., 2019) have identical objectives, and earlier surveys by Rennard et al. (2023); Lin and Ng (2019); Liu et al. (2024) have effectively summarized these studies. Therefore, we only mention it here to clarify document-level conceptualization for readers and will not go into further detail in later sections to avoid overlap.

System Level: Finally, system-level conceptualization aims to simplify the understanding of a complex system by abstracting its behavior and functionality into a higher-level representation. It is derived from the design of operating systems (Doane et al., 1990) and is under-studied in the domain of NLP. The only representative example is a recent work by Subramonian et al. (2023), where the authors provide a systematic categorization of NLP tasks based on their objectives and characteristics while neglecting the detailed format of input/output and the datasets on which the tasks are evaluated. Due to the limited number of works available, we will not survey this type of conceptualization.

In later sections, we focus specifically on entity and event-level conceptualizations and propose a taxonomy to categorize works into three categories. To ensure that our search for papers is comprehensive and objective in relation to our target scope, we propose the following three objectives for selecting the most relevant papers. First, we aim

for papers that adhere to the paradigm of linking different instances together and use concepts as representations of the formed clusters. We also seek papers that aim to establish hierarchies between different entities and events. Finally, we look for papers that directly seek abstractions of entities or events via concepts. Our proposed taxonomy is primarily categorized into resources, methods, and downstream applications of conceptualization, as this is the most straightforward structure for readers to grasp the topic.

3 Tasks and Datasets

We first survey available datasets and benchmarks, as well as their associated tasks, for these two types of conceptualizations. Statistical comparisons between different resources are shown in Table 1. For datasets that also serve as evaluation benchmarks, we mark their associated tasks with classification task (CLS) and generation task (GEN).

3.1 Concept Linking Task

Most conceptualizations can be formulated as a concept linking task, where the goal is to link an instance i to a concept c such that i can be semantically represented by c . It is challenging due to the infinite number of possible instance-concept pairs. Previous approaches, such as those by Brauer et al. (2010); Yates et al. (2015), have attempted to further restrict the task to linking instances to a limited set of strict ontologies using heuristic or statistical methods. The task can also be formulated with a generative objective, which requires a model to generate c directly given i as input.

Type	Dataset	#Instance	#Concept.	Tasks
<i>Entity</i>	WordNet	82,115	84,428	N/A
	Probase	10,378,743	16,285,393	N/A
	Probase+	10,378,743	21,332,357	N/A
	YAGO	143,210	352,297	N/A
	DBPedia	1,000,000	1,000,000	N/A
	ConceptNet	21,000,000	8,000,000	N/A
	COPEN	24,000	393	CLS
<i>Event</i>	Abs.ATM.	21,493	503,588	CLS, GEN
	Abs.Pyr.	17,000	220,797	CLS, GEN
	CANDLE	21,442	6,181,391	N/A

Table 1: Statistical comparisons between different datasets with entity and event level conceptualizations.

3.2 Entity-level Conceptualization Datasets

To conceptualize different entities into concepts, multiple large-scale concept taxonomies have been constructed as resources for this type of conceptualization. WordNet (Miller, 1995) is the first and most well-known concept taxonomy, which is a large lexical database of English. It is a network of concepts, where each concept is a set of synonyms. Probase (Wu et al., 2012; Liang et al., 2017) is a later built concept taxonomy, which is a large-scale probabilistic taxonomy of concepts. It is constructed by analyzing a large amount of web pages and search logs. YAGO (Suchanek et al., 2007) is a semantic knowledge base, which is a large-scale concept taxonomy of entities and events. It is constructed by extracting information from Wikipedia (Merity et al., 2017) and WordNet. DBPedia (Auer et al., 2007) is a large-scale knowledge base which is built by extracting structured information from Wikipedia. It also contains structured conceptual knowledge about entities and events. ConceptNet (Speer et al., 2017) is the most recent concept taxonomy, featuring a large-scale semantic network of concepts. It is constructed by extracting structured information from various sources, including Wikipedia, WordNet, and Open Mind Common Sense (Singh et al., 2002). Recently, Peng et al. (2022) introduced COPEN, an entity level conceptualization benchmark that is constructed by probing language models to retrieve concepts of an entity from a pre-defined set of concepts. All of them are important knowledge bases that are rich in entity conceptualizations.

3.3 Event-level Conceptualization Datasets

Compared to abstracting entities, there are fewer resources available for event-level conceptualizations. The most notable is the AbstractATOMIC dataset (He et al., 2024), which was constructed

by filtering head events from the ATOMIC dataset and identifying instance candidates within each event using syntactic parsing and human-defined rules. These instances are matched against Probase and WordNet to acquire candidate concepts using GlossBERT (Huang et al., 2019), which are then verified by a supervised model and human annotations. AbsPyramid (Wang et al., 2024c) extends the AbstractATOMIC pipeline to ASER (Zhang et al., 2020, 2022), a large-scale eventuality knowledge graph, by incorporating candidate concepts generated by ChatGPT to complement Probase and WordNet. It also extends coverage to verbs in addition to nouns and events, and broadens the domain of events from social aspects to all aspects. Both datasets provide rich event conceptualizations sourced from diverse origins.

4 Conceptualization Acquisition Methods

Next, we survey methods for performing or collecting entity and event-level conceptualizations. We categorize them into three paradigms: extraction, retrieval, and generative-based methods, which are briefly demonstrated in Figure 3. We provide more discussions in Appendix A.

4.1 Extraction-Based Methods

Extracting concepts from text is the earliest paradigm for systematically collecting conceptualizations (Montgomery, 1982; Gelfand et al., 1998). It typically involves first extracting all possible concepts from the text, followed with identifying the relationships between these concepts. In this process, concepts are recognized either by looking for the most frequent words or by matching against a predefined list of patterns, such as “is a,” “is a type of”, etc. Instances are then matched by looking for the subject of these patterns in the text, which forms instance-conceptualization pairs. The main advantages of extraction-based methods (Wang et al., 2016; Parameswaran et al., 2010; Rajagopal et al., 2013; Hovy et al., 2009; Krishnan et al., 2017; Pasca, 2009) are easy implementation, high processing speed, and free of training data. This has facilitated the development of many large-scale concept taxonomies and knowledge bases, such as WordNet (Miller, 1995), ConceptNet (Speer et al., 2017; Speer and Havasi, 2013, 2012b,a), Probase (Wu et al., 2012; Liang et al., 2017), and DBpedia (Auer et al., 2007; Bizer et al., 2009). However, these methods, while successful in ex-

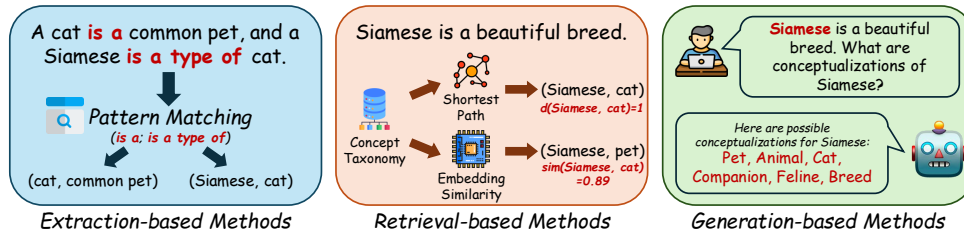


Figure 3: Conceptual demonstration of different types of methods in performing or collecting entity and event level conceptualizations. Instance and conceptualization pairs can be obtained at the end of each type of method.

tracting conceptual relationships from text, are limited by text quality, reliance on predefined concepts, lack of semantic understanding, difficulty handling ambiguous words, and poor generalization to new domains or unseen concepts.

4.2 Retrieval-Based Methods

4.2.1 Semantic-Based Retrieval

Semantic-based retrieval methods aim to obtain conceptualizations by looking at the semantic similarity between the input instance and the concepts in a pre-defined concept taxonomy. It typically involves representing both the instance and a set of concepts into a shared semantic space and calculating the similarity between them. One representative approach is to use WordNet (Miller, 1995), a large lexical database of English words, to calculate semantic similarity between two words as their shortest path in the WordNet hierarchy (Liu et al., 2004). Other methods (Natsev et al., 2007; Song et al., 2011, 2015; Koopman et al., 2012; Zheng and Yu, 2015; Chen et al., 2018; Hua et al., 2015) also share similar aspirations and define their own way of calculating such similarities. However, these methods are usually limited by the need for comprehensive and accurate knowledge bases, high computational costs, the inability to handle unseen concepts, and the loss of important semantic context, prompting the development of neural-based retrieval methods.

4.2.2 Neural-Based Retrieval

Neural-based retrieval methods overcome previous limitations by leveraging neural networks (or language models) to learn the semantic representations of the input instance and the concepts in the knowledge base or concept taxonomy. Then, the similarity between the input instance and the concepts can be calculated based on the learned representation embeddings. This approach can be benefitted by the advancement in language modeling, such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), and DeBERTa (He et al., 2021, 2023). The most representative work in neural-based con-

cept retrieval is AbstractATOMIC (He et al., 2024). It uses GlossBERT (Huang et al., 2019) to encode concepts (from WordNet and Probase) and instances (extracted from events in ATOMIC (Sap et al., 2019)) into embeddings and leverage cosine similarity and human annotations to collect conceptualizations in a large scale manner. Other methods (Wang et al., 2024c; Zhang et al., 2020, 2022; Lu et al., 2023; Wang et al., 2023b; Gao et al., 2022; Becker et al., 2021) similarly adopt different strategies in leveraging LMs as encoders, expanding the coverage of instances, training retrieval models. Despite their promising results, these methods are limited by their need for extensive labeled data, reliance on the completeness and accuracy of the knowledge base, and inability to retrieve new concepts that are out of training data.

