



ASER: Towards Large-scale Commonsense Knowledge Acquisition via Higher-order Selectional Preference over Eventualities

Hongming Zhang¹, Xin Liu¹, Haojie Pan¹, Haowen Ke, Jiefu Ou, Tianqing Fang, Yangqiu Song

CSE, HKUST, HKSAR, China

Abstract

Commonsense knowledge acquisition and reasoning have long been a core artificial intelligence problem. However, in the past, there has been a lack of scalable methods to collect commonsense knowledge. In this paper, we propose to develop principles for collecting commonsense knowledge based on selectional preference, which is a common phenomenon in human languages that has been shown to be related to semantics. We generalize the definition of selectional preference from one-hop linguistic syntactic relations to higher-order relations over linguistic graphs. Unlike previous commonsense knowledge definition (e.g., ConceptNet), the selectional preference (SP) knowledge only relies on statistical distribution over linguistic graphs, which can be efficiently and accurately acquired from the unlabeled corpus with modern tools, rather than human-defined relations. As a result, acquiring SP knowledge is a much more scalable way of acquiring commonsense knowledge. Following this principle, we develop a large-scale eventuality (a linguistic term covering activity, state, and event)-based knowledge graph ASER, where each eventuality is represented as a dependency graph, and the relation between them is a discourse relation defined in shallow discourse parsing. The higher-order selectional preference over collected linguistic graphs reflects various kinds of commonsense knowledge. For example, dogs are more likely to bark than cats as the eventuality “dog barks” appears 14,998 times in ASER while “cat barks” only appears 6 times. “Be hungry” is more likely to be the reason rather than result of “eat food” as the edge ⟨“be hungry,” Cause, “eat food”⟩ appears in ASER while ⟨“eat food,” Cause, “be hungry”⟩ does not. Moreover, motivated by the observation that humans understand events by abstracting the observed events to a higher level and can thus transferring their knowledge to new events, we propose a conceptualization module on top of the collected knowledge to significantly boost the coverage of ASER. In total, ASER contains 648 million edges between 438 million eventualities. After conceptualization with Probase, a selectional preference based concept-instance relational knowledge base, our concept graph contains 15 million conceptualized eventualities and 224 million edges between them. Detailed analysis is provided to demonstrate its quality. All the collected data, APIs, and tools that can help convert collected SP knowledge into the format of ConceptNet are available at <https://github.com/HKUST-KnowComp/ASER>.

© 2011 Published by Elsevier Ltd.

Keywords: Commonsense Acquisition, Selectional Preference, Eventualities

Email addresses: hzhanga1@cse.ust.hk (Hongming Zhang), xliucr@cse.ust.hk (Xin Liu), hpanad@cse.ust.hk (Haojie Pan), hkeaa@cse.ust.hk (Haowen Ke), jouaa@connect.ust.hk (Jiefu Ou), tfangaa@cse.ust.hk (Tianqing Fang), yqsong@cse.ust.hk (Yangqiu Song)

¹The first three authors make equally important contributions to this work. Detailed contributions are in the end of this paper.

1. Introduction

Knowledge is crucial to understanding natural language. When reading, in addition to linguistic knowledge of the vocabulary and grammar of a language, readers need to have knowledge about the structure of texts, knowledge about the subject, and background or commonsense knowledge about the world in order to comprehend the text. For example, when a user says “I am hungry” to a chatbot at 1:00 pm, the chatbot should be able to understand that the user may want to have lunch rather than breakfast and recommend some nearby restaurants. This requires the chatbot to understand the complex commonsense knowledge about user’s states and potential consequent activities (i.e., being hungry can motivate the user to eat) and the implications of location and time (i.e., compared with breakfast, lunch is more likely to appear at 1:00 pm. Thus the chatbot should recommend some real food rather than just a cup of coffee).

Commonsense reasoning has long been a challenging problem in the artificial intelligence field. As discussed in [1], commonsense knowledge refers to “millions of basic facts and understanding possessed by most people.” Unlike factual knowledge like “London is the capital of UK,” which is always true, commonsense knowledge is often not inevitably true and only reflects a kind of contextual preference. For example, in most cases, rocks are not used for eating, but some birds do eat rocks to digest. Such kind of knowledge is also called factoids [2, 3]. To effectively represent such preference-like commonsense knowledge, selectional preference [4] was proposed, which was traditionally defined on top of single dependency connections (e.g., *nsubj*, *doobj*, and *amod*). Given a word and a dependency relation, humans have preferences for which words are likely to be connected. For instance, when seeing the verb “sing,” it is highly plausible that its object is “a song,” and when seeing the noun “air,” it is highly plausible that its modifier is “fresh.” However, such selectional preference can only represent commonsense inside an event or state (e.g., which event/state is more likely to happen) and cannot represent commonsense between events/states. One such example is discussed by [5] and similar examples are frequently observed in the Winograd Schema Challenge [6]:

- The soldiers fired at the women, and we saw several of them fall.

To resolve the pronoun “them” in the above example, Wilks argued that machines need to access the *partial information* “hurt things tending to fall down,” which can be translated into the following form: (hurt, X) $\xrightarrow{\text{ResultIn}}$ (X, fall).

In history, many efforts have been devoted to acquiring commonsense knowledge in the form of multi-relational factoids. For example, the Cyc project initiated in the 1980s [7] and ConceptNet (originated from Open Mind Common Sense, OMCS) initiated in 2002 [1], tried to use experts or ordinary people to annotate commonsense knowledge collectively. However, as aforementioned, two properties of commonsense knowledge determine that we cannot acquire all commonsense knowledge with such approaches. First, the scale of commonsense knowledge could be enormous and it is infeasible to perform crowd-sourcing for commonsense knowledge acquisition on such a huge scale. Second, commonsense knowledge is often a kind of preference rather than fixed fact, and thus it is not suitable to represent commonsense knowledge with fixed triplets (e.g., <“rock,” NotUsedFor, “eat”>) as used in Cyc and ConceptNet. Recently, pre-trained language representation models (e.g., BERT [8] and RoBERTa [9]) have been developed to acquire rich human knowledge implicitly and have demonstrated promising results on many downstream tasks. However, as shown in LAMA [10] and TransOMCS [11], even though these models are good at capturing factual knowledge about named entities, they still struggle at capturing commonsense knowledge, especially those complex commonsense knowledge between eventualities (a linguistic term covering activities, states, and events after [12], e.g., “I am hurt”). One possible explanation is that compared with tokens or named entities, the distribution of eventualities is generally much more sparse. More importantly, as discussed by [1], much commonsense knowledge, which is trivial for humans, is typically not discussed in our daily language at all. As a result, even though these deep pre-trained language representation models are good at acquiring knowledge from textual data, they could not effectively acquire or reason commonsense knowledge they rarely or never see in the form of word sequences.

To explore a scalable way of acquiring commonsense knowledge, in this paper, we propose an approach to constructing a large-scale weighted eventuality knowledge graph, ASER (Activities, States, Events, and their Relations), by extending the traditional definition of selectional preference to higher-order selectional preference over eventualities. The eventualities (i.e., nodes of ASER) are extracted using selected dependency patterns. The edges are based on discourse relations (e.g., *Result*) in discourse analysis. As shown in Figure 1, both nodes and edges are associated with frequency-based weights to reflect higher-order selectional preferences given a specific linguistic (either dependency or discourse) pattern. As discussed by [4, 13], such frequency distribution can serve as a good fit for

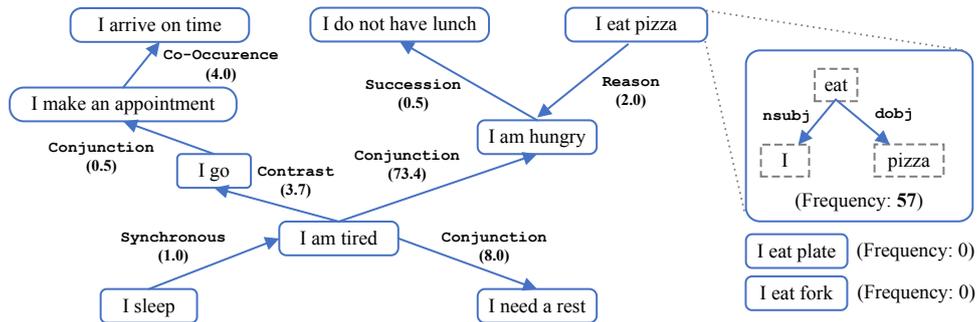


Figure 1: ASER Demonstration. Eventualities are connected with weighted directed edges. Each eventuality is a dependency graph.

humans’ selectional preference, which is indeed the commonsense knowledge. An example is shown in Figure 1. In ASER, “I eat plate” and “I eat fork” never appear in ASER while “I eat pizza” appears 57 times. We can infer that “plate” and “fork” are not subjects that can be eaten while “pizza” is. Similarly, the frequencies of edges can be used to reflect higher-order selectional preference between eventualities. For example, by observing that ⟨“Person be hungry”-Result-“Person eat”⟩ appear at least 12 times while ⟨“Person be hungry”-Reason-“Person eat”⟩ appears only once, we can know that “Person be hungry” is more likely to result in rather than be caused by “Person eat.” We argue that the higher-order selectional preference in ASER can be scalable and effective to represent previously defined commonsense knowledge types in ConceptNet [1] and potentially many other types of commonsense knowledge.

To build such a large-scale eventuality knowledge graph, we first leverage unsupervised algorithms and existing tools (e.g., dependency/discourse parsing) to extract eventualities and their relations from raw documents. For the eventuality extraction, considering that the English language’s syntax is relatively fixed and consistent across domains and topics, instead of defining complex triggers and role structures of events, we use syntactic patterns to extract all possible eventualities. We do not distinguish between semantic senses or categories of particular triggers or arguments in eventualities but treat all extracted words with their dependency relations as hyperedge in a graph to define an eventuality as a primitive semantic unit in our knowledge graph. For eventuality relation extraction, we adopt an end-to-end discourse parser [14] to determine the discourse relations between eventuality spans automatically and then create edges based on the predicted relation. Compared with previous commonsense knowledge acquisition methods, acquiring selectional preference knowledge with linguistic patterns and discourse relation prediction models is much cheaper and scalable. Thus, it can be used to extract large-scale selectional preference knowledge from the unlabeled corpus. After that, to overcome the challenge that a large portion of the commonsense knowledge is rarely expressed in textual corpus and motivated by the observation [15] that human beings often conceptualize the events to a more abstract level such that they can be applied to new events, we propose to leverage existing conceptualization techniques [16, 17] to automatically generalize the knowledge we observed and extracted to those unseen eventualities.

As a result, we create ASER, which contains 438,648,952 unique eventualities and 648,514,465 edges. Table 1 provides a size comparison between three variations of ASER² (i.e., core, full, and concept) and existing eventuality-related (or simply verb-centric) knowledge bases. Essentially, they are not large enough as modern knowledge graphs and inadequate for capturing the richness and complexity of eventualities and their relations. FrameNet [18] is considered the earliest knowledge base defining events and their relations. It provides annotations about relations among about 1,000 human-defined eventuality frames, which contain 27,691 eventualities. However, given the fine-grained definition of frames, the scale of the annotations is limited. ACE [19] (and its follow-up evaluation TAC-KBP [20]) reduces the number of event types and annotates more examples in each of the event types. PropBank [21] and NomBank [22] build frames over syntactic parse trees, and focus on annotating popular verbs and nouns. TimeBank focuses only on temporal relations between verbs [23]. While the aforementioned knowledge bases are annotated

²ASER (core) includes all extracted eventualities that appear more than once, ASER (full) includes all extracted eventualities, and ASER (concept) includes all conceptualized eventualities.