4.3 Generative-Based Methods

4.3.1 Fine-Tuning-Based Generative Methods

Fine-tuning-based generative methods aim to take an entity or event as input and generate the concept directly via a fine-tuned generative language model. This approach allows the model to generate conceptualizations for new instances and offers maximum flexibility of the input. Several methods (Peng et al., 2022; Yuan et al., 2023; He et al., 2024; Wang et al., 2024c,b, 2023b) have adopted this paradigm in training generative conceptualizers, based on models such as GPT2 (Radford et al., 2019), BART (Lewis et al., 2020), and T5 (Raffel et al., 2020), for automated conceptualization acquisition. These methods typically train LMs on human-annotated or pre-existing conceptualization resources and yield outstanding results. However, fine-tuning-based generative methods are limited by their high computational cost, time-consuming and resource-intensive data collection, uncertain performance across diverse domains, and relatively low quality of novel concepts compared to human annotations. While these are common limitations associated with fine-tuned generative models, zero-shot generative methods using powerful LLMs and

advanced prompting techniques potentially address these issues.

4.3.2 Zero-Shot Generative Methods

Finally, zero-shot generative-based methods leverage powerful LLMs (Brown et al., 2020; OpenAI, 2022, 2023; Reid et al., 2024; Touvron et al., 2023a,b) to generate the concept directly from an input instance. They rely on the vast amount of internal knowledge within the model and human-crafted prompts to efficiently distill conceptualizations and abstract knowledge from the models. This is particularly useful when training data is scarce or when the domain is new and there are no existing training data available. Existing methods (Wang et al., 2024a,c; Zheng et al., 2023; Zhao et al., 2024) all share similar aspirations in collecting conceptualizations. The benefits are significant, as these methods can collect conceptualizations efficiently and at low cost without specific fine-tuning. The resulting conceptualization knowledge base are thus scalable and downstream models trained on them typically have improved generalization ability to new instances and domains. However, to ensure high-quality generated conceptualizations, it is recommended to implement quality control mechanisms such as human evaluation or discriminators as post-filters. Recent studies (Wang et al., 2024a; Fang et al., 2024) have shown that commonsense plausibility estimators (Liu et al., 2023b) are effective for such quality control.

5 Downstream Applications

We then survey downstream tasks that can benefit from applying conceptualizations to provide readers with a general picture of what can be achieved and how to benefit from integrating conceptualizations. An overview of performances by different methods that leverage conceptualization, evaluated on various benchmarks, are shown in Figure 4.

5.1 Commonsense Reasoning

Commonsense reasoning is the ability to make inferences about the world based on common knowledge, which involves reasoning about everyday events and situations (Davis, 1990; Davis and Marcus, 2015). In this section, we discuss how conceptualizations benefit models in performing commonsense reasoning tasks.

Generative Commonsense Inference Modeling:

The task of generative commonsense inference

modeling (COMET; (Bosselut et al., 2019; Hwang et al., 2021)) aims to complete an inferential commonsense knowledge given a head event and a commonsense relation. State-of-the-art methods for COMET mainly fine-tune language models on large-scale commonsense knowledge bases, which suffer from data sparsity and lack of diversity in commonsense knowledge. Although transfer from LLMs helps (West et al., 2022, 2023), distilled knowledge tends to be too easy for models to learn and converge to trivial inferences. To address these issues, Wang et al. (2023b) proposed to leverage conceptualization as knowledge augmentation tools to improve COMET. Conceptualizations are first derived from head events to obtain abstracted events. Then, the tail of the original commonsense knowledge is placed back to the abstracted event to form abstracted commonsense knowledge. These derived abstract knowledge are then integrated with the original knowledge in commonsense knowledge bases to enrich the diversity of commonsense knowledge. Experiments show consistent improvement in models’ performances. Wang et al. (2024a) further show that, by instantiating conceptualizations in abstract knowledge back to other novel instances, models can be further improved by training with newly instantiated knowledge. Liu et al. (2023a) also proposed a task that aims to generate diverse sentences describing concept relationships in various everyday scenarios. Conceptualizations and associated abstract knowledge can further boost models’ performances on this task.

Commonsense Question Answering: The task of commonsense question answering aims to answer questions that require commonsense knowledge. Various benchmarks and datasets have been proposed to evaluate LMs’ performances, such as Abductive NLI (aNLI; (Bhagavatula et al., 2020)), CommonsenseQA (CSQA; (Talmor et al., 2019)), PhysicalQA (PIQA; (Bisk et al., 2020)), SocialQA (SIQA; (Sap et al., 2019)), and Winogrande (WG; (Sakaguchi et al., 2021)). To obtain a generalizable model for commonsense question answering, the most effective pipeline fine-tunes language models on QA pairs synthesized from knowledge in commonsense knowledge bases (Ma et al., 2021; Shi et al., 2023; Wang et al., 2023a). The head h_o and relation r of a (h_o, r, t) triple are transformed into a question using natural language prompts, with the tail t serving as the correct an-

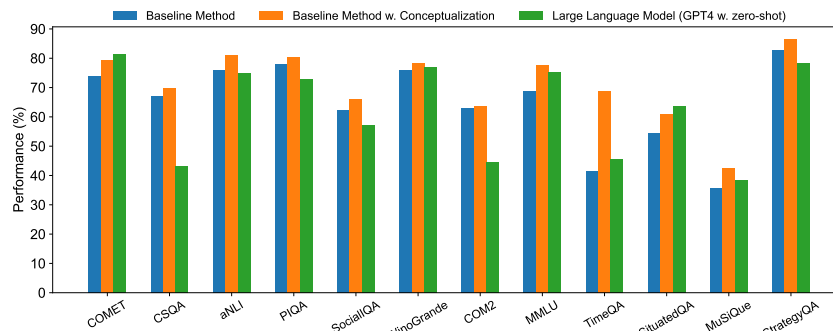


Figure 4: Empirical benefits of conceptualization in methods across various benchmarks compared to baselines.

499 swer option. Distractors or negative examples are
500 generated by randomly sampling tails from triples
501 that do not share common keywords with the head.
502 To leverage conceptualization into the QA synthesis
503 process, Wang et al. (2023a); Fang et al. (2024)
504 have proposed two strategies: On the one hand,
505 they improve distractor sampling by incorporat-
506 ing conceptualizations of head events into com-
507 mon words of the question, thereby enabling selec-
508 tion of more relevant distractors that improve the
509 model’s ability to discern correct answers from dis-
510 tractors. On the other hand, abstract knowledge de-
511 rived from head events are integrated into original
512 synthesized QA pairs, akin to COMET, to enrich
513 the training data with diverse information, thereby
514 enhancing the model’s generalization capability in
515 commonsense question answering tasks. Experimen-
516 tal results show that the proposed strategies
517 significantly improve the performance of common-
518 sense question answering with conceptualization.

519 **5.2 Complex and Factual Reasoning**

520 Complex reasoning refers to the ability to solve
521 intricate problems that necessitate multiple steps
522 of reasoning, which involves reasoning upon intri-
523 cate scenarios, which may encompass multiple en-
524 tities, events, and relations. Fang et al. (2024) pro-
525 posed to synthesize complex queries based on com-
526 monsense knowledge triples from ATOMIC. Both
527 human-defined rules and tails generated by large
528 language models are utilized to generate these com-
529 plex queries. The model is subsequently trained
530 on these complex queries to enhance its capabil-
531 ity to solve complex reasoning problems. In this
532 context, conceptualizations of head events can be
533 used as augmentations to generate more diverse
534 and complex queries (Cui et al., 2017). This can
535 assist the model in learning to solve more intricate
536 problems. Simultaneously, conceptualizations of
537 head events can also be used to generate more in-
538 formative distractors. This can aid the model in

learning to distinguish more effectively between
correct answers and distractors.

541 Zheng et al. (2023) also developed a prompting
542 method to improve the performance of LLMs on
543 general and factual QA tasks. It involves instruct-
544 ing the model with a simple zero-shot prompt to
545 consider each question abstractly by generating and
546 probing relevant concepts, then using this knowl-
547 edge in the prompt to generate the answer. This
548 simple prompting method has been shown to signif-
549 icantly improve the performance of large language
550 models on general QA tasks, including MMLU
551 (Physics and Chemistry) (Hendrycks et al., 2021),
552 TimeQA (Chen et al., 2021), StrategyQA (Geva
553 et al., 2021), and MuSiQue (Trivedi et al., 2022).
554 This work is interesting as it demonstrates that a
555 simple prompting method can significantly enhance
556 the performance of LLMs on general QA tasks.

557 **5.3 Others**

558 Aside from those two types of tasks, the line of
559 works focusing on ultra-fine entity (Choi et al.,
560 2018; Li et al., 2022; Dai and Zeng, 2023; Jiang
561 et al., 2023; Li et al., 2023; Feng et al., 2023a; Dai
562 et al., 2021; Liu et al., 2021; Onoe et al., 2021)
563 and event typing (Zhou et al., 2023; Pepe et al.,
564 2022; Chen et al., 2020) can also be benefited by
565 conceptualization. These tasks aim to type named
566 entities, nominal nouns, and pronouns into a set of
567 free-form phrases. Conceptualizations can serve as
568 a bridge between the surface form and the target
569 type, which is crucial for these tasks.

570 **6 Future Directions and Conclusions**

571 Finally, we conclude our work by discussing two
572 interesting future directions.

573 **6.1 Controllable Generation**

574 Firstly, we envision that conceptualization can as-
575 sist controllable text generation (Feng et al., 2023b;
576 Huang et al., 2023; Zhang et al., 2024). In some

formulations, the task requires the model to generate a brief piece of text that remains consistent within a specific context or scope (Meng et al., 2022). Conceptualizations can be applied as additional supervision signals or constraints that guide the model to generate text whose conceptualizations align with those in the input theme, thereby enhancing the controllability of the generated text. This could be achieved by training a pair of conceptualization generator and discriminator, which could be used to generate the conceptualizations and evaluate their consistency between input and output text. Conceptualization can also serve as data augmentation tools to provide more training data, preferably guided with human annotation or large language models as loose teachers, for training more robust text generators that better align with the controllable targeting data.