	# Eventuality	# Relation	# Relation Types
FrameNet	27,691	1,709	7
ACE	3,290	0	0
PropBank	112,917	0	0
NomBank	114,576	0	0
TimeBank	7,571	8,242	1
OMCS (Only include edges about eventualities)	74,989	116,097	4
Event2Mind	24,716	57,097	3
ProPora	2,406	16,269	1
ATOMIC	309,515	877,108	9
ATOMIC-2020	638,128	1,331,113	23
GLUECOSE	286,753	304,099	10
Knowlywood	964,758	2,644,415	4
ASER (core)	52,940,258	52,296,498	14
ASER (full)	438,648,952	648,514,465	14
ASER (concept)	15,640,017	224,213,142	14

Table 1: Size comparison of ASER and existing eventuality-related resources. # Eventuality, # Relation, and # Relation Types are the number of eventualities, relations between these eventualities, and relation types. For KGs containing knowledge about both entity and eventualities, we report the statistics about the eventualities subset. ASER (core) filters out eventualities that appear only once and thus has better accuracy while ASER (full) can cover more knowledge. ASER (concept) runs conceptualization to aggregate diverse relations from a much cleaner ASER, resulting in a much denser commonsense knowledge graph.

by domain experts, OMCS/ConceptNet³ [1], Event2Mind [24], ProPora [25], ATOMIC [26], ATOMIC-2020 [27], and GLUECOSE [28] leveraged crowdsourcing platforms or the general public to annotate commonsense knowledge about eventualities, in particular the relations among them. Furthermore, KnowlyWood [29] uses semantic parsing to extract activities (verb+object) from movie/TV scenes and novels to build four types of relations (parent, previous, next, similarity) between activities using inference rules. Compared with all these eventuality-related KGs, ASER is larger by one or more orders of magnitude in terms of the numbers of eventualities and relations it contains.

In summary, our contributions are as follows.

1. **Representation of commonsense knowledge with higher-order selectional preference:** We extend the original definition of selectional preference to higher-order selectional preference between eventualities, and show that we can cheaply acquire selectional preference knowledge from the unlabeled corpus and convert such knowledge into commonsense knowledge in the format of other commonsense knowledge bases such as ConceptNet and ATOMIC.
2. **Definition of ASER:** We define a brand new knowledge graph (KG) where the primitive units of semantics are eventualities. We organize our KG as a relational graph of hyperedges. Each eventuality instance is a hyperedge connecting several vertices, which are words. A relation between two eventualities in our KG represents one of the 14 relation types defined in PDTB [30] or a co-occurrence relation.
3. **Scalable Extraction of ASER:** We perform eventuality extraction over large-scale corpora. We design several high-quality patterns based on dependency parsing results to extract all eventualities that match these patterns and then apply a discourse parsing system to extract the eventuality relations. In the end, we leverage a conceptualization module to generalize the extracted knowledge to unseen eventualities.
4. **Inference over ASER:** We also provide several ways of commonsense inference over ASER. We show that both eventuality and relation retrieval over one-hop or multi-hop relations can be modeled as conditional probability inference problems.
5. **Evaluation and Applications of ASER:** We conduct an extensive evaluation to demonstrate the quality of extracted eventuality knowledge and the transferability from such linguistic-based knowledge to commonsense knowledge.

³Following the original definition, we only select the four relations (“HasPrerequisite,” “HasFirstSubevent,” “HasSubEvent,” and “HasLast-SubEvent”) that involve eventualities.

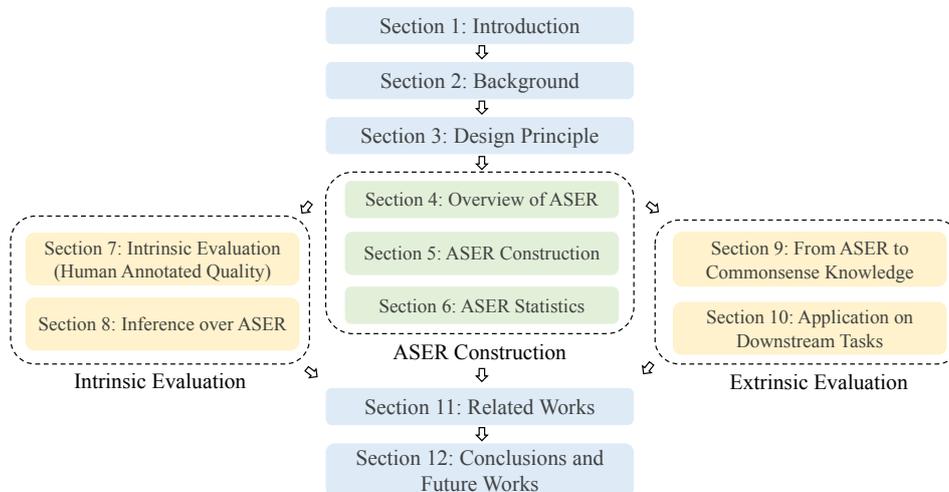


Figure 2: Paper Road Map.

The paper organization is presented in Figure 2. After the introduction section, we introduce background knowledge about previous works on large-scale knowledge bases construction in Section 2. And then, in Section 3, we introduce the design principles of ASER. Based on these principles, we present the construction details of ASER in Section 4, 5, and 6. Specifically, in Section 4, we show the overall framework and all used notations in ASER. After that, we discuss how to extract those eventualities, eventuality relations, and concept-level eventuality knowledge in Section 5. All statistics including the number of unique eventualities and edges are presented in Section 6. After constructing ASER, we conduct both intrinsic and extrinsic evaluations and analyses to analyze the quality of our knowledge base. In Section 7, we randomly sample eventualities and edges from ASER and invite crowdsourcing annotators from the Amazon Mechanical Turk to annotate the quality. To better understand the ASER knowledge, we conduct an in-depth inspection of potential inference on ASER knowledge in Section 8. To prove that ASER can indeed cover rich commonsense knowledge, we conduct experiments in Section 9 to demonstrate the transferability from SP knowledge to commonsense knowledge defined in other human-defined commonsense knowledge bases such as ConceptNet [1] and ATOMIC [26]. After that, in Section 10, we show that the knowledge in ASER could be helpful for downstream tasks such as commonsense reading comprehension and daily dialogue generation. In Section 11, we introduce the related works about commonsense knowledge acquisition, traditional syntactic-based information extraction, and conceptualization. In the end, we conclude this paper and introduce all the released resources with Section 12.

2. Background

In his conceptual semantics theory, Ray Jackendoff, a Rumelhart Prize⁴ winner, describes semantic meaning as “a finite set of mental primitives and a finite set of principles of mental combination [31].” The primitive units of semantic meanings include *Thing* (or *Object*), *Activity*,⁵ *State*, *Event*, *Place*, *Path*, *Property*, *Amount*, etc. Understanding the semantics related to the world requires the understanding of these units and their relations. Traditionally, linguists and domain experts built knowledge graphs (KGs)⁶ to formalize these units and enumerate categories (or senses) and relations of them. Typical KGs include WordNet [34] for words, FrameNet [18] for events, and Cyc [7]

⁴The David E. Rumelhart Prize is funded for contributions to the theoretical foundations of human cognition.

⁵In his original book, he called it *Action*. But given the other definitions and terminologies we adopted [32, 12], it means *Activity*.

⁶Traditionally, people used the term “knowledge base” to describe the database containing human knowledge. In 2012, Google released its knowledge graph where vertices and edges in a knowledge base are emphasized. We discuss in the context of the knowledge graph, as our knowledge is also constructed as a complex graph. For more information about terminologies, please refer to [33].

and ConceptNet [1] for commonsense knowledge. However, their small scales restricted their usage in real-world applications.

Nowadays, with the growth of web contents, computational power, and the availability of crowdsourcing platforms, many modern and large-scale KGs, such as Freebase [35], KnowItAll [36], TextRunner [37], YAGO [38, 39], BabelNet [40], DBpedia [41], NELL [42], Probase [43], and Google Knowledge Vault [44], have been built based on semi-automatic mechanisms. Most of these KGs are designed and constructed based on facts about *Things* or *Objects*, such as instances and their concepts, named entities and their categories, as well as their properties and relations. On top of them, a lot of semantic understanding problems such as question answering [45] can be supported by grounding natural language texts on knowledge graphs, e.g., asking a bot for the nearest restaurants for lunch. Nevertheless, these KGs may fall short in circumstances that require not only fact knowledge about *Things* or *Objects*, but also the commonsense knowledge about *Activities*, *States*, and *Events*. Consider the aforementioned utterance that a human would talk to the bot at 1 PM: “I am hungry,” which may also imply one’s need for a restaurant recommendation. This, however, will not be possible unless the bot can identify that the consequence of being hungry would be “having lunch” at noon.

In this paper, we propose to leverage higher-order selectional preference to discover and store commonsense knowledge about *Activities* (or process, e.g., “I sleep”), *States* (e.g., “I am hungry”), *Events* (e.g., “I make a call”), and their *Relations* (e.g., “I am hungry” may result in “I have lunch”), for which we call ASER. In fact, *Activities*, *States*, and *Events*, which are expressed by verb-related clauses, are all eventualities following the commonly adopted terminology and categorization proposed by Mourelatos [32] and Bach [12]. Previous literature on eventualities mostly focuses on extracting eventualities from text with pre-defined event schemas, which enumerates triggers with senses and arguments with roles, defined in FrameNet [18] or ACE [19]. However, as the pre-defined event ontology is often domain-specific and small (e.g., ACE contains 33 event types), the extracted events cannot cover all commonsense. Different from them, instead of using a small event ontology, we use patterns over the dependency graphs, which could contain multiple words and dependency edges, to extract eventualities. Any events that satisfy the pre-defined patterns will be extracted, and thus we achieve much broader coverage. Besides the eventuality extraction, extracting relations between eventualities is another vital research problem in the NLP community. For example, HieVe [46] focuses on extracting super-sub event relations, and TimeBank [23] focuses on the temporal relations. These works typically focus on identifying implicit relations, which is a very challenging task, and the state-of-the-art models can only achieve 59.5 F1 [47] and 75.5 F1 [48] on the HieVe and TimeBank datasets, respectively. As a result, current models are still not ready to be used to extract high-quality relations between events. At the same time, the current state-of-the-art model on implicit discourse relations can only achieve the accuracy of 57.26 [49] on CoNLL-2015 dataset [50]. As an alternative, we discard the implicit relations and only focus on explicit discourse relations between events. By doing so, we sacrifice the recall but make sure the high-quality of the collected knowledge. For example, the used discourse parser proposed by [14] can guarantee 90.14% F1 on explicit discourse relation classification. Simultaneously, we try to scan a huge corpus to guarantee the resulting knowledge graph’s overall coverage.

3. Design Principles

As aforementioned, ASER is a large-scale eventuality-based knowledge graph. Here by eventuality, we mean *Activities*, *States*, and *Events*, which are defined based on the commonly adopted terminology and categorization proposed by Mourelatos [32] and Bach [12]:

- **Activity:** An activity is also called a process [32, 12]. Both activity and event are occurrences (actions) described by active verbs. An example is “The coffee machine is brewing coffee.”
- **State:** A state is usually described by a stative verb and cannot be qualified as actions. A typical state expression is “The coffee machine is ready for brewing coffee.”
- **Event:** An event is defined as an occurrence that is inherently countable [32]. For example, “The coffee machine brews a cup of coffee once more” is an event because it admits a countable noun “a cup” and cardinal count adverbials “once,” while “The coffee machine brews coffee” is not an event with an imperfective aspect which is not countable.