Similarly, it may also benefit hallucination reduction (Choubey et al., 2023; Dale et al., 2023; Ji et al., 2023b; Sun et al., 2023). Hallucination (Ji et al., 2023a) refers to generating text that is unsupported by the input context, such as introducing information that is not present in the context or even contradicts it. In many reasoning scenarios, hallucination can be detrimental to the model’s performance, and neutralizing it is crucial for ensuring the reliability of the generated text. Towards this objective, conceptualization can be similarly applied as external signals to verify the generated text and ensure its accuracy. By measuring the semantic distance of conceptualizations between the given input and generated contents, hallucinations can possibly be detected by finding clearly unrelated concepts appearing at both ends. Empirical metrics to measure such distance can be the shortest path length of concepts in taxonomies such as WordNet (Miller, 1995) and Probase (Wu et al., 2012), or even embedding similarity between different concepts. However, it’s important to build a comprehensive set of conceptualizations of a given text to support such a verification process, as incomplete conceptualizations may cause erroneously detected hallucinations due to human-caused errors. We leave detailed implementations to future work.

6.2 Modeling Changes in Distribution

Conceptualization also plays a pivotal role in building reasoning systems that can capture situational changes in distribution to achieve System II reasoning (Sloman, 1996; Kahneman, 2011). Among the several components that make up System II rea-

soning, a key element is the ability to reason with situational changes in distribution (Bengio et al., 2021, 2019). These changes are triggered by environmental factors and actions by the agents themselves or others, especially when dealing with non-stationarities (Bengio, 2017). This ability can be achieved by dynamically recombining existing concepts in the given environment or action and learning from the resultant situational changes (Lake and Baroni, 2018; Bahdanau et al., 2019; de Vries et al., 2019). For instance, consider the event “PersonX is driving a car on a sunny day.” A change in the weather from sunny to rainy could cause a different outcome, such as “PersonX becomes more cautious and drives slower.” This illustrates that a change in weather conditions can lead to a change in the driver’s behavior, representing an environmental change that triggers situational changes within the distribution of different weather conditions. In this process, the model is required to infer different changes that can possibly occur within a single event as the context, and reason about the potential outcome of each change. To model the distribution of different changes within an event, conceptualization can be used to represent the different states of the environment or action (Wang and Song, 2024). The model can then reason about the changes in distribution by manipulating the granularity of conceptualized changes. This type of distributional conceptualization not only provides an ontology for modeling the distribution of different changes within an event, but also assists the model in reasoning about the potential outcomes with appropriate abstract knowledge. Future works can leverage LLMs to curate benchmark datasets via sequential conceptualization generation and develop advanced systems for System II reasoning.

6.3 Conclusions

In conclusion, this work surveys conceptualizations by proposing a four-level hierarchical definition and reviewing representative works in acquiring, leveraging, and applying entity and event-level conceptualization to downstream reasoning tasks. We also propose several intriguing ideas related to conceptualizations that may inspire further research. We hope our work paves the way for more research works toward generalizable machine intelligence through conceptualization and fosters the development of more advanced systems that can capture, organize, and learn world knowledge through connection between concepts, much like humans do.

679 Limitations

680 The main limitations of our survey are two-fold.
681 First, due to the vast amount of literature on con-
682 ceptualization and conceptual knowledge across
683 various datasets, we only cover the most represen-
684 tative works that stand out for their exceptional
685 value and uniqueness in our taxonomy. Most of the
686 papers are sourced from ACL Anthology¹, ACM
687 Digital Library², and proceedings of leading arti-
688 ficial intelligence and machine learning conferences.
689 Consequently, it is possible that some other related
690 works are not included, but we aim to cover them
691 in future versions. Second, our survey specifically
692 focuses on entity and event level conceptualiza-
693 tion, leaving document/paragraph level and system
694 level conceptualization unaddressed. However, it is
695 impossible to survey everything within one single
696 submission. Future research can expand the scope
697 of our survey to include more types of conceptual-
698 izations and modalities, such as categorization in
699 the vision modality (Chen and Wang, 2004).

700 Ethics Statement

701 Our paper presents a comprehensive survey of con-
702 ceptualization, with a specific focus on entity and
703 event levels. All datasets and models reviewed in
704 this survey are properly cited and are available un-
705 der free-access licenses for research purposes. We
706 did not conduct additional dataset curation or hu-
707 man annotation work. Therefore, to the best of our
708 knowledge, this paper does not yield any ethical
709 concerns.

710 References

711 Danis Alukaev, Semen Kiselev, Ilya Pershin, Bulat
712 Ibragimov, Vladimir Ivanov, Alexey Kornae, and
713 Ivan Titov. 2023. [Cross-modal conceptualization in
714 bottleneck models](#). In *Proceedings of the 2023 Con-
715 ference on Empirical Methods in Natural Language
716 Processing, EMNLP 2023, Singapore, December 6-
717 10, 2023*, pages 5241–5253. Association for Compu-
718 tational Linguistics.

719 Sören Auer, Christian Bizer, Georgi Kobilarov, Jens
720 Lehmann, Richard Cyganiak, and Zachary G. Ives.
721 2007. [Dbpedia: A nucleus for a web of open data](#). In
722 *The Semantic Web, 6th International Semantic Web
723 Conference, 2nd Asian Semantic Web Conference,
724 ISWC 2007 + ASWC 2007, Busan, Korea, Novem-
725 ber 11-15, 2007*, volume 4825 of *Lecture Notes in
726 Computer Science*, pages 722–735. Springer.

¹<https://aclanthology.org/>

²<https://dl.acm.org/>

Dzmitry Bahdanau, Shikhar Murty, Michael 727
Noukhovitch, Thien Huu Nguyen, Harm de Vries, 728
and Aaron C. Courville. 2019. [Systematic gener- 729
alization: What is required and can it be learned?](#) 730
In *7th International Conference on Learning 731
Representations, ICLR 2019, New Orleans, LA, USA, 732
May 6-9, 2019*. OpenReview.net. 733

Jiaxin Bai, Xin Liu, Weiqi Wang, Chen Luo, and 734
Yangqiu Song. 2023. [Complex query answering on 735
eventuality knowledge graph with implicit logical 736
constraints](#). In *Advances in Neural Information Pro- 737
cessing Systems 36: Annual Conference on Neural 738
Information Processing Systems 2023, NeurIPS 2023, 739
New Orleans, LA, USA, December 10 - 16, 2023*. 740

Maria Becker, Katharina Korfhage, and Anette Frank. 741
2021. [COCO-EX: A tool for linking concepts from 742
texts to conceptnet](#). In *Proceedings of the 16th Con- 743
ference of the European Chapter of the Association 744
for Computational Linguistics: System Demonstra- 745
tions, EACL 2021, Online, April 19-23, 2021*, pages 746
119–126. Association for Computational Linguistics. 747

Yoshua Bengio. 2017. [The consciousness prior](#). *CoRR*, 748
abs/1709.08568. 749

Yoshua Bengio, Yann LeCun, and Geoffrey E. Hin- 750
ton. 2021. [Deep learning for AI](#). *Commun. ACM*, 751
64(7):58–65. 752

Yoshua Bengio et al. 2019. From system 1 deep learning 753
to system 2 deep learning. In *Neural Information 754
Processing Systems*. 755

Chandra Bhagavatula, Ronan Le Bras, Chaitanya 756
Malaviya, Keisuke Sakaguchi, Ari Holtzman, Han- 757
nah Rashkin, Doug Downey, Wen-tau Yih, and Yejin 758
Choi. 2020. [Abductive commonsense reasoning](#). In 759
*8th International Conference on Learning Represen- 760
tations, ICLR 2020, Addis Ababa, Ethiopia, April 761
26-30, 2020*. OpenReview.net. 762

Ning Bian, Xianpei Han, Le Sun, Hongyu Lin, Yaojie 763
Lu, and Ben He. 2023. [Chatgpt is a knowledgeable 764
but inexperienced solver: An investigation of com- 765
monsense problem in large language models](#). *CoRR*, 766
abs/2303.16421. 767

Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng 768
Gao, and Yejin Choi. 2020. [PIQA: reasoning about 769
physical commonsense in natural language](#). In *The 770
Thirty-Fourth AAAI Conference on Artificial Intelli- 771
gence, AAAI 2020, The Thirty-Second Innovative Ap- 772
plications of Artificial Intelligence Conference, IAAI 773
2020, The Tenth AAAI Symposium on Educational 774
Advances in Artificial Intelligence, EAAI 2020, New 775
York, NY, USA, February 7-12, 2020*, pages 7432– 776
7439. AAAI Press. 777

Christian Bizer, Jens Lehmann, Georgi Kobilarov, Sören 778
Auer, Christian Becker, Richard Cyganiak, and Se- 779
bastian Hellmann. 2009. [Dbpedia - A crystallization 780
point for the web of data](#). *J. Web Semant.*, 7(3):154– 781
165. 782