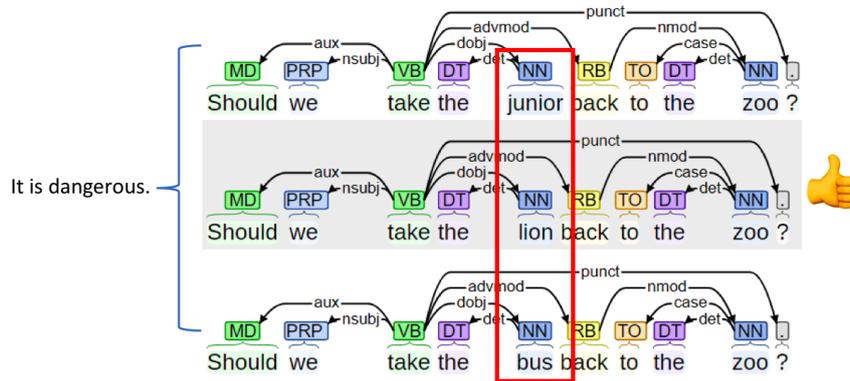


Figure 3: Principle Demonstration. When we fix the grammar, humans’ preferences over the linguistic description are reflecting the commonsense. The examples are based on ones used in [53].

Unlike the previous works [51, 52], we do not distinguish activities (or processes), states, and events. Instead, we use dependency patterns to represent all the eventualities that can be activities, states, and events and also discourse relations [30] such as *COMPARISON*, *Contrast* and *CONTINGENCY*, *Cause* as the relation types between eventualities based on the following two design principles.

3.1. The Lower Bound of a Semantic Theory

As discussed by the lower bound of a semantic theory [53], understanding human language requires both knowledge about the language (i.e., grammar) and knowledge about the world. As a result, if we fix the grammar structure of the linguistic descriptions, their difference will be mostly captured by the semantics. An example is shown in Figure 3. There are three sentences that share the same grammar structure but describe different events, which may have different reasons, effects, and sub-events. Given that the previous context is “It is dangerous,” humans normally will prefer the second sentence to appear in this context because a lion is a dangerous animal. And such preference can reflect the commonsense knowledge we are looking for. Motivated by this, we propose to use dependency patterns to categorize eventualities and discourse relations as the relations between eventualities. As a result, the frequency of eventualities and edges can be naturally used to represent humans’ preferences when the grammar structure is fixed.

Historically, such grammar-based semantics is called selectional preference [54], which is a relaxation of selectional restrictions [53]. Initially, the research on selectional preference focuses on the IsA hierarchy in WordNet [34] and verb-object dependency relations. Later on, the idea of selectional preference was extended to verb-subject dependency relations. Several first-order selectional preference examples are as follows.

- $SP(\text{Cat}, \text{IsA}, \text{Animal}) > SP(\text{Cat}, \text{IsA}, \text{Plant})$
- $SP(\text{Eat}, \text{dobj}, \text{Food}) > SP(\text{Eat}, \text{dobj}, \text{Rock})$
- $SP(\text{Sing}, \text{nsubj}, \text{Singer}) > SP(\text{Sing}, \text{nsubj}, \text{House})$

Recently, to represent more complex commonsense knowledge, the principle of selectional preference was extended to the second-order [13]. The motivation is that humans tend to have a strong preference over the property of certain verbs’ subjects and objects. For example, we can formalize the commonsense that the subject of eat is more likely to be “hungry” rather than “tasty” with the following second-order selectional preference:

- $SP(\text{Eat}, \text{Nsubj-}amod, \text{Hungry}) > SP(\text{Eat}, \text{Nsubj-}amod, \text{Tasty})$

More examples are shown in Table 2. Higher plausibility scores indicate that the annotators have a stronger preference for the combination. For the first-order selectional preference, people are most likely to select “meal” rather than “mail” as the object of “eat.” Similarly, we can see that “heavy” is a common property of the *object* of “lift.” As shown by the experiments in [13], such selectional preference knowledge is crucial for solving commonsense reasoning tasks such as Winograd Schema Challenge [6].

(a) dobj		(b) nsubj		(c) amod	
SP Pair	Plausibility	SP Pair	Plausibility	SP Pair	Plausibility
(eat, meal)	10.00	(singer, sing)	10.00	(fresh, air)	9.77
(close, door)	8.50	(law, permit)	7.78	(new, method)	8.89
(convince, people)	7.75	(women, pray)	5.83	(young, people)	6.82
(touch, food)	5.50	(realm, remain)	3.06	(medium, number)	4.09
(hate, investment)	4.00	(victim, contain)	2.22	(immediate, food)	2.50
(confront, impulse)	2.78	(bar, act)	1.39	(eager, price)	1.36
(eat, mail)	0.00	(textbook, eat)	0.00	(secret, wind)	0.75

(d) dobj_amod		(e) nsubj_amod	
SP Pair	Plausibility	SP Pair	Plausibility
(lift, heavy <i>object</i>)	9.17	(friendly <i>subject</i> , smile)	10.00
(design, new <i>object</i>)	8.00	(evil <i>subject</i> , attack)	9.00
(recall, previous <i>object</i>)	7.05	(recent <i>subject</i> , demonstrate)	6.00
(attack, small <i>object</i>)	5.23	(random <i>subject</i> , bear)	4.00
(drag, drunk <i>object</i>)	4.25	(happy <i>subject</i> , steal)	2.25
(inform, weird <i>object</i>)	3.64	(stable <i>subject</i> , understand)	1.75
(earn, rubber <i>object</i>)	0.63	(sunny <i>subject</i> , make)	0.56

Table 2: Examples of first-order and second-order selectional preference and their plausibility ratings provided by human annotators [13]. *Object* and *subject* are place holders to help understand the second-hop selectional preference relations.

In this work, we further extend the idea of selectional preference to discourse relations between eventualities, which is denoted as higher-order selectional preference over eventualities.

3.2. The Need of Aggregating “Partial Information” in Commonsense Reasoning

As discussed by [5], to effectively represent the selectional preference over linguistic relations and use that knowledge for language inference, we need to do aggregation over the “partial information,” which may “not be invariably true” but “tends to be of a very high degree of generality indeed” [1]. For example, Wilks used the following sentence as an example.

- The soldiers fired at the women, and we saw several of them fall.

We know that *them* should refer to *women* rather than *soldiers* because we have the partial information that “hurt things tending to fall down.” Formally, it can be translated into the following form:

- (hurt, X) $\xrightarrow{\text{ResultIn}}$ (X, fall).

There are many ways to find such representations of knowledge, e.g., first or even second-order logic. However, existing logic-based or semantic frame-based methods such as combinatory categorial grammar [55] or semantic role labeling [18, 56] require large amounts of annotation. Moreover, the semantic roles defined in labeled frames [18, 56] are too coarse-grained to support fine-grained conceptual reasoning. An efficient way of acquiring such partial information is to do the aggregation over collected selectional preference about the instance-level eventualities and their conceptualizations. We have shown that using such higher-order selectional preference, we can solve a subset of Winograd Schema Challenge (WSC) [6] with 70% accuracy [13]. For example, to solve the WSC example:

- The fish ate the worm. It was tasty.
- The fish ate the worm. It was hungry.

we can merge all subjects and object, and get the following frequency information:

- Frequency(‘X eats Y’, co-occur, ‘X is hungry’) = 18 and Frequency(‘X eats Y’, co-occur, ‘Y is hungry’) = 1;
- Frequency(‘X eats Y’, co-occur, ‘X is tasty’) = 0 and Frequency(‘X eats Y’, co-occur, ‘Y is tasty’) = 7.

These numbers reflect the aforementioned second-order selectional preferences based on which we can solve the questions. Although such aggregation has been shown to be useful for the Winograd Schema Challenge, the collected partial information can be too coarse. It only aggregates all information to be X or Y. However, in real-world applications, we also need to know the following question for fine-grained concepts other than humans

- Frequency(‘Company acquires Startups’, ResultIn, ‘Stock increases’)=?

Therefore, a principled way of performing the conceptualization of instances and partial concept information aggregation is needed. Thus, we propose to leverage another existing knowledge base, Probase [43], to perform conceptualization [16, 17] over entities that Probase can recognize. For example, after observing that both “having a cat” and “having a dog” can cause “being happy,” and with the help of Probase, we can conceptualize and aggregate “having a cat” and “having a dog” to be “having a pet,” we can then conclude that “having a pet” can cause “being happy.”

One thing worth mentioning is that the main methodology of ASER is that after the aggregation, the more heavily weighted (i.e., frequent) eventualities or edges make more sense than the less heavily weighted ones. As a result, when we conduct the conceptualization, we do not need to consider the context because other eventualities exist, and after the aggregation, the more heavily weighted ones will still make more sense. For example, given the eventuality “I eat apple,” we do not need to worry about which one of “fruit” and “company” we should conceptualize “apple” to because we will see many other eventualities related to fruit such as “I eat banana” and “I eat orange.” In the end, the overall weight of “I eat fruit” will still be much higher than “I eat company.”

4. Overview of ASER

ASER is a hybrid graph combining a hypergraph $\{\mathcal{V}, \mathcal{E}\}$ where each hyperedge is constructed over vertices, and a traditional graph $\{\mathcal{E}, \mathcal{R}\}$ where each edge is built among eventualities. For example, $E_h=(I, am, hungry)$ and $E_t=(I, eat, anything)$ are eventualities, where we omit the internal dependency structures for brevity. They have a relation $\langle E_h, Result, E_t \rangle$, where Result is the relation type. We devise the formal definition of ASER as below.

Definition 1. ASER KG is a hybrid graph \mathcal{H} of eventualities E 's. Each **eventuality** E is a hyperedge linking to a set of vertices v 's. Each vertex v is a **word** in the vocabulary. We define $v \in \mathcal{V}$ in the vertex set and $E \in \mathcal{E}$ in the hyperedge set. $\mathcal{E} \subseteq \mathcal{P}(\mathcal{V}) \setminus \{\emptyset\}$ is a subset of the power set of \mathcal{V} . We also define a **relation** $R_{i,j} \in \mathcal{R}$ between two eventualities E_i and E_j , where \mathcal{R} is the relation set. Each relation has a **type** $T \in \mathcal{T}$ where \mathcal{T} is the type set. Overall, we have ASER KG $\mathcal{H} = \{\mathcal{V}, \mathcal{E}, \mathcal{R}, \mathcal{T}\}$.

4.1. Eventuality

Unlike named entities or concepts, which are noun phrases, eventualities are usually expressed as verb phrases, which are more complicated in structure. Our definition of eventualities is built upon the following two assumptions: (1) Syntactic English patterns are relatively fixed and consistent; (2) The eventuality's semantic meaning is determined by the words it contains. To avoid the extracted eventualities being too sparse, we use words fitting specific patterns rather than a whole sentence to represent an eventuality. Also, to make sure the extracted eventualities have complete semantics, we retain all necessary words extracted by patterns rather than those simple verbs or verb-object pairs in sentences. The selected patterns are shown in Table 3. For example, for the eventuality (dog, bark), we have a relation nsubj between the two words to indicate that there is a subject-of-a-verb relation in between. We now formally define an eventuality as follows.

Definition 2. An eventuality E is a hyperedge linking multiple words $\{v_1, \dots, v_N\}$, where N is the number of words in eventuality E . Here, $v_1, \dots, v_N \in \mathcal{V}$ are all in the vocabulary. A pair of words in E (v_i, v_j) may follow a syntactic relation $e_{i,j}$. The weight of E , denoted as $w_E^{(e)}$, is defined by the frequencies of appearance in the whole corpora.

Pattern	Code	Example
n_1 -nsubj- v_1	s-v	“The dog barks”
n_1 -nsubj- v_1 -dobj- n_2	s-v-o	“I love you”
n_1 -nsubj- v_1 -xcomp- a	s-v-a	“He felt ill”
n_1 -nsubj- v_1 -xcomp- v_2	s-v-v	“I want to go”
n_1 -nsubj-(v_1 -iobj- n_2)-dobj- n_3	s-v-o-o	“You give me the book”
n_1 -nsubj- v_1 -xcomp- v_2 -dobj- n_2	s-v-v-o	“I want to eat the apple”
n_1 -nsubj-(v_1 -dobj- n_2)-xcomp- v_2 -dobj- n_3	s-v-o-v-o	“I ask you to help us”
n_1 -nsubj-(v_1 -dobj- n_2)-xcomp-(v_2 -iobj- n_3)-dobj- n_4	s-v-o-v-o-o	“president urges the congress to make her citizen”
n_1 -nsubj- a_1 -cop- be	s-be-a	“The dog is cute”
n_1 -nsubj- n_2 -cop- be	s-be-o	“He is a boy”
n_1 -nsubj- v_1 -xcomp- n_2 -cop- be	s-v-be-o	“I want to be a hero”
n_1 -nsubj- v_1 -xcomp- a_1 -cop- be	s-v-be-a	“I want to be slim”
n_1 -nsubj-(v_1 -iobj- n_2)-xcomp- n_3 -cop- be	s-v-o-be-o	“I want her to be hero”
n_1 -nsubj-(v_1 -iobj- n_2)-xcomp- a_1 -cop- be	s-v-o-be-a	“I want her to be happy”
$there$ -expl- be -nsubj- n_1	there-be-o	“There is an apple”
n_1 -nsubjpass- v_1	spass-v	“The bill is paid”
n_1 -nsubjpass- v_1 -dobj- n_2	spass-v-o	“He is served water”
n_1 -nsubjpass- v_1 -xcomp- v_2 -dobj- n_2	spass-v-v-o	“He is asked to help us”

Table 3: Selected eventuality patterns (“v” stands for normal verbs other than “be,” “be” stands for “be” verbs, “n” stands for nouns, “a” stands for adjectives, and “p” stands for prepositions.), Code (to save space, we create a unique code for each pattern and will use that in the rest of this paper), and the corresponding examples.

We use patterns from dependency parsing to extract eventualities E ’s from unstructured large-scale corpora. Here $e_{i,j}$ is one of the relations that dependency parsing may return. Although the recall is sacrificed in this way, our patterns are of high precision, and we use massive corpora to extract as many eventualities as possible. This strategy is also shared with many other modern KGs [36, 37, 42, 43].

4.2. Eventuality Relation

For relations among eventualities, as introduced in Section 1, we follow PDTB’s [30] definition of relations between sentences or clauses but simplify them to eventualities. Following the CoNLL 2015 discourse parsing shared task [50], we select 14 discourse relation types and an additional co-occurrence relation to build our knowledge graph.

Definition 3. A relation R between a pair of eventualities E_h and E_t has one of the following types $T \in \mathcal{T}$ and all types can be grouped into five categories: **Temporal** (including Precedence, Succession, and Synchronous), **Contingency** (including Reason, Result, and Condition), **Comparison** (including Contrast and Concession), **Expansion** (including Conjunction, Instantiation, Restatement, Alternative, ChosenAlternative, and Exception), and **Co-Occurrence**. The detailed definitions of these relation types are shown in Table 4. The weight of $R = \langle E_h, T, E_t \rangle$, which is denoted as $w_R^{(r)}$, is defined by the sum of weights of $\langle E_h, T, E_t \rangle$ that appear in the whole corpora.

4.3. ASER Conceptualization

As aforementioned, to overcome the challenge that trivial commonsense is often omitted in humans’ communication, we propose to leverage the conceptualization to generalize the knowledge about observed eventualities to unseen ones. For each eventuality $E \in \mathcal{E}$, whose weight is $w_E^{(e)}$, and we can conceptualize E to E' with confidence $w_{E,E'}^{(c)}$, we will get a new conceptualized eventuality E' with the weight $w_{E'}^{(e)} = w_E^{(e)} \cdot w_{E,E'}^{(c)}$. Similarly, assume that an edge $R \in \mathcal{R}$ is $\langle E_h, T, E_t \rangle$ and its weight is $w_R^{(r)}$, and E_h and E_t can be conceptualized to E'_h and E'_t with the confidence $w_{E_h,E'_h}^{(c)}$ and $w_{E_t,E'_t}^{(c)}$, respectively. We can then get a new conceptualized edge $\langle E'_h, T, E'_t \rangle$ with the weight $w_R^{(r)} \cdot w_{E_h,E'_h}^{(c)} \cdot w_{E_t,E'_t}^{(c)}$. Details about how to leverage an external hypernym knowledge base to get the conceptualized eventualities and determine the confidence scores are presented in Section 5.

Relation	Explanation
$\langle E_h, \text{Precedence}, E_t \rangle$	E_h happens before E_t .
$\langle E_h, \text{Succession}, E_t \rangle$	E_h happens after E_t .
$\langle E_h, \text{Synchronous}, E_t \rangle$	E_h happens at the same time as E_t .
$\langle E_h, \text{Reason}, E_t \rangle$	E_h happens because E_t happens.
$\langle E_h, \text{Result}, E_t \rangle$	If E_h happens, it will result in the happening of E_t .
$\langle E_h, \text{Condition}, E_t \rangle$	Only when E_t happens, E_h can happen.
$\langle E_h, \text{Contrast}, E_t \rangle$	E_h and E_t share a predicate or property and have significant difference on that property.
$\langle E_h, \text{Concession}, E_t \rangle$	E_h should result in the happening of E' , but E_t indicates the opposite of E' happens.
$\langle E_h, \text{Conjunction}, E_t \rangle$	E_h and E_t both happen.
$\langle E_h, \text{Instantiation}, E_t \rangle$	E_t is a more detailed description of E_h .
$\langle E_h, \text{Restatement}, E_t \rangle$	E_t restates the semantics meaning of E_h .
$\langle E_h, \text{Alternative}, E_t \rangle$	E_h and E_t are alternative situations of each other.
$\langle E_h, \text{ChosenAlternative}, E_t \rangle$	E_h and E_t are alternative situations of each other, but the subject prefers E_h .
$\langle E_h, \text{Exception}, E_t \rangle$	E_t is an exception of E_h .
$\langle E_h, \text{Co-Occurrence}, E_t \rangle$	E_h and E_t appear in the same sentence.

Table 4: Eventuality relation types between two eventualities E_h and E_t and explanations.

4.4. KG Storage

In total, we use the following three tables of the SQLite database to store ASER.

- **Eventuality:** As aforementioned, all eventualities in ASER are dependency graphs, where vertices are the words and edges are dependency relations. We generate unique “eids” for eventualities by hashing their words, pos-tags, and dependencies and store eventualities in an *Eventuality* table with SQLite database where “eids” is the key, and patterns, verb(s), skeleton words, words, pos-tags, dependencies, and frequencies are the other attribute columns.
- **Concept:** To effectively distinguish the eventualities before and after the conceptualization, we store eventualities created by the conceptualization step in another *Concept* table and denote the id as “cid.” As the dependency edges are inherited from the original eventualities, we only hash the conceptualized words to generate the “cids.” For each conceptualized eventuality, we store its “cid,” pattern, frequency, and “eids” of the original eventualities.
- **Relations:** We store the relations between eventualities in the *Relations* table. For each pair of eventualities (i.e., E_h and E_t), if there is at least an edge between them, we will create an instance and generate a “rid” for them by hashing the concatenation of their “eids.” For the storage efficiency and retrieval feasibility, we store all edges and the associated weights between E_h and E_t as well as the eventuality ids of E_h and E_t in that instance.

5. ASER Construction

In this section, we introduce the ASER construction details.

5.1. System Overview

The overall framework of our extraction system is shown in Figure 4. After collecting the raw corpora, we first preprocess the texts with the dependency parser. Then we perform eventuality extraction with pattern matching. We collect sentences and adjacent sentence pairs that contain more than two eventualities into an instance collection. After that, we extract discourse relations from these candidate instances with the help of an explicit discourse parser [14]. Considering that the discourse parser’s discourse argument span might not be identical to the extracted eventualities, we apply token-based Simpson’s similarity between the arguments spans and eventualities to determine whether the discourse arguments are enough to represent the meaning of the extracted eventualities. We only keep the extraction

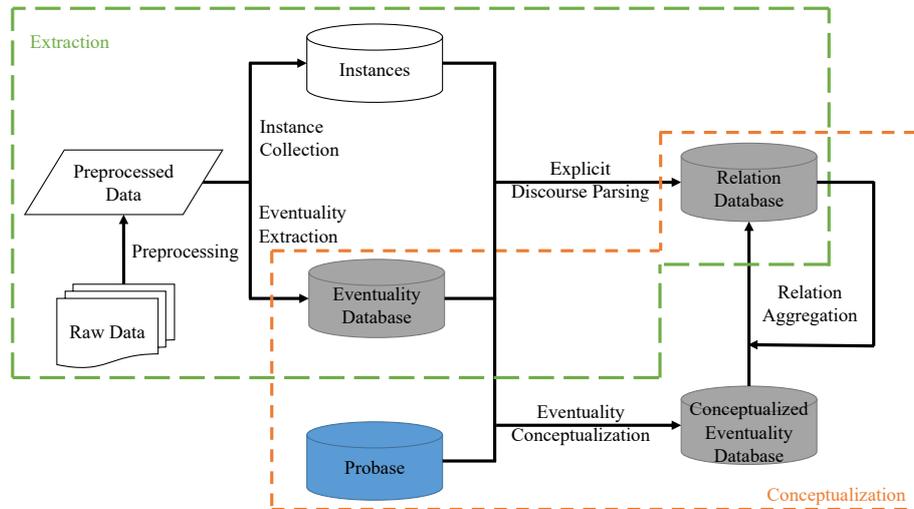


Figure 4: ASER construction framework. The extraction and the conceptualization process are shown in the orange dash-dotted and green dashed boxes, respectively. The blue database is Probase and three gray databases are the resulted ASER.

results with the Simpson’s similarity larger than 0.8. After the initial ASER construction, we leverage the *IsA* relations between nouns and named entities from Probase [43] to conduct the conceptualization. In the end, we aggregate relations between conceptualized eventualities by retrieving head and tail eventualities from the conceptualized eventuality database and the eventuality relation database. In the following sub-sections, we will introduce each part of the system separately.

5.2. Corpora

To ensure the broad coverage of ASER, we select corpora from different resources (reviews, news, forums, social media, movie subtitles, e-books) as the raw data. The details of these datasets are as follows.