783	Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. COMET: commonsense transformers for automatic knowledge graph construction . In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers</i> , pages 4762–4779. Association for Computational Linguistics.	840
784		841
785		842
786		843
787		844
788		845
789		846
790		
791		
792	John D Bransford and Jeffery J Franks. 1971. The abstraction of linguistic ideas. <i>Cognitive psychology</i> , 2(4):331–350.	
793		
794		
795	Falk Brauer, Michael Huber, Gregor Hackenbroich, Ulf Leser, Felix Naumann, and Wojciech M. Barczynski. 2010. Graph-based concept identification and disambiguation for enterprise search . In <i>Proceedings of the 19th International Conference on World Wide Web, WWW 2010, Raleigh, North Carolina, USA, April 26-30, 2010</i> , pages 171–180. ACM.	
796		
797		
798		
799		
800		
801		
802	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners . In <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual</i> .	
803		
804		
805		
806		
807		
808		
809		
810		
811		
812		
813		
814		
815		
816		
817	Susan Carey. 1991. Knowledge acquisition: Enrichment or conceptual change. <i>The epigenesis of mind: Essays on biology and cognition</i> , pages 257–291.	
818		
819		
820	Chunkit Chan, Jiayang Cheng, Weiqi Wang, Yuxin Jiang, Tianqing Fang, Xin Liu, and Yangqiu Song. 2024a. Exploring the potential of chatgpt on sentence level relations: A focus on temporal, causal, and discourse relations . In <i>Findings of the Association for Computational Linguistics: EACL 2024, St. Julian's, Malta, March 17-22, 2024</i> , pages 684–721. Association for Computational Linguistics.	
821		
822		
823		
824		
825		
826		
827		
828	Chunkit Chan, Cheng Jiayang, Yauwai Yim, Zheyue Deng, Wei Fan, Haoran Li, Xin Liu, Hongming Zhang, Weiqi Wang, and Yangqiu Song. 2024b. Negotiationtom: A benchmark for stress-testing machine theory of mind on negotiation surrounding . <i>CoRR</i> , abs/2404.13627.	
829		
830		
831		
832		
833		
834	Lihan Chen, Jiaqing Liang, Chenhao Xie, and Yanghua Xiao. 2018. Short text entity linking with fine-grained topics . In <i>Proceedings of the 27th ACM International Conference on Information and Knowledge Management, CIKM 2018, Torino, Italy, October 22-26, 2018</i> , pages 457–466. ACM.	
835		
836		
837		
838		
839		
	Muhao Chen, Hongming Zhang, Haoyu Wang, and Dan Roth. 2020. What are you trying to do? semantic typing of event processes . In <i>Proceedings of the 24th Conference on Computational Natural Language Learning, CoNLL 2020, Online, November 19-20, 2020</i> , pages 531–542. Association for Computational Linguistics.	840
		841
		842
		843
		844
		845
		846
	Wenhu Chen, Xinyi Wang, and William Yang Wang. 2021. A dataset for answering time-sensitive questions . In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual</i> .	847
		848
		849
		850
		851
		852
	Yihan Chen, Benfeng Xu, Quan Wang, Yi Liu, and Zhendong Mao. 2024. Benchmarking large language models on controllable generation under diversified instructions . In <i>Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada</i> , pages 17808–17816. AAAI Press.	853
		854
		855
		856
		857
		858
		859
		860
		861
		862
	Yixin Chen and James Ze Wang. 2004. Image categorization by learning and reasoning with regions . <i>J. Mach. Learn. Res.</i> , 5:913–939.	863
		864
		865
	Jiayang Cheng, Lin Qiu, Tsz Ho Chan, Tianqing Fang, Weiqi Wang, Chunkit Chan, Dongyu Ru, Qipeng Guo, Hongming Zhang, Yangqiu Song, Yue Zhang, and Zheng Zhang. 2023. Storyanalogy: Deriving story-level analogies from large language models to unlock analogical understanding . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 11518–11537. Association for Computational Linguistics.	866
		867
		868
		869
		870
		871
		872
		873
		874
		875
	Eunsol Choi, Omer Levy, Yejin Choi, and Luke Zettlemoyer. 2018. Ultra-fine entity typing . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers</i> , pages 87–96. Association for Computational Linguistics.	876
		877
		878
		879
		880
		881
		882
	Prafulla Kumar Choubey, Alexander R. Fabbri, Jesse Vig, Chien-Sheng Wu, Wenhao Liu, and Nazneen Rajani. 2023. Cape: Contrastive parameter ensembling for reducing hallucination in abstractive summarization . In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 10755–10773. Association for Computational Linguistics.	883
		884
		885
		886
		887
		888
		889
		890
	Wanyun Cui, Yanghua Xiao, Haixun Wang, Yangqiu Song, Seung-won Hwang, and Wei Wang. 2017. KBQA: learning question answering over QA corpora and knowledge bases . <i>Proc. VLDB Endow.</i> , 10(5):565–576.	891
		892
		893
		894
		895
	Hongliang Dai, Yangqiu Song, and Haixun Wang. 2021. Ultra-fine entity typing with weak supervision from a	896
		897

898	masked language model. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021</i> , pages 1790–1799. Association for Computational Linguistics.	Stephanie M. Doane, James W. Pellegrino, and Roberta L. Klatzky. 1990. Expertise in a computer operating system: Conceptualization and performance. <i>Hum. Comput. Interact.</i> , 5(2-3):267–304.	956 957 958 959
905	Hongliang Dai and Ziqian Zeng. 2023. From ultra-fine to fine: Fine-tuning ultra-fine entity typing models to fine-grained. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 2259–2270. Association for Computational Linguistics.	Jonathan St BT Evans. 2003. In two minds: dual-process accounts of reasoning. <i>Trends in cognitive sciences</i> , 7(10):454–459.	960 961 962
912	Dhairya Dalal, Paul Buitelaar, and Mihael Arcan. 2023. Calm-bench: A multi-task benchmark for evaluating causality-aware language models. In <i>Findings of the Association for Computational Linguistics: EACL 2023, Dubrovnik, Croatia, May 2-6, 2023</i> , pages 296–311. Association for Computational Linguistics.	Tobias Falke and Iryna Gurevych. 2019. Fast concept mention grouping for concept map-based multi-document summarization. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)</i> , pages 695–700. Association for Computational Linguistics.	963 964 965 966 967 968 969 970 971
918	David Dale, Elena Voita, Loïc Barrault, and Marta R. Costa-jussà. 2023. Detecting and mitigating hallucinations in machine translation: Model internal workings alone do well, sentence similarity even better. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 36–50. Association for Computational Linguistics.	Tobias Falke, Christian M. Meyer, and Iryna Gurevych. 2017. Concept-map-based multi-document summarization using concept coreference resolution and global importance optimization. In <i>Proceedings of the Eighth International Joint Conference on Natural Language Processing, IJCNLP 2017, Taipei, Taiwan, November 27 - December 1, 2017 - Volume 1: Long Papers</i> , pages 801–811. Asian Federation of Natural Language Processing.	972 973 974 975 976 977 978 979 980
927	Ernest Davis. 1990. <i>Representations of commonsense knowledge</i> . notThenot Morgan Kaufmann series in representation and reasoning. Morgan Kaufmann.	Tianqing Fang, Zeming Chen, Yangqiu Song, and Antoine Bosselut. 2024. Complex reasoning over logical queries on commonsense knowledge graphs. <i>CoRR</i> , abs/2403.07398.	981 982 983 984
930	Ernest Davis and Gary Marcus. 2015. Commonsense reasoning and commonsense knowledge in artificial intelligence. <i>Commun. ACM</i> , 58(9):92–103.	Tianqing Fang, Weiqi Wang, Sehyun Choi, Shibo Hao, Hongming Zhang, Yangqiu Song, and Bin He. 2021a. Benchmarking commonsense knowledge base population with an effective evaluation dataset. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021</i> , pages 8949–8964. Association for Computational Linguistics.	985 986 987 988 989 990 991 992 993
933	Harm de Vries, Dzmitry Bahdanau, Shikhar Murty, Aaron C. Courville, and Philippe Beaudoin. 2019. CLOSURE: assessing systematic generalization of CLEVR models. In <i>Visually Grounded Interaction and Language (ViGIL), NeurIPS 2019 Workshop, Vancouver, Canada, December 13, 2019</i> .	Tianqing Fang, Hongming Zhang, Weiqi Wang, Yangqiu Song, and Bin He. 2021b. DISCOS: bridging the gap between discourse knowledge and commonsense knowledge. In <i>WWW '21: The Web Conference 2021, Virtual Event / Ljubljana, Slovenia, April 19-23, 2021</i> , pages 2648–2659. ACM / IW3C2.	994 995 996 997 998 999
939	Zheyang Deng, Weiqi Wang, Zhaowei Wang, Xin Liu, and Yangqiu Song. 2023. Gold: A global and local-aware denoising framework for commonsense knowledge graph noise detection. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 3591–3608. Association for Computational Linguistics.	Amir Feder, Nadav Oved, Uri Shalit, and Roi Reichart. 2021. Causalm: Causal model explanation through counterfactual language models. <i>Comput. Linguistics</i> , 47(2):333–386.	1000 1001 1002 1003
946	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)</i> , pages 4171–4186. Association for Computational Linguistics.	Yanlin Feng, Adithya Pratapa, and David R. Mortensen. 2023a. Calibrated seq2seq models for efficient and generalizable ultra-fine entity typing. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 15550–15560. Association for Computational Linguistics.	1004 1005 1006 1007 1008 1009 1010
955		Yuxi Feng, Xiaoyuan Yi, Xiting Wang, Laks V. S. Lakshmanan, and Xing Xie. 2023b. Dunst: Dual noisy self	1011 1012

1013	training for semi-supervised controllable text generation . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 8760–8785. Association for Computational Linguistics.	<i>Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual</i> , pages 17427–17440.	1069 1070
1014			
1015			
1016			
1017			
1018			
1019	Silin Gao, Jena D. Hwang, Saya Kanno, Hiromi Wakaki, Yuki Mitsufuji, and Antoine Bosselut. 2022. Comfact: A benchmark for linking contextual commonsense knowledge . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022</i> , pages 1656–1675. Association for Computational Linguistics.	Eduard H. Hovy, Zornitsa Kozareva, and Ellen Riloff. 2009. Toward completeness in concept extraction and classification . In <i>Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing, EMNLP 2009, 6-7 August 2009, Singapore, A meeting of SIGDAT, a Special Interest Group of the ACL</i> , pages 948–957. ACL.	1071 1072 1073 1074 1075 1076 1077
1020			
1021			
1022			
1023			
1024			
1025			
1026			
1027	Boris Gelfand, Marilyn Wulfecker, and WF Punch. 1998. Automated concept extraction from plain text. In <i>AAAI 1998 Workshop on Text Categorization</i> , pages 13–17.	Wen Hua, Zhongyuan Wang, Haixun Wang, Kai Zheng, and Xiaofang Zhou. 2015. Short text understanding through lexical-semantic analysis . In <i>31st IEEE International Conference on Data Engineering, ICDE 2015, Seoul, South Korea, April 13-17, 2015</i> , pages 495–506. IEEE Computer Society.	1078 1079 1080 1081 1082 1083
1028			
1029			
1030			
1031	Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? A question answering benchmark with implicit reasoning strategies . <i>Trans. Assoc. Comput. Linguistics</i> , 9:346–361.	Luyao Huang, Chi Sun, Xipeng Qiu, and Xuanjing Huang. 2019. Glossbert: BERT for word sense disambiguation with gloss knowledge . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 3507–3512. Association for Computational Linguistics.	1084 1085 1086 1087 1088 1089 1090 1091 1092
1032			
1033			
1034			
1035			
1036	Fausto Giunchiglia and Toby Walsh. 1992. A theory of abstraction . <i>Artif. Intell.</i> , 57(2-3):323–389.	Xuancheng Huang, Zijun Liu, Peng Li, Tao Li, Maosong Sun, and Yang Liu. 2023. An extensible plug-and-play method for multi-aspect controllable text generation . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 15233–15256. Association for Computational Linguistics.	1093 1094 1095 1096 1097 1098 1099 1100
1037			
1038	Mutian He, Tianqing Fang, Weiqi Wang, and Yangqiu Song. 2024. Acquiring and modeling abstract commonsense knowledge via conceptualization . <i>Artificial Intelligence</i> , page 104149.	Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. 2021. (comet-) atomic 2020: On symbolic and neural commonsense knowledge graphs . In <i>Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021</i> , pages 6384–6392. AAAI Press.	1101 1102 1103 1104 1105 1106 1107 1108 1109 1110 1111
1039			
1040			
1041			
1042	Mutian He, Yangqiu Song, Kun Xu, and Dong Yu. 2020. On the role of conceptualization in commonsense knowledge graph construction . <i>CoRR</i> , abs/2003.03239.	Raghav Jain, Daivik Sojitra, Arkadeep Acharya, Sriparna Saha, Adam Jatowt, and Sandipan Dandapat. 2023. Do language models have a common sense regarding time? revisiting temporal commonsense reasoning in the era of large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 6750–6774. Association for Computational Linguistics.	1112 1113 1114 1115 1116 1117 1118 1119 1120
1043			
1044			
1045			
1046	Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2023. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023a. Survey of hallucination in natural language generation . <i>ACM Comput. Surv.</i> , 55(12):248:1–248:38.	1121 1122 1123 1124 1125
1047			
1048			
1049			
1050			
1051			
1052	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: decoding-enhanced bert with disentangled attention . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.		
1053			
1054			
1055			
1056			
1057			
1058	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding . In <i>9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021</i> . OpenReview.net.		
1059			
1060			
1061			
1062			
1063			
1064	Yining Hong, Li Yi, Josh Tenenbaum, Antonio Torralba, and Chuang Gan. 2021. PTR: A benchmark for part-based conceptual, relational, and physical reasoning . In <i>Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information</i>		
1065			
1066			
1067			
1068			

1240	2023, Toronto, Canada, July 9-14, 2023, pages 4719–	Knowledge-driven data construction for zero-shot	1297
1241	4731. Association for Computational Linguistics.	evaluation in commonsense question answering. In	1298
1242	Jiacheng Liu, Wenya Wang, Dianzhuo Wang, Noah A.	<i>Thirty-Fifth AAAI Conference on Artificial Intelli-</i>	1299
1243	Smith, Yejin Choi, and Hannaneh Hajishirzi. 2023b.	<i>gence, AAAI 2021, Thirty-Third Conference on In-</i>	1300
1244	Vera: A general-purpose plausibility estimation	<i>novative Applications of Artificial Intelligence, IAAI</i>	1301
1245	model for commonsense statements . In <i>Proceedings</i>	<i>2021, The Eleventh Symposium on Educational Ad-</i>	1302
1246	<i>of the 2023 Conference on Empirical Methods in Natu-</i>	<i>vanced in Artificial Intelligence, EAAI 2021, Vir-</i>	1303
1247	<i>ral Language Processing, EMNLP 2023, Singapore,</i>	<i>tual Event, February 2-9, 2021, pages 13507–13515.</i>	1304
1248	<i>December 6-10, 2023, pages 1264–1287. Association</i>	AAAI Press.	1305
1249	for Computational Linguistics.		
1250	Qing Liu, Hongyu Lin, Xinyan Xiao, Xianpei Han,	Joshua Maynez, Priyanka Agrawal, and Sebastian	1306
1251	Le Sun, and Hua Wu. 2021. Fine-grained entity	Gehrmann. 2023. Benchmarking large language	1307
1252	typing via label reasoning . In <i>Proceedings of the</i>	model capabilities for conditional generation . In	1308
1253	<i>2021 Conference on Empirical Methods in Natural</i>	<i>Proceedings of the 61st Annual Meeting of the As-</i>	1309
1254	<i>Language Processing, EMNLP 2021, Virtual Event</i>	<i>sociation for Computational Linguistics (Volume 1:</i>	1310
1255	<i>/ Punta Cana, Dominican Republic, 7-11 November,</i>	<i>Long Papers), ACL 2023, Toronto, Canada, July 9-14,</i>	1311
1256	<i>2021, pages 4611–4622. Association for Computa-</i>	<i>2023, pages 9194–9213. Association for Computa-</i>	1312
1257	<i>tional Linguistics.</i>	<i>tional Linguistics.</i>	1313
1258	Ran Liu, Ming Liu, Min Yu, He Zhang, Jianguo Jiang,	Tao Meng, Sidi Lu, Nanyun Peng, and Kai-Wei Chang.	1314
1259	Gang Li, and Weiqing Huang. 2024. Sumsurvey:	2022. Controllable text generation with neurally-	1315
1260	An abstractive dataset of scientific survey papers for	decomposed oracle . In <i>Advances in Neural Infor-</i>	1316
1261	long document summarization . In <i>Findings of the As-</i>	<i>mation Processing Systems 35: Annual Conference</i>	1317
1262	<i>sociation for Computational Linguistics, ACL 2024,</i>	<i>on Neural Information Processing Systems 2022,</i>	1318
1263	<i>Bangkok, Thailand and virtual meeting, August 11-</i>	<i>NeurIPS 2022, New Orleans, LA, USA, November 28</i>	1319
1264	<i>16, 2024, pages 9632–9651. Association for Computa-</i>	<i>- December 9, 2022.</i>	1320
1265	<i>tional Linguistics.</i>		
1266	Shuang Liu, Fang Liu, Clement T. Yu, and Weiyi Meng.	Stephen Merity, Caiming Xiong, James Bradbury, and	1321
1267	2004. An effective approach to document retrieval	Richard Socher. 2017. Pointer sentinel mixture mod-	1322
1268	via utilizing wordnet and recognizing phrases . In	els . In <i>5th International Conference on Learning</i>	1323
1269	<i>SIGIR 2004: Proceedings of the 27th Annual Inter-</i>	<i>Representations, ICLR 2017, Toulon, France, April</i>	1324
1270	<i>national ACM SIGIR Conference on Research and</i>	<i>24-26, 2017, Conference Track Proceedings. Open-</i>	1325
1271	<i>Development in Information Retrieval, Sheffield, UK,</i>	<i>Review.net.</i>	1326
1272	<i>July 25-29, 2004, pages 266–272. ACM.</i>		
1273	Yihong Liu, Haotian Ye, Leonie Weissweiler, Philipp	Thomas Mesnard, Cassidy Hardin, Robert Dadashi,	1327
1274	Wicke, Renhao Pei, Robert Zangenfeind, and Hin-	Surya Bhupatiraju, Shreya Pathak, Laurent Sifre,	1328
1275	rich Schütze. 2023c. A crosslingual investigation of	Morgane Rivi�ere, Mihir Sanjay Kale, Juliette Love,	1329
1276	conceptualization in 1335 languages . In <i>Proceed-</i>	Pouya Tafti, L�eonard Hussenot, Aakanksha Chowdh-	1330
1277	<i>ings of the 61st Annual Meeting of the Association</i>	ery, Adam Roberts, Aditya Barua, Alex Botev, Alex	1331
1278	<i>for Computational Linguistics (Volume 1: Long Pa-</i>	Castro-Ros, Ambrose Slone, Am�elie H�eliou, Andrea	1332
1279	<i>pers), ACL 2023, Toronto, Canada, July 9-14, 2023,</i>	Tacchetti, Anna Bulanova, Antonia Paterson, Beth	1333
1280	<i>pages 12969–13000. Association for Computational</i>	Tsai, Bobak Shahriari, Charline Le Lan, Christo-	1334
1281	<i>Linguistics.</i>	pher A. Choquette-Choo, Cl�ement Crepy, Daniel Cer,	1335
1282	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	Daphne Ippolito, David Reid, Elena Buchatskaya,	1336
1283	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	Eric Ni, Eric Noland, Geng Yan, George Tucker,	1337
1284	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	George-Christian Muraru, Grigory Rozhdestvenskiy,	1338
1285	Roberta: A robustly optimized BERT pretraining	Henryk Michalewski, Ian Tenney, Ivan Grishchenko,	1339
1286	approach . <i>CoRR</i> , abs/1907.11692.	Jacob Austin, James Keeling, Jane Labanowski,	1340
1287	Mengying Lu, Yuquan Wang, Jifan Yu, Yexing Du,	Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan,	1341
1288	Lei Hou, and Juanzi Li. 2023. Distantly supervised	Jeremy Chen, Johan Ferret, Justin Chiu, and et al.	1342
1289	course concept extraction in moocs with academic	2024. Gemma: Open models based on gemini re-	1343
1290	discipline . In <i>Proceedings of the 61st Annual Meet-</i>	search and technology . <i>CoRR</i> , abs/2403.08295.	1344
1291	<i>ing of the Association for Computational Linguis-</i>	George A. Miller. 1995. Wordnet: A lexical database	1345
1292	<i>tics (Volume 1: Long Papers), ACL 2023, Toronto,</i>	for english . <i>Commun. ACM</i> , 38(11):39–41.	1346
1293	<i>Canada, July 9-14, 2023, pages 13044–13059. Asso-</i>	Christine A Montgomery. 1982. Concept extrac-	1347
1294	<i>ciation for Computational Linguistics.</i>	tion . <i>American journal of computational linguistics</i> ,	1348
1295	Kaixin Ma, Filip Ilievski, Jonathan Francis, Yonatan	8(2):70–73.	1349
1296	Bisk, Eric Nyberg, and Alessandro Oltramari. 2021.	Gregory Murphy. 2004. <i>The big book of concepts</i> . MIT	1350
		press.	1351
		Apostol Natsev, Alexander Haubold, Jelena Tesic, Lex-	1352
		ing Xie, and Rong Yan. 2007. Semantic concept-	1353
		based query expansion and re-ranking for multimedia	1354