- **Yelp:** Yelp is a social media platform where users can write reviews for businesses (e.g., restaurants). The latest release of the Yelp dataset⁷ contains over five million reviews.
- **New York Times (NYT):** The NYT [57] corpus contains over 1.8 million news articles from the NYT throughout 20 years (1987 - 2007).
- **Wiki:** Wikipedia is one of the largest free knowledge datasets. To build ASER, we select the English version of Wikipedia.⁸
- **Reddit:** Reddit is one of the largest online forums. In this work, we select the anonymized post records⁹ over one period month.
- **Movie Subtitles:** The movie subtitles corpus was collected by [58], and we select the English subset, which contains subtitles for more than 310K movies.
- **E-books:** The last resource we include is the free English electronic books from Project Gutenberg.¹⁰

We merge these resources as a whole to perform the knowledge extraction. The detailed statistics are presented in Table 5.

5.3. Preprocessing

For each document, we aim to extract eventualities, relations between eventualities, conceptualized eventualities, and relations between conceptualized eventualities. Based on the consideration of the text parsing complexity and

⁷<https://www.yelp.com/dataset/challenge>

⁸<https://dumps.wikimedia.org/enwiki/>

⁹<https://www.reddit.com/r/datasets/comments/3bxlg7>

¹⁰<https://www.gutenberg.org/>

Name	# Sentences	# Tokens	Corpus Size	Category
YELP	54.5 M	838.8 M	2.5G	Reviews
NYT	49.8 M	1,179.4 M	3G	News
Wiki	110.6 M	2,435.4 M	13G	Knowledge
Reddit	253.6 M	3,371.3 M	21G	Forum
Subtitles	444.6 M	3,229.4 M	13G	Movie Scripts
E-books	210.6 M	3,610.0 M	21G	Stories
Overall	1,123.7 M	14,664.2 M	73.5G	-

Table 5: Statistics of used corpora. (M means millions and G means Gigabytes.)

quality, we parse each paragraph¹¹ instead of a whole document with the CoreNLP tool¹² to acquire the lemmatized tokens, pos-tags, named entities, the dependency graph, and the constituency tree. Before parsing, we replace URLs with a special token $\langle URL \rangle$ and drop tables in Reddit data.

5.4. Eventuality Extraction

To ensure that all the extracted eventualities are semantically complete without being too complicated, we design 18 patterns to extract the eventualities via pattern matching. Each of the patterns contains three kinds of dependency edges: positive dependency edges, optional dependency edges, and negative dependency edges. All the positive edges are shown in Table 3. Six more dependency relations (*advmod*, *amod*, *nummod*, *aux*, *compound*, and *neg*) are optional dependency edges that can associate with any of the selected patterns. We omit all optional edges in the table because they are the same for all patterns. All other dependency edges are considered negative dependency edges, designed to ensure all the extracted eventualities are semantically complete and all the patterns are exclusive with each other. Take sentence “I have a book” as an example, we will only select $\langle \text{“I,” “have,” “book”} \rangle$ rather than $\langle \text{“I,” “have”} \rangle$ as the valid eventuality, because “have”-*dobj*-“book” is a negative dependency edge for pattern “s-v.”

To extract eventualities from sentence s , considering that s may contain multiple eventualities, we first split it into simple clauses based on the constituency tree. To do so, besides the commonly used *SBAR* node, we also follow previous discourse parsing systems [14] to use a connective classifier to detect possible separators. As a result, we split sentences based on both the subordinate conjunctions and connectives. After that, for each verb v in sentence s , we find the dependency graph \mathcal{D} of the simple clause that contains v . We then try to match \mathcal{D} with all patterns one by one. For each pattern, we put the verb v as the starting point (i.e., v_1 in the pattern) and then try to find all the positive dependency edges. If we can find all the positive dependency edges around the center verb, these matched edges and words linked by these edges are considered as potential edges and words of a valid eventuality. Next, other edges and words are added via optional dependency edges. In the end, we will check if any negative dependency edge can be found in the dependency graph. If not, we will keep current edges and words as a valid eventuality. Otherwise, we will disqualify it. The pseudo-code of the eventuality extraction algorithm is in Algorithm 1. The time complexity of eventuality extraction is $O(|S| \cdot |\overline{\mathcal{D}}| \cdot |\overline{\mathcal{V}^{(v)}}|)$ where $|S|$ is the number of sentences, $|\overline{\mathcal{D}}|$ is the average number of dependency edges in a dependency parse tree, and $|\overline{\mathcal{V}^{(v)}}|$ is the average number of verbs in a sentence.

5.5. Eventuality Relation Extraction

We then introduce how to extract the relations between eventualities. Specifically, we employ an end-to-end discourse parser to extract the discourse relations. The discourse parser’s job is to parse a piece of text into a set of discourse relations between two adjacent or non-adjacent discourse units. Take the sentence “I have a story book, but it is not interesting.” as an example. Ideally, a good discourse parser extracts “I have a story book” as *arg1*, “it is not interesting” as *arg2*, “but” as the connective, and annotate the relation as “Contrast.” In our current pipeline, we use the state-of-the-art discourse parser [14], which is pre-trained on the CoNLL 2015 Shared Task data (PDTB) [50]. From CoNLL 2015 results,¹³ we can find out that this discourse parser can achieve 90.00% and 90.79% F1 scores on

¹¹As the discourse parser extracts discourse relations by the constituency tree of a sentence or trees of adjacent sentences, parsing sentences one by one would miss or misclassify some discourse relations.

¹²<https://stanfordnlp.github.io/CoreNLP>

¹³<https://www.cs.brandeis.edu/~clp/conll15st/results.html>

Algorithm 1 Eventuality Extraction with One Pattern p

INPUT: Parsed dependency graph \mathcal{D} , center verb v , positive dependency edges $\mathcal{P}_p^{(pos)}$, optional edges $\mathcal{P}_p^{(opt)}$, and negative edges $\mathcal{P}_p^{(neg)}$.

OUTPUT: extracted eventuality E .

```

1: Initialize eventuality edge list  $\mathcal{D}'$ .
2: Set the center verb  $v$  as the  $v_1$  in the pattern  $p$ 
3: for each connection  $d$  (a relation and the associated word) in positive dependency edges  $\mathcal{P}_p^{(pos)}$  do
4:   if find  $d$  in  $\mathcal{D}$  then
5:     Append  $d$  in  $\mathcal{D}'$ .
6:   else
7:     Return null.
8:   end if
9: end for
10: for each connection  $d$  in optional dependency edges  $\mathcal{P}_p^{(opt)}$  do
11:   if find  $d$  in  $\mathcal{D}$  then
12:     Append  $d$  in  $\mathcal{D}'$ .
13:   end if
14: end for
15: for each connection  $d$  in negative dependency edges  $\mathcal{P}_p^{(neg)}$  do
16:   if find  $d$  in  $\mathcal{D}$  then
17:     Return null.
18:   end if
19: end for
20: Build eventuality instance  $E$  from  $\mathcal{D}'$ .
21: Return  $E$ .

```

the test data from PDTB and the blind test data from Wikinews respectively on the explicit relation classification, but performance drops to 42.72% and 34.45% on the implicit relation classification. Hence, to guarantee the extraction quality, we only consider the explicit discourse relations. In explicit discourse parsing, there are two situations: both arguments are in the same sentence or not. Statistics show that less than 0.1% arguments are located in non-adjacent sentences in the explicit scenario, so we simply assume that the first argument is located in the same sentence (SS) or the previous sentence (PS). Specifically, the explicit discourse parser is consist of five components: (1) connective extractor to identify whether a word is a possible connective, (2) arg1 position classifier to decide whether the arg1 is located in the same sentence as the connective c or the previous sentence of c ; (3) SS argument extractor to extract the spans of two arguments in the same sentence; (4) PS argument extractor to extract the spans of two arguments in adjacent sentences; (5) explicit relation classifier to classify the relation type of c . Extractors in this system are essentially binary classifiers to identify whether a word is a connective or a part of any argument. The pseudo-code of eventuality relation extraction algorithm is shown in Algorithm 2.

As the extracted arguments might not be identical as the extracted eventualities, we use the Simpson's similarity to determine whether the discourse relations between arguments can be assigned to the extracted eventualities:

$$w_{A,E}^{(sim)} = \text{Simpson}(A, E) = \frac{|\mathcal{W}_A \cup \mathcal{W}_E|}{\min\{|\mathcal{W}_A|, |\mathcal{W}_E|\}}, \quad (1)$$

where A is an argument, E is an eventuality, \mathcal{W}_A and \mathcal{W}_E are token sets of A and E , $|\cdot|$ is the size of a token set. If the similarity $\text{Simpson}(A, E) \geq 0.8$, we consider the argument-level relations relevant to A can be assigned to the eventuality E with a weight $w_{A,E}^{(sim)}$, which is inversely proportional to the size of all matched eventualities $|\mathcal{E}|$. It is worth noting that Eq. (1) allows one argument A , which could include multiple eventualities as long as all tokens in eventualities can be covered by A . In this situation, the weight of the relation between the eventuality $E_1 \in \mathcal{E}_h$ from arg1 and the eventuality $E_2 \in \mathcal{E}_t$ from arg2 is inversely proportional to the product of extracted eventuality sizes from two arguments $|\mathcal{E}_h| \cdot |\mathcal{E}_t|$. Section 5.7 provides detailed descriptions.

Algorithm 2 Eventuality Relation Extraction**INPUT:** Parsed constituency trees \mathcal{K}_1 and \mathcal{K}_2 from adjacent sentences.**OUTPUT:** Extracted relations \mathcal{R} .

```

1: Initialize relation list  $\mathcal{R}$  as empty.
2: Extract possible connectives  $\mathcal{C}$  by a connective extractor given  $\mathcal{K}_1$  and  $\mathcal{K}_2$ .
3: for each possible connective  $c \in \mathcal{C}$  do
4:   if two arguments of  $c$  in the same sentence then
5:     Extract  $A_1$  and  $A_2$  by a SS arguments extractor given  $c$  and the sentence.
6:   else
7:     Extract  $A_1$  by a PS argument1 extractor given  $c$  and  $\mathcal{K}_1$ .
8:     Extract  $A_2$  by a PS argument2 extractor given  $c$  and  $\mathcal{K}_2$ .
9:   end if
10:  if  $A_1$  is not null and  $A_2$  is not null then
11:    Classify the relation  $y$  by a explicit relation classifier given  $c$ ,  $\mathcal{K}_1$ , and  $\mathcal{K}_2$ 
12:    Find eventualities  $\mathcal{E}_h$  that are extracted from  $A_1$ .
13:    Find eventualities  $\mathcal{E}_t$  that are extracted from  $A_2$ .
14:    Set weight  $w$  as  $1/(|\mathcal{E}_h| \cdot |\mathcal{E}_t|)$ .
15:    for each eventuality  $E_h$  in  $\mathcal{E}_h$  do
16:      for each eventuality  $E_t$  in  $\mathcal{E}_t$  do
17:        Build relation instance  $R = \langle E_h, y, E_t \rangle$  with a weight  $w$ .
18:        Append  $R$  in  $\mathcal{R}$ .
19:      end for
20:    end for
21:  end if
22: end for
23: Return  $\mathcal{R}$ .

```

5.6. Enriching ASER with Conceptualization

We then introduce the conceptualization details. For each noun or pronoun in the extracted eventualities, we will try to conceptualized it to a higher level with the following steps. If it is a named entity, we will conceptualized it to the corresponding NER tags. Specifically, we include the 13 NER types: “Time,” “Date,” “Duration,” “Money,” “Percent,” “Number,” “Country,” “State or Province,” “City,” “Nationality,” “Person,” “Religion,” “URL.” If it is a personal pronoun (e.g., “I,” “you,” or “they”), we will conceptualize it to “PersonX.”¹⁴ As all aforementioned conceptualization is designed by experts, we set the conceptualization probability to be 1. If it is a regular noun, we will try to conceptualize it with Probase [43]. Specifically, for each noun, we will retrieve its top-five hypernyms (i.e., concepts) and the associated probability from Probase.