1355	retrieval. In <i>Proceedings of the 15th International Conference on Multimedia 2007, Augsburg, Germany, September 24-29, 2007</i> , pages 991–1000. ACM.	1412
1356		1413
1357		1414
1358	Yasumasa Onoe, Michael Boratko, Andrew McCallum, and Greg Durrett. 2021. Modeling fine-grained entity types with box embeddings . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021</i> , pages 2051–2064. Association for Computational Linguistics.	1415
1359		
1360		
1361		
1362		
1363		
1364		
1365		
1366		
1367	OpenAI. 2022. Chatgpt: Optimizing language models for dialogue . <i>OpenAI</i> .	
1368		
1369	OpenAI. 2023. GPT-4 technical report . <i>CoRR</i> , abs/2303.08774.	
1370		
1371	Shiyan Ou, Viktor Pekar, Constantin Orasan, Christian Spurk, and Matteo Negri. 2008. Development and alignment of a domain-specific ontology for question answering . In <i>Proceedings of the International Conference on Language Resources and Evaluation, LREC 2008, 26 May - 1 June 2008, Marrakech, Morocco</i> . European Language Resources Association.	
1372		
1373		
1374		
1375		
1376		
1377		
1378	Aditya G. Parameswaran, Hector Garcia-Molina, and Anand Rajaraman. 2010. Towards the web of concepts: Extracting concepts from large datasets . <i>Proc. VLDB Endow.</i> , 3(1):566–577.	
1379		
1380		
1381		
1382	Marius Pasca. 2009. Outclassing wikipedia in open-domain information extraction: Weakly-supervised acquisition of attributes over conceptual hierarchies . In <i>EACL 2009, 12th Conference of the European Chapter of the Association for Computational Linguistics, Proceedings of the Conference, Athens, Greece, March 30 - April 3, 2009</i> , pages 639–647. The Association for Computer Linguistics.	
1383		
1384		
1385		
1386		
1387		
1388		
1389		
1390	Hao Peng, Xiaozhi Wang, Shengding Hu, Hailong Jin, Lei Hou, Juanzi Li, Zhiyuan Liu, and Qun Liu. 2022. COPEN: probing conceptual knowledge in pre-trained language models . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022</i> , pages 5015–5035. Association for Computational Linguistics.	
1391		
1392		
1393		
1394		
1395		
1396		
1397		
1398	Sveva Pepe, Edoardo Barba, Rexhina Blloshmi, and Roberto Navigli. 2022. STEPS: semantic typing of event processes with a sequence-to-sequence approach . In <i>Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022</i> , pages 11156–11164. AAAI Press.	
1399		
1400		
1401		
1402		
1403		
1404		
1405		
1406		
1407		
1408	Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is chatgpt a general-purpose natural language processing task solver? In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 1339–1384. Association for Computational Linguistics.	1412
1409		1413
1410		1414
1411		1415
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners . <i>OpenAI blog</i> , 1(8):9.	1416
		1417
		1418
		1419
	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 . <i>J. Mach. Learn. Res.</i> , 21:140:1–140:67.	1420
		1421
		1422
		1423
		1424
	Dheeraj Rajagopal, Erik Cambria, Daniel Olsher, and Kenneth Kwok. 2013. A graph-based approach to commonsense concept extraction and semantic similarity detection . In <i>22nd International World Wide Web Conference, WWW '13, Rio de Janeiro, Brazil, May 13-17, 2013, Companion Volume</i> , pages 565–570. International World Wide Web Conferences Steering Committee / ACM.	1425
		1426
		1427
		1428
		1429
		1430
		1431
		1432
	Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy P. Lillicrap, Jean-Baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, Ioannis Antonoglou, Rohan Anil, Sebastian Borgeaud, Andrew M. Dai, Katie Millican, Ethan Dyer, Mia Glaese, Thibault Sottiaux, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, James Molloy, Jilin Chen, Michael Isard, Paul Barham, Tom Hennigan, Ross McIlroy, Melvin Johnson, Johan Schalkwyk, Eli Collins, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Clemens Meyer, Gregory Thornton, Zhen Yang, Henryk Michalewski, Zaheer Abbas, Nathan Schucher, Ankesh Anand, Richard Ives, James Keeling, Karel Lenc, Salem Haykal, Siamak Shakeri, Pranav Shyam, Aakanksha Chowdhery, Roman Ring, Stephen Spencer, Eren Sezener, and et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context . <i>CoRR</i> , abs/2403.05530.	1433
		1434
		1435
		1436
		1437
		1438
		1439
		1440
		1441
		1442
		1443
		1444
		1445
		1446
		1447
		1448
		1449
		1450
		1451
		1452
	Virgile Rennard, Guokan Shang, Julie Hunter, and Michalis Vazirgiannis. 2023. Abstractive meeting summarization: A survey . <i>Trans. Assoc. Comput. Linguistics</i> , 11:861–884.	1453
		1454
		1455
		1456
	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: an adversarial winograd schema challenge at scale . <i>Commun. ACM</i> , 64(9):99–106.	1457
		1458
		1459
		1460
	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social iqa: Commonsense reasoning about social interactions . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 4462–4472. Association for Computational Linguistics.	1461
		1462
		1463
		1464
		1465
		1466
		1467
		1468
		1469

1470	Haochen Shi, Weiqi Wang, Tianqing Fang, Baixuan Xu,	Arjun Subramonian, Xingdi Yuan, Hal Daumé III, and	1525
1471	Wenxuan Ding, Xin Liu, and Yangqiu Song. 2023.	Su Lin Blodgett. 2023. It takes two to tango: Navigating conceptualizations of NLP tasks and measurements of performance . In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 3234–3279. Association for Computational Linguistics.	1526
1472	QADYNAMICS: training dynamics-driven synthetic QA diagnostic for zero-shot commonsense question answering . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 15329–15341. Association for Computational Linguistics.		1527
1473			1528
1474			1529
1475			1530
1476			1531
1477			
1478	Push Singh, Thomas Lin, Erik T. Mueller, Grace Lim,	Fabian M. Suchanek, Gjergji Kasneci, and Gerhard	1532
1479	Travell Perkins, and Wan Li Zhu. 2002. Open mind common sense: Knowledge acquisition from the general public . In <i>On the Move to Meaningful Internet Systems, 2002 - DOA/CoopIS/ODBASE 2002 Confederated International Conferences DOA, CoopIS and ODBASE 2002 Irvine, California, USA, October 30 - November 1, 2002, Proceedings</i> , volume 2519 of <i>Lecture Notes in Computer Science</i> , pages 1223–1237. Springer.	Weikum. 2007. Yago: a core of semantic knowledge . In <i>Proceedings of the 16th International Conference on World Wide Web, WWW 2007, Banff, Alberta, Canada, May 8-12, 2007</i> , pages 697–706. ACM.	1533
1480			1534
1481			1535
1482			1536
1483		Bin Sun, Yitong Li, Fei Mi, Fanhu Bie, Yiwei Li,	1537
1484		and Kan Li. 2023. Towards fewer hallucinations in knowledge-grounded dialogue generation via augmentative and contrastive knowledge-dialogue . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 1741–1750. Association for Computational Linguistics.	1538
1485			1539
1486			1540
1487			1541
1488	Steven A Sloman. 1996. The empirical case for two systems of reasoning. <i>Psychological bulletin</i> , 119(1):3.		1542
1489			1543
			1544
			1545
1490	Yangqiu Song, Haixun Wang, Zhongyuan Wang, Hong-	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and	1546
1491	song Li, and Weizhu Chen. 2011. Short text conceptualization using a probabilistic knowledgebase . In <i>IJCAI 2011, Proceedings of the 22nd International Joint Conference on Artificial Intelligence, Barcelona, Catalonia, Spain, July 16-22, 2011</i> , pages 2330–2336. IJCAI/AAAI.	Jonathan Berant. 2019. Commonsenseqa: A question answering challenge targeting commonsense knowledge . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)</i> , pages 4149–4158. Association for Computational Linguistics.	1547
1492			1548
1493			1549
1494			1550
1495			1551
1496			1552
1497	Yangqiu Song, Shusen Wang, and Haixun Wang. 2015. Open domain short text conceptualization: A generative + descriptive modeling approach . In <i>Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015</i> , pages 3820–3826. AAAI Press.		1553
1498			1554
1499			1555
1500		Qingyu Tan, Hwee Tou Ng, and Lidong Bing. 2023. Towards benchmarking and improving the temporal reasoning capability of large language models . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 14820–14835. Association for Computational Linguistics.	1556
1501			1557
1502			1558
1503			1559
1504	Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge . In <i>Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA</i> , pages 4444–4451. AAAI Press.		1560
1505			1561
1506			1562
1507			1563
1508		Joshua B Tenenbaum, Charles Kemp, Thomas L Griffiths, and Noah D Goodman. 2011. How to grow a mind: Statistics, structure, and abstraction. <i>science</i> , 331(6022):1279–1285.	1564
1509			1565
1510	Robyn Speer and Catherine Havasi. 2012a. Conceptnet 5. Tiny Trans. Comput. Sci. , 1.		1566
1511			1567
1512	Robyn Speer and Catherine Havasi. 2012b. Representing general relational knowledge in conceptnet 5 . In <i>Proceedings of the Eighth International Conference on Language Resources and Evaluation, LREC 2012, Istanbul, Turkey, May 23-25, 2012</i> , pages 3679–3686. European Language Resources Association (ELRA).	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models . <i>CoRR</i> , abs/2302.13971.	1568
1513			1569
1514			1570
1515			1571
1516			1572
1517			1573
1518	Robyn Speer and Catherine Havasi. 2013. Conceptnet 5: A large semantic network for relational knowledge . In Iryna Gurevych and Jungi Kim, editors, <i>The People’s Web Meets NLP, Collaboratively Constructed Language Resources</i> , Theory and Applications of Natural Language Processing, pages 161–176. Springer.	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esioibu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	1574
1519			1575
1520			1576
1521			1577
1522			1578
1523			1579
1524			1580
			1581
			1582

1583	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models . <i>CoRR</i> , abs/2307.09288.	
1584		
1585		
1586		
1587		
1588		
1589		
1590		
1591		
1592		
1593		
1594		
1595		
1596		
1597		
1598	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. Musique: Multi-hop questions via single-hop question composition . <i>Trans. Assoc. Comput. Linguistics</i> , 10:539–554.	
1599		
1600		
1601		
1602	Shuting Wang, Alexander Ororbia, Zhaohui Wu, Kyle Williams, Chen Liang, Bart Pursel, and C. Lee Giles. 2016. Using prerequisites to extract concept maps from textbooks . In <i>Proceedings of the 25th ACM International Conference on Information and Knowledge Management, CIKM 2016, Indianapolis, IN, USA, October 24-28, 2016</i> , pages 317–326. ACM.	
1603		
1604		
1605		
1606		
1607		
1608		
1609	WeiQi Wang, Tianqing Fang, Wenxuan Ding, Baixuan Xu, Xin Liu, Yangqiu Song, and Antoine Bosselut. 2023a. CAR: conceptualization-augmented reasoner for zero-shot commonsense question answering . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 13520–13545. Association for Computational Linguistics.	
1610		
1611		
1612		
1613		
1614		
1615		
1616		
1617	WeiQi Wang, Tianqing Fang, Chunyang Li, Haochen Shi, Wenxuan Ding, Baixuan Xu, Zhaowei Wang, Jiaxin Bai, Xin Liu, Jiayang Cheng, Chunkit Chan, and Yangqiu Song. 2024a. CANDLE: iterative conceptualization and instantiation distillation from large language models for commonsense reasoning . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> . Association for Computational Linguistics.	
1618		
1619		
1620		
1621		
1622		
1623		
1624		
1625		
1626		
1627		
1628	WeiQi Wang, Tianqing Fang, Baixuan Xu, Chun Yi Louis Bo, Yangqiu Song, and Lei Chen. 2023b. CAT: A contextualized conceptualization and instantiation framework for commonsense reasoning . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 13111–13140. Association for Computational Linguistics.	
1629		
1630		
1631		
1632		
1633		
1634		
1635		
1636		
1637	WeiQi Wang and Yangqiu Song. 2024. MARS: Benchmarking the metaphysical reasoning abilities of language models with a multi-task evaluation dataset . <i>CoRR</i> , abs/2406.02106.	
1638		
1639		
1640		
	Wenbo Wang, Yang Gao, Heyan Huang, and Yuxiang Zhou. 2019. Concept pointer network for abstractive summarization . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 3074–3083. Association for Computational Linguistics.	1641
		1642
		1643
		1644
		1645
		1646
		1647
		1648
		1649
	Yi Ru Wang, Jiafei Duan, Dieter Fox, and Siddhartha S. Srinivasa. 2023c. NEWTON: are large language models capable of physical reasoning? In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 9743–9758. Association for Computational Linguistics.	1650
		1651
		1652
		1653
		1654
		1655
		1656
	Zhaowei Wang, Wei Fan, Qing Zong, Hongming Zhang, Sehyun Choi, Tianqing Fang, Xin Liu, Yangqiu Song, Ginny Y. Wong, and Simon See. 2024b. Absinstruct: Eliciting abstraction ability from llms through explanation tuning with plausibility estimation . <i>CoRR</i> , abs/2402.10646.	1657
		1658
		1659
		1660
		1661
		1662
	Zhaowei Wang, Haochen Shi, Weiqi Wang, Tianqing Fang, Hongming Zhang, Sehyun Choi, Xin Liu, and Yangqiu Song. 2024c. AbsPyramid: Benchmarking the abstraction ability of language models with a unified entailment graph . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 3991–4010, Mexico City, Mexico. Association for Computational Linguistics.	1663
		1664
		1665
		1666
		1667
		1668
		1669
		1670
	Zihao Wang, Weizhi Fei, Hang Yin, Yangqiu Song, Ginny Y. Wong, and Simon See. 2023d. Wasserstein-fisher-rao embedding: Logical query embeddings with local comparison and global transport . In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 13679–13696. Association for Computational Linguistics.	1671
		1672
		1673
		1674
		1675
		1676
		1677
		1678
	Zihao Wang, Yangqiu Song, Ginny Y. Wong, and Simon See. 2023e. Logical message passing networks with one-hop inference on atomic formulas . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	1679
		1680
		1681
		1682
		1683
		1684
	Zihao Wang, Hang Yin, and Yangqiu Song. 2021. Benchmarking the combinatorial generalizability of complex query answering on knowledge graphs . In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual</i> .	1685
		1686
		1687
		1688
		1689
		1690
		1691
	Peter West, Chandra Bhagavatula, Jack Hessel, Jena D. 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, NAACL 2022, Seattle</i> ,	1692
		1693
		1694
		1695
		1696
		1697
		1698
		1699

1700	WA, United States, July 10-15, 2022, pages 4602–		
1701	4625. Association for Computational Linguistics.		
1702	Peter West, Ronan Le Bras, Taylor Sorensen,		
1703	Bill Yuchen Lin, Liwei Jiang, Ximing Lu, Khyathi		
1704	Chandu, Jack Hessel, Ashutosh Baheti, Chandra		
1705	Bhagavatula, and Yejin Choi. 2023. Novacomnet: Open commonsense foundation models with symbolic knowledge distillation.		
1706	In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 1127–1149. Association for Computational Linguistics.		
1707			
1708			
1709			
1710			
1711	Wentao Wu, Hongsong Li, Haixun Wang, and		
1712	Kenny Qili Zhu. 2012. Probase: a probabilistic taxonomy for text understanding.		
1713	In <i>Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2012, Scottsdale, AZ, USA, May 20-24, 2012</i> , pages 481–492. ACM.		
1714			
1715			
1716			
1717	Shan Yang, Yongfei Zhang, Guanglin Niu, Qinghua		
1718	Zhao, and Shiliang Pu. 2021. Entity concept-enhanced few-shot relation extraction.		
1719	In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 2: Short Papers), Virtual Event, August 1-6, 2021</i> , pages 987–991. Association for Computational Linguistics.		
1720			
1721			
1722			
1723			
1724			
1725			
1726	Andrew Yates, Nazli Goharian, and Ophir Frieder. 2015. Extracting adverse drug reactions from social media.		
1727	In <i>Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA</i> , pages 2460–2467. AAAI Press.		
1728			
1729			
1730			
1731	Changlong Yu, Xin Liu, Jefferson Maia, Yang Li,		
1732	Tianyu Cao, Yifan Gao, Yangqiu Song, Rahul		
1733	Goutam, Haiyang Zhang, Bing Yin, and Zheng Li. 2024. Cosmo: A large-scale e-commerce common sense knowledge generation and serving system at amazon.		
1734	In <i>Companion of the 2024 International Conference on Management of Data, SIGMOD/PODS '24</i> , page 148–160, New York, NY, USA. Association for Computing Machinery.		
1735			
1736			
1737			
1738			
1739			
1740	Changlong Yu, Weiqi Wang, Xin Liu, Jiaxin Bai,		
1741	Yangqiu Song, Zheng Li, Yifan Gao, Tianyu Cao,		
1742	and Bing Yin. 2023. Folkscope: Intention knowledge graph construction for e-commerce commonsense discovery.		
1743	In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 1173–1191. Association for Computational Linguistics.		
1744			
1745			
1746			
1747			
1748	Chenhan Yuan, Qianqian Xie, Jimin Huang, and Sophia		
1749	Ananiadou. 2024. Back to the future: Towards explainable temporal reasoning with large language models.		
1750	In <i>Proceedings of the ACM on Web Conference 2024, WWW 2024, Singapore, May 13-17, 2024</i> , pages 1963–1974. ACM.		
1751			
1752			
1753			
1754	Siyu Yuan, Deqing Yang, Jinxi Liu, Shuyu Tian, Jiaqing		
1755	Liang, Yanghua Xiao, and Rui Xie. 2023. Causality-aware concept extraction based on knowledge-guided		
1756	prompting. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 9255–9272. Association for Computational Linguistics.	1757	
		1758	
		1759	
		1760	
		1761	
		1762	
		1763	
		1764	
		1765	
		1766	
		1767	
		1768	
		1769	
		1770	
		1771	
		1772	
		1773	
		1774	
		1775	
		1776	
		1777	
		1778	
		1779	
		1780	
		1781	
		1782	
		1783	
		1784	
		1785	
		1786	
		1787	
		1788	
		1789	
		1790	
		1791	
		1792	
		1793	
		1794	
		1795	
		1796	
		1797	
		1798	
		1799	
		1800	
		1801	
		1802	

Appendices

A Conceptualization Acquisition Methods

In this appendix, we elaborate further on different methods of acquiring conceptualization and provide detailed explanations of their weaknesses.