Given an eventuality E with m tokens t_1, t_2, \dots, t_m to be mapped into concept tokens, we conceptualize it to a conceptualized eventuality C with the probability:

$$\Pr(C|E) = \prod_{i=1}^m \Pr(t_i^{(c)}|t_i^{(e)}). \quad (2)$$

Here $t_i^{(c)}$ is the corresponding token-level concept for token $t_i^{(e)}$. And $\Pr(t_i^{(c)}|t_i^{(e)})$ is the likelihood for $\langle t_i^{(e)}, \text{IsA}, t_i^{(c)} \rangle$ provided by Probase or 1.0 if $t_i^{(e)}$ can be conceptualized with rules. For each conceptualized eventuality C , we would have a list of eventualities \mathcal{E}_C that can be conceptualized to it. We can then compute the overall weight of C with Eq. (3), where $w_E^{(e)}$ is the weight of E :

$$w_C^{(c)} = \sum_{E \in \mathcal{E}_C} \Pr(C|E) \cdot w_E^{(e)}. \quad (3)$$

¹⁴If there are multiple people in the same edge, we will distinguish them with “PersonX” and “PersonY” etc.

Conceptualized ASER

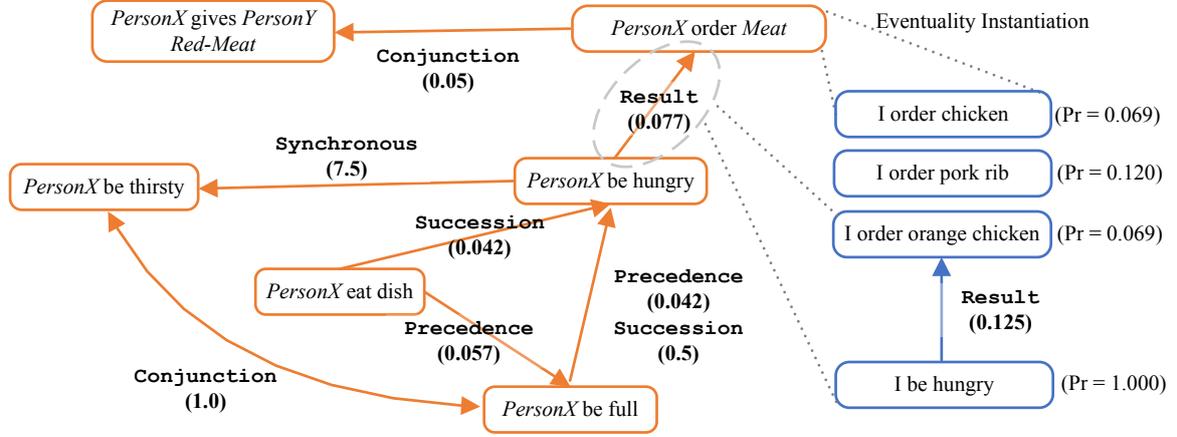


Figure 5: Demonstration of the conceptualized ASER. The eventualities are conceptualized and connected with weighted edges. Each concept contains its projections to specific eventualities.

We then introduce how to construct the edges between a conceptualized eventuality C and an original eventuality E . For any $E' \in \mathcal{E}_C$, if there is an edge $\langle E', T, E \rangle$ or $\langle E, T, E' \rangle$, we can then construct a new edge $\langle C, T, E \rangle$ or $\langle E, T, C \rangle$ with the weight based on Eq. (4) or Eq. (5), respectively, where $w_R^{(r)}$ means of weight of the relation R .

$$w_{\langle C, T, E \rangle}^{(r)} = \sum_{E' \in \mathcal{E}_C} \Pr(C|E') \cdot w_{\langle E', T, E \rangle}^{(r)}, \quad (4)$$

$$w_{\langle E, T, C \rangle}^{(r)} = \sum_{E' \in \mathcal{E}_C} w_{\langle E, T, E' \rangle}^{(r)} \cdot \Pr(C|E'). \quad (5)$$

As each conceptualized eventuality is correlated with a set of original eventualities, we need to aggregate the edges between the original eventualities to build the connections between the conceptualized ones. Formally, given two conceptualized eventualities C_h and C_t , we first retrieve all related original edges $\{\langle E_h, T, E_t \rangle | E_h \in \mathcal{E}_{C_h}, E_t \in \mathcal{E}_{C_t}\}$. Then we calculate the weight $\Pr(C_h|E_h) \cdot w_{\langle E_h, T, E_t \rangle}^{(r)} \cdot \Pr(C_t|E_t)$ for each related edge. Finally, we aggregate all weights to construct the weight as Eq. (6) for the edge between C_h and C_t associated with the relation type T .

$$w_{\langle C_h, T, C_t \rangle}^{(r)} = \sum_{E_h \in \mathcal{E}_{C_h}} \sum_{E_t \in \mathcal{E}_{C_t}} \Pr(C_h|E_h) \cdot w_{\langle E_h, T, E_t \rangle}^{(r)} \cdot \Pr(C_t|E_t). \quad (6)$$

An illustration of the conceptualized ASER is shown in Figure 5. We can get the conceptualized eventuality “*PersonX* be hungry” from “I am hungry,” “they are hungry,” and other extracted eventualities with $\Pr(C|E) = 1.000$ because their subjects (pronouns or names) are mapped to the token-level concept “*PersonX*.” As a comparison, “*PersonX* order Meat” is not a deterministic eventuality: it can be conceptualized from “I order chicken” with $\Pr(\text{PersonX order Meat} | \text{I order chicken}) = \Pr(\text{PersonX} | \text{I}) \cdot \Pr(\text{Meat} | \text{chicken}) = 1.000 \cdot 0.069 = 0.069$, “I order pork rib” with $\Pr(\text{PersonX order Meat} | \text{I order pork rib}) = \Pr(\text{PersonX} | \text{I}) \cdot \Pr(\text{Meat} | \text{pork rib}) = 1.000 \cdot 0.120 = 0.120$, or other extracted ones. Based on Eq. (3), after aggregating all weights together, we can get the concept weights for “*PersonX* be hungry” and “*PersonX* order Meat” are 1389.000 and 27.705, respectively. As for the relations between the two conceptualized eventualities, we find $w_{\langle \text{I am hungry}, \text{Result}, \text{I order orange chicken} \rangle}^{(r)} = 0.125$ and $w_{\langle \text{I am too hungry}, \text{Result}, \text{I order the fried chicken} \rangle}^{(r)} =$

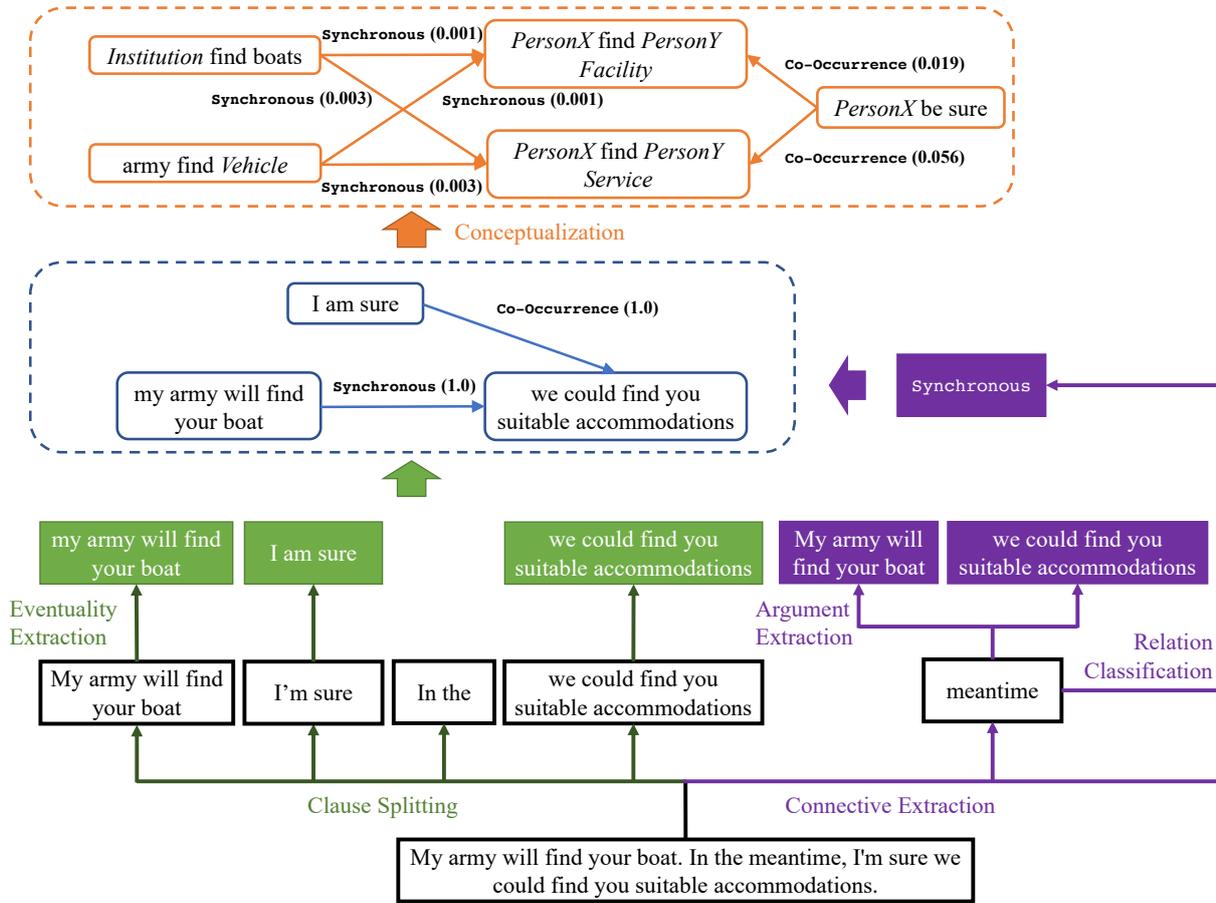


Figure 6: ASER building example. The eventuality extraction, the relation extraction, and the conceptualization process are shown in green, violet, and orange colors, respectively. For the clear representation, for each conceptualized eventuality, we only show the skeleton words and hide the optional ones. We only show the two conceptualized eventualities with the highest probabilities for each extracted eventuality. As a demonstration, the weights of eventualities and relations are calculated from this example instead of the whole KG.

1.000, so the relation weight in the concept-level is calculated as follows:

$$\begin{aligned}
 &W_{(PersonX \text{ be hungry}, Result, PersonX \text{ order Meat})}^{(r)} \\
 &= \Pr(PersonX \text{ be hungry} | I \text{ am hungry}) \cdot 0.125 \cdot \Pr(PersonX \text{ order Meat} | I \text{ order orange chicken}) \\
 &\quad + \Pr(PersonX \text{ be hungry} | I \text{ am too hungry}) \cdot 1.000 \cdot \Pr(PersonX \text{ order Meat} | I \text{ order the fried chicken}) \\
 &= 1.000 \cdot 0.125 \cdot 0.069 + 1.000 \cdot 1.000 \cdot 0.069 \\
 &= 0.077.
 \end{aligned}$$

Similarly, we can calculate all weights among conceptualized eventualities associated with different relation types. One thing worth mentioning is that the relation weights depend not only on the relation weights in the extracted knowledge bases but also on the conceptualization probabilities.