A.1 Extraction Based Methods

For methods that follow the concept extraction paradigm, Wang et al. (2016) proposed a framework to optimize both tasks simultaneously, leading to stronger performances even compared to supervised concept extraction methods. Parameswaran et al. (2010) also proposed a market-basket-based solution, which adapts statistical measures of support and confidence to design a concept extraction algorithm that achieved high precision in concept extraction. Rajagopal et al. (2013) proposed a solution to extract concepts from common-sense text, which uncovers many novel pieces of knowledge that cannot be found in the original corpora. Hovy et al. (2009); Krishnan et al. (2017); Pasca (2009) similarly proposed their solutions for large-scale concept extraction for more efficient data mining.

While these methods have been successful in extracting concepts and relationships from text, they have several limitations. First, they are heavily dependent on the quality of the text and the predefined list of concepts. If the text is noisy or contains many irrelevant words, the performance of these methods can degrade significantly, and the resulting extracted concepts may also tend to be noisy. Second, it's important to note that these methods primarily rely on parsing or pattern matching techniques on text and do not capture semantic information from the text. This potentially makes extracted concepts represented as isolated entities without any context or relationships and could result in mis-extraction of concepts or relationships, especially when the text contains ambiguous or polysemous words. For example, the word "bank" can refer to a financial institution, a river bank, or a memory bank, and without proper context, it's difficult to determine the correct meaning of it, thus leading to incorrect concept extraction. A low-performance parser, if wrongly parsing these words, may also lead to noisy results. Lastly, these methods are not able to generalize well to unseen concepts or text patterns that are not present in the predefined list of concepts. This limits their applicability to

new domains or tasks that require the extraction of novel concepts or relationships. For example, to extract concepts from medical or legal domain text, specific patterns or extraction rules need to be designed, which may not be present when extracting normal conversational text.

A.2 Retrieval Based Methods

A.2.1 Semantic-Based Retrieval

To perform semantic-based retrieval, (Natsev et al., 2007) proposed several approaches for semantic concept-based query expansion and re-ranking in multimedia retrieval, achieving consistent performance improvement compared to text retrieval and multimodal retrieval baseline. (Song et al., 2011, 2015) improved text understanding by using a probabilistic knowledge base based on concepts and developed a Bayesian inference mechanism to conceptualize words and short text. Experimental results show significant improvements on text clustering compared to purely statistical methods and methods that use existing knowledge bases. (Koopman et al., 2012) proposed a corpus-driven approach, adapted from LSA, to retrieve medical concepts with semantic similarity measures. (Zheng and Yu, 2015) similarly used topic modeling and key concept retrieval methods to construct queries from electronic health records, which significantly improves the retrieval of tailored online consumer-oriented health education materials.

Although these methods have shown promising results in various domains, they have several limitations. First, the performance of semantic-based retrieval heavily relies on the quality of the knowledge base or concept taxonomy. In other words, it requires the knowledge base to be comprehensive, accurate, hierarchical, and up-to-date. There are very few knowledge bases that meet all these requirements, and constructing such a knowledge base is a non-trivial task. With incomplete knowledge bases, which are common in practice, the performance of semantic-based retrieval methods can be significantly degraded. Second, semantic-based retrieval methods are usually computationally expensive, as they require calculating the similarity between the input instance and all concepts in the knowledge base. This can induce exponentially increasing computational cost as the size of the knowledge base grows. When dealing with large-scale applications, this even becomes infeasible. Though caching and indexing techniques can be

used to speed up the retrieval process, they are not always effective and cannot generalize well when unseen concepts or instances are encountered. Third, semantic-based retrieval methods still do not consider the semantic context of the input instance. A straightforward formulation is that the model treats the input instance as a bag of words and ignores the word order and syntactic structure. This can lead to a loss of important semantic information, especially when the input instance is long and complex. In this case, the semantic similarity between the input instance and the concepts in the knowledge base may not reflect the true semantic relevance.

A.2.2 Neural-Based Retrieval

For neural-based retrieval, aside from He et al. (2024), (Lu et al., 2023) similarly proposes a novel three-stage framework, which leverages the power of pre-trained language models explicitly and implicitly and employs discipline-embedding models with a self-train strategy based on label generation refinement across different domains.

To deal with the large amount of unlabeled data after human annotation, (Wang et al., 2023b) further proposed a semi-supervised method to unlabel the data with a supervised trained conceptualization discriminator. The discriminator is trained to rate the plausibility of unlabeled conceptualization and the model will be further refined by training on a concatenation of labeled and unlabeled data. This results in a significant improvement in the performance of the conceptualization discriminator, thus enhancing the quality of the retrieved concepts.

Despite these promising results in concept retrieval, neural-based retrieval methods have several limitations. First, these methods are usually data-hungry and require a large amount of labeled data for training. This can be a bottleneck in practice, as labeling data is often expensive and time-consuming. Human annotations are usually required to collect such data, and for models to be generalizable across different domains, the labeled data should be diverse and representative. This is even more costly and challenging to obtain. Second, neural-based retrieval methods still rely on the coverage and quality of the knowledge base or concept taxonomy. If the knowledge base is incomplete or inaccurate, the performance of neural-based retrieval methods can be significantly affected. Moreover, they cannot generate new concepts or instances that are not in the knowledge

base, which limits their generalization ability.

A.3 Generative-Based Methods

A.3.1 Fine-Tuning-Based Generative Methods

While most fine-tuning based methods are explicitly discussed in the main body, we explain their limitations here. First, these methods are usually computationally expensive, as they require fine-tuning a large pre-trained language model on a specific dataset. Both the fine-tuning and the training data collection process can be time-consuming and resource-intensive. Extensive crowd-sourcing or human annotations are usually required to collect high-quality training data, which can be costly and challenging to obtain when the domain coverage scales up. Second, the feasibility of fine-tuning-based generative methods on other domains, such as medical or legal text, is still an open question. The performance of these methods heavily relies on the quality and diversity of the training data, and it's not clear how well they can generalize to new domains or tasks as text understanding abilities vary across different domains. For social commonsense, pre-trained language models have shown strong performance possibly due to a large overlap in the training data distribution, but for other domains, the performance is still unclear. Lastly, although existing studies have shown that fine-tuning based generators can deliver novel concepts that are not in the training data, such a ratio is relatively low and the quality of the generated concepts is still not as good as human annotated ones. This is expected as the models are fitted into the distribution of the training data, and it's hard for them to generate concepts that are out of the distribution.

A.3.2 Zero-Shot Generative Methods

Zero-shot generative methods aim to generate the desired output for any task's input without any task-specific fine-tuning. A very representative example of such generative models is the recently popularized LLMs (OpenAI, 2022, 2023; Touvron et al., 2023a,b; Mesnard et al., 2024; Reid et al., 2024). These models have been pre-trained on very large corpora, including those from the web, Wikipedia, books, and more, and have shown strong performance in various natural language processing tasks, including text generation (Maynez et al., 2023; Chen et al., 2024), temporal reasoning (Tan et al., 2023; Yuan et al., 2024), causal reasoning (Chan et al., 2024a; Dalal et al., 2023; Jin et al., 2023), commonsense reasoning (Jain et al., 2023; Bian

2003 et al., 2023; Fang et al., 2021b,a; Deng et al., 2023),
2004 logical reasoning (Wang et al., 2023d,e, 2021; Bai
2005 et al., 2023), and more (Qin et al., 2023; Cheng
2006 et al., 2023; Chan et al., 2024b).

2007 In the context of conceptualization acquisition,
2008 zero-shot generative methods aim to generate con-
2009 ceptualizations for instances without any instance-
2010 conceptualization pairs in the training data. Wang
2011 et al. (2024a) proposed a few-shot knowledge dis-
2012 tillation method to distill conceptualizations and
2013 associated abstract inferential knowledge from a
2014 large language model to a large-scale knowledge
2015 base. Wang et al. (2024c) also proposed acquiring
2016 conceptualizations for entities and events in ASER
2017 by instructing ChatGPT with a few-shot prompt.
2018 They further designed an instruction-tuning based
2019 method to evoke more conceptualizations from
2020 large language models by fine-tuning them with
2021 explanations on how the conceptualization is de-
2022 rived from the instance and their plausible reason-
2023 ing chains (Wang et al., 2024b). Zheng et al. (2023)
2024 proposed a simple prompting technique, inspired
2025 by chain-of-thought reasoning, that enables LLMs
2026 to do conceptualizations to derive high-level con-
2027 cepts and first principles from instances containing
2028 specific details. Zhao et al. (2024) advanced this
2029 idea by proposing to extract predictive high-level
2030 features (concepts) from a large language model’s
2031 hidden layer activations.

2032 The benefits of these methods are twofold. First,
2033 such generation can introduce conceptualizations
2034 at a very low cost, as the models are pre-trained
2035 and do not require any task-specific fine-tuning.
2036 The only burden seems to be deployment and in-
2037 ference cost, which require a large amount of com-
2038 putational resources and time for large-scale gen-
2039 eration. However, compared to all previous fine-
2040 tuning-based methods, zero-shot generative meth-
2041 ods are much more efficient and scalable, as they do
2042 not require any training data or fine-tuning process.
2043 Second, zero-shot generative methods have shown
2044 strong generalization capabilities to new instances
2045 and domains. They can generate conceptualiza-
2046 tions for instances that are not in the training data
2047 and have shown strong performance in various con-
2048 ceptualization acquisition tasks. This is particularly
2049 useful when the training data is scarce or when the
2050 domain is new, and there are no existing training
2051 data available. Since these large language models
2052 are pre-injected with vast amounts of knowledge,
2053 this makes generalization possible.