5.7. ASER Building Example

At the end of this section, we use an example to demonstrate the whole extraction pipeline. As shown in Figure 6, given a text “My army will find your boat. In the meantime, I’m sure we could find you suitable accommodations.”¹⁵

¹⁵This case comes from Movie Subtitles.

our system will first detect the possible connective “meantime” and split this text into four simple clauses: “My army will find your boat,” “In the,” “I’m sure,” and “we could find you suitable accommodations” with the constituency parsing. After that, our system will leverage the patterns designed in Table 3 to extract eventualities from the raw text by Algorithm 1. Simultaneously, two arguments “My army will find your boat” and “we could find you suitable accommodations” are extracted by argument extractors. The discourse parsing system predicts the corresponding discourse relation as Synchronous. As the first and last extracted eventualities can perfectly match the extracted arguments, we then create an edge \langle “my army will find your boat,” Synchronous, “we could find you suitable accommodations” \rangle . We also create an edge \langle “I am sure,” Co-Occurrence, “we could find you suitable accommodations” \rangle because the two eventualities appear in the same sentence. After extracting the original eventualities and edges, we then try to expand it with the conceptualization.¹⁶ For example, “I am sure” can be directly conceptualized as “*PersonX* be sure” directly because “I” is a personal pronoun. As both of the other two eventualities contain regular nouns (i.e., “army”), these eventualities can be conceptualized to multiple eventualities. After checking Probase, we find out that “army” can be conceptualized to “*Institution*” and “*Organization*” with the weights 0.058 and 0.038, “boat” can be conceptualized to “*Vehicle*” and “*Item*” with the weights 0.059 and 0.049, “accommodation” can be conceptualized to “*Service*” and “*Facility*” with the weights 0.056 and 0.019, respectively. We show the two most likely results for each original eventuality (if it has multiple possible conceptualization results) in Figure 6. In the end, we can construct edges between conceptualized eventualities, where the weights are the product of conceptualization probabilities, e.g., \langle “*Institution* find boats,” Synchronous, “*PersonX* find *PersonY* *Service*” \rangle with the weight $0.058 \times 0.056 = 0.003$, \langle “*PersonX* be sure,” Co-Occurrence, “*PersonX* find *PersonY* *Facility*” \rangle with the weight $1.000 \times 0.019 = 0.019$.

Pattern	ASER (full)		ASER (core)		ASER (concept)
	# Eventuality	# Unique	# Eventuality	# Unique	# Unique
s-v	351,082,855	100,645,728	260,663,083	14,337,769	1,022,415
s-v-o	284,103,317	159,948,356	139,031,585	18,100,360	8,252,653
s-v-a	11,546,768	6,149,584	5,951,980	752,468	139,087
s-v-v	24,549,946	11,129,566	14,624,526	1,591,424	216,413
s-v-o-o	6,154,685	3,789,253	2,765,728	460,526	514,084
s-v-v-o	29,445,708	18,659,717	12,720,497	2,187,577	1,482,783
s-v-o-v-o	3,863,478	2,674,229	1,462,883	288,326	522,613
s-v-o-v-o-o	91,532	59,290	40,428	8,499	18,461
s-be-a	79,235,136	29,845,112	52,068,570	3,733,978	465,747
s-be-o	98,411,474	53,503,410	49,979,659	6,337,042	2,312,209
s-v-be-a	1,927,990	982,438	1,035,864	123,263	29,738
s-v-be-o	2,322,890	1,574,896	909,250	184,298	139,239
s-v-o-be-a	277,087	191,973	100,917	18,793	6,151
s-v-o-be-o	307,031	231,289	95,815	22,411	32,796
there-be-o	16,021,849	6,642,438	10,013,628	953,041	39,500
spass-v	61,524,872	38,270,144	25,935,769	3,498,516	276,817
spass-v-o	5,519,982	4,129,709	1,677,244	330,229	154,410
spass-v-v-o	257,004	221,820	46,475	11,738	14,901
Overall	976,643,604	438,648,952	579,123,901	52,940,258	15,640,017

Table 6: Statistics of the eventuality extraction. # Eventuality and # Unique mean the total number and the unique number of extracted eventualities using corresponding patterns or conceptualized eventualities from them.

6. ASER Statistics

In total, we collect 976,643,604 eventualities from the raw documents. We filter those low-frequency eventualities that only appear once and retain 52,940,258 unique eventualities in ASER (core). From Table 6, we can find the “s-v”

¹⁶In the real system, we first extract the original ASER, and then apply the conceptualization step over the whole KG. The presented single sentence example is just for the demonstration.

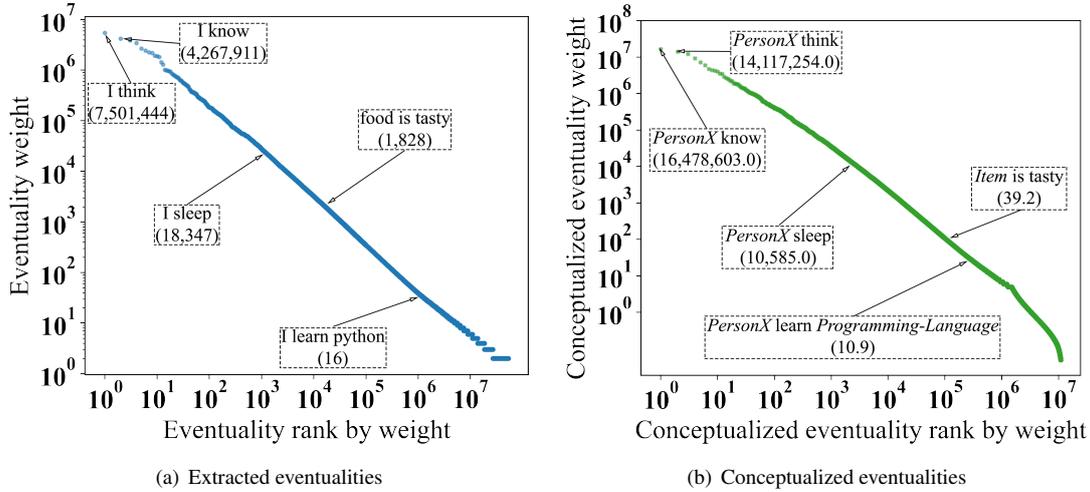


Figure 7: Eventuality distributions.

Relation	ASER (full)	ASER (core)	ASER (concept)
Precedence	14,058,213	1,790,016	4,798,015
Succession	4,939,291	663,183	1,963,820
Synchronous	19,464,898	3,123,042	8,013,943
Reason	9,775,829	2,205,076	6,439,128
Result	16,153,925	2,012,311	6,718,666
Condition	18,052,484	3,160,271	8,063,967
Contrast	59,333,901	8,655,661	24,978,311
Concession	5,684,395	477,155	1,499,276
Conjunction	82,121,343	13,978,907	45,597,200
Instantiation	1,278,381	18,496	93,266
Restatement	1,304,095	65,753	242,301
Alternative	3,539,892	583,174	123,883
ChosenAlternative	647,228	35,406	1,843,140
Exception	106,000	20,155	93,412
Co-Occurrence	412,054,590	49,232,161	113,744,814
Overall	648,514,465	86,020,767	224,213,142

Table 7: Statistics of the eventuality relation extraction.

and “s-v-o” are the most frequent patterns. On the other hand, even though those complex patterns appear relatively less frequently, thanks to the large scale of ASER, they still appear thousands to millions of times.

The original eventuality distribution is presented in Figure 7(a). In general, the distribution follows Zipf’s law, where only a small number of eventualities appear many times while the majority of eventualities appear only a few times. To better illustrate the distribution of eventualities, we also show several representative eventualities along with their weights, and we have two observations. First, eventualities which can be used in general cases, like “I think (7,501,444)” and “I know (4,267,911)” appear much more times than other eventualities. Second, eventualities in ASER are more closely related to our daily life like “I sleep (18,347)” or “food is tasty (1,828)” rather than domain-specific ones such as “I learn python (16).”

To achieve the balance between the quality and quantity of conceptualization results, we apply the conceptualization over eventualities whose frequencies are no less than five. As shown in Table 6, after the conceptualization, we get 15,640,017 more unique eventualities. It is obvious that patterns with more nouns (e.g., “s-v-o,” “s-v-v-o,” “s-be-

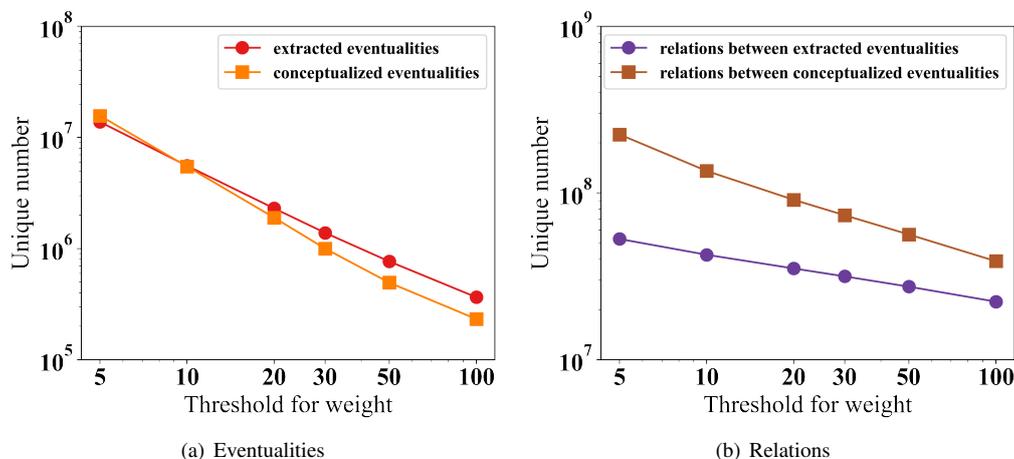


Figure 8: Distributions of eventualities and relations.

o”) dominate the conceptualized eventualities. The reason is that the conceptualization is only designed for nouns, and each noun phrase would be replaced with a general noun phrase if such hypernym relation appears in Probase. For conceptualized eventualities, we can observe a similar distribution in Figure 7(b). The top three conceptualized eventualities are “*PersonX* know” (16,478,603.0), “*PersonX* think” (14,117,254.0), and “*PersonX* say” (12,113,913.0). Although “I think” (7,501,444) appears the most in the raw data (“I know” appears 3,447,429 times in the raw data), but “you think” (1,444,333), “he thinks” (314,806), “they thinks” (205,432), “we think” (196,633), “it thinks” (174,729), “she thinks” (142,628), etc. appear much less than “you know” (4,726,264), “he knows” (396,013), “they know” (260,656), “we know” (457,115), “it knows” (247,409), “she knows” (190,803), etc., respectively. Finally, the weight of “*PersonX* know” exceeds that of “*PersonX* think.”

For relations, as shown in Table 7, we collect 648,514,465 unique relations from six data resources across different categories. To reduce noises in parsing and extraction, we also filter out relations that $\sum_{T' \in \mathcal{T}} w_{(E_h, T', E_t)}^{(r)} \leq 1$ where E_h and E_t are the head eventuality and the tail eventuality. Furthermore, if the head or the tail is filtered out by eventuality filtering, the relation is also dropped. Finally, we keep 86,020,767 unique relations in ASER (core), among which there are 36,788,606 relations belonging to 14 discourse relation types depending on the connectives and arguments, like *Conjunction* (e.g., “and”), *Contrast* (e.g., “but”), *Condition* (e.g., “if”), *Synchronous* (e.g., “meanwhile”), *Reason* (e.g., “because”), *Result* (e.g., “so”). When we filter out more low-frequency eventualities, the number of relations decreases slightly. For example, when we keep high-frequency eventualities whose frequencies are no less than five, 26.0% of eventualities (13,766,746) and 61.5% of relations (88,629,385) are preserved. We apply the conceptualization over these preserved eventualities and relations based on quantity and quality considerations. Finally, we obtain 15,640,017 unique conceptualized eventualities and 224,213,142 relations between these conceptualized eventualities. In total, we have about 26 times more relations between conceptualized eventualities than original eventualities.

To better understand the distributions of extracted and conceptualized knowledge, we show the number of eventualities and edges over different filtering thresholds in Figure 8(a) and Figure 8(b), respectively. For the extracted knowledge, the number of eventualities and relations decreases exponentially when the threshold ranges from 5 to 100. For the conceptualized knowledge, the rate of diminishing is even larger. When the threshold is less than 10, the conceptualized eventuality size is greater than the original size. But it is significantly less than the size of extracted eventualities as the threshold is larger than 10. On the other hand, the number of conceptualized eventuality relations consistently exceeds the original relation size, which results in a denser conceptualized knowledge graph.

7. Intrinsic Evaluation

In this section, we leverage human annotation to evaluate the quality of ASER from the following perspectives:

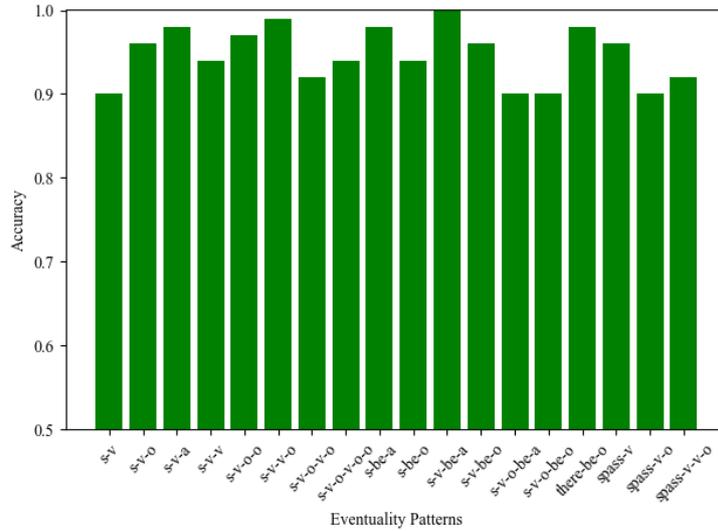


Figure 9: Human annotation of eventuality extraction quality.

- Eventuality Extraction:** We first evaluate how well the extracted eventualities can represent the original sentence’s semantics. For example, if the original sentence is “The kid goes to study,” eventuality “kid-go-to-study” with pattern “s-v-v” can fully represent the semantics, but eventuality “kid-go” with the pattern “s-v” cannot. We show the percentage of all extracted eventualities that can fully represent the original sentences’ core semantic based on different eventuality patterns.
- Lower-order Selectional Preference:** Besides the extraction quality, we also care about how well the eventuality statistics in ASER can reflect human’s selectional preference. For example, the frequency of “I eat food” should be higher than “I eat house.” As such preference appears inside eventualities, we denote them as the lower-order selectional preference.
- Discourse Extraction:** After evaluating the eventualities, we assess how well the extracted edges can correctly represent the discourse relations in the original sentence. For example, assume that the original sentence is “he went to school while I was still preparing the breakfast.” and we have successfully extracted two eventualities “he went to school” and “I was preparing breakfast,” the correct discourse relation between them should be “Synchronous” rather than “Contrast.” In this evaluation, we report the accuracy based on different discourse relations.
- Higher-order Selectional Preference:** Last but not least, we annotate whether edge frequencies in ASER can reflect the higher-order selectional preference among eventualities or not. For example, the frequency of “I am hungry”-Result-“I eat food” should be larger than “I am hungry”-Reason-“I eat food.”

Evaluation details and result analysis are as follows.

7.1. Eventuality Extraction

To evaluate the correctness of the selected eventuality patterns and the effectiveness of the extraction algorithm, we first employ the Amazon Mechanical Turk platform (MTurk)¹⁷ to evaluate the quality of eventuality extraction. We randomly select 50 extracted eventualities for each eventuality pattern and then provide these extracted eventualities along with their original sentences to the annotators. For each pair of eventuality and sentences, the annotators are asked to label whether the extracted eventuality phrase can fully and precisely represent the original sentences’ semantic meaning. If so, they should label them with “Valid.” Otherwise, they should label it with “Not Valid.” For

¹⁷<https://www.mturk.com/>

7.2. Low-order Selectional Preference

To evaluate whether the eventuality frequencies in ASER can reflect human’s low-order selectional preference, we first compare the plausibility of more frequent eventualities versus less frequent ones. For each eventuality pattern, we randomly select 50 eventuality pairs such that they only have a one-word difference but with a significant frequency difference. Specifically, we require the frequency of the high-frequency one to be larger than five, which is the medium frequency of all eventualities, and the frequency of the high-frequency one must be at least five times larger than the frequency of the low-frequency one. For each eventuality, we invite six annotators from MTruk to ask them which one of the eventualities seems more plausible to them. If more annotators agree that the more frequent one makes more sense, we will label that pair as a positive correlation. On the other hand, if more annotators agree that the less frequent one makes more sense, we will label that pair as a negative correlation. If the voting draws, we will label it as similar. As this evaluation fails to consider the eventualities with the frequency zero (i.e., they do not exist in ASER) and whether an eventuality exists or not is also a good preference indicator, we add another evaluation to prove that. For each eventuality pattern, we randomly select 50 eventualities. Then for each of the eventualities, we randomly select a negative example by randomly changing a word inside the eventuality with another word of the same POS tag label such that their grammar structure is the same. We also conduct filtering to guarantee the negative examples do not appear in ASER. Last but not least, to show the influence of the conceptualization, we conduct the aforementioned two experiments on both the original ASER before the conceptualization and the final one after the conceptualization.¹⁸

We present the annotation results in Figure 10. The green color indicates the number of eventuality pairs that the more frequent eventuality makes more sense, and the purple color indicates the number of eventualities pairs the less frequent eventuality makes more sense. From the result in Figure 10 (a), we can see that more than 70% of the eventuality pairs as positively correlated, which is consistent with the previous study on the correlation between frequency and selectional preference [13]. At the same time, we also observe that about 30% of the less frequent eventualities are also quite plausible, which is mainly because the frequency of an eventuality is also severely influenced by the rareness of the words inside the eventuality. For example, the eventuality “I eat avocado” appears much less than “I eat apple” because avocado is much rarer than apple rather than “I eat apple” makes more sense than “eat avocado.” The results in Figure 10 (b) help prove that the low-frequent eventualities still contain rich low-order selectional preference because, for more than 90% of the pairs, the randomly extracted pairs in ASER makes more sense than those out of ASER. Furthermore, the experimental results in Figure 10 (c) and (d) show that even though the conceptualization process significantly improves the coverage of ASER, it would not hurt the overall quality. This is mainly because, during the conceptualization step, we carefully design the new weights based on the original weight and the confidence scores provided by Probase [43].

7.3. Relation Extraction

Besides the eventuality extraction, we also care about the extraction quality of the discourse relations between eventualities. For each relation type, we randomly select 50 edges and the corresponding sentences. We generate a question for each pair of them by asking the annotators if they think the extracted discourse relation can represent the correct relation in the original sentence. If so, they should label it as “Valid.” Otherwise, they should label it as “Not Valid.” Similar to the eventuality extraction experiment, we invite six annotators for each edge. If more than four of them agree that the extracted relation is “Valid,” we will consider it to be “Valid.”

From the results in Figure 11, we can see that the overall accuracy is about 80%, which is consistent with the reported performance of the used discourse relations extraction system [14]. Besides that, we also notice that the model performance varies on different relation types. For example, the model tends to perform well on simple types such as “Reason” and “Alternative” because the popular connectives (i.e., “because” and “or”) are less ambiguous. As a comparison, when the connective is more ambiguous (e.g., “while” for “Synchronous”), the overall performance will drop.

7.4. Higher-order Selectional Preference

Finally, we evaluate whether the edge frequency in ASER can be used to reflect human’s high-order selectional preference about eventualities. Similar to the evaluation on the lower-order selectional preference, we conduct two

¹⁸For the experiment on the ASER after the conceptualization, we only sample the conceptualized eventualities and ignore the original ones.

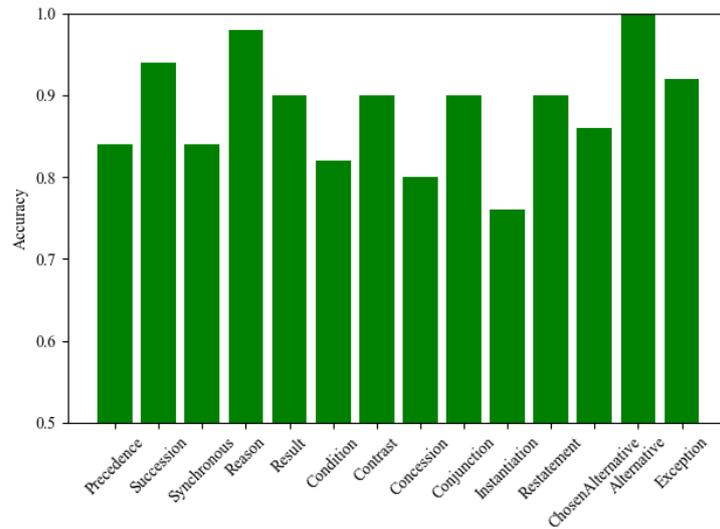


Figure 11: Human annotation of discourse relation extraction quality.

experiments (i.e., (1) High frequency vs. Low Frequency; (2) Exist vs. None-exist) on ASER before and after the conceptualization.¹⁹ For the “High frequency vs. Low Frequency” experiment, we randomly sample 50 edge pairs for each relation type such that the two edges in each pair share the same head eventuality, relation type, but different tail eventuality (e.g., ⟨“I am hungry,” Result, “I eat food”⟩ versus ⟨“I am hungry,” Result, “I exercise”⟩). More importantly, the two sampled edges should have significantly different frequencies. Specifically, we require the frequency of the high-frequency one to be larger than five, and the frequency of the high-frequency one must be at least five times larger than the frequency of the low frequency one. For the “Exist vs. None-exist” experiment, for each relation type, we first randomly sample 50 edges. Then for each edge, we randomly replace a single word of the tail eventuality such that the new tail is very similar to the original one but the created edge does not exist in ASER.

The annotations results are presented in Figure 12. In general, we can make similar observations as the lower SP that the correlation is more significant when we compare the existing and non-existing edges. Besides that, the experiments on the conceptualized ASER help demonstrate that the conceptualization module will not influence the overall quality of edge frequencies.

8. Inference over ASER

In this section, we first introduce two kinds of inferences (eventuality retrieval and relation retrieval) based on ASER. For each of them, inferences over both one-hop and multi-hop are provided. Complexities of these two retrieval algorithms are $O(n^k)$, where n is the number of average adjacent eventualities per eventuality and k is the number of hops. In this section, we show how to conduct these inferences over one-hop and two-hop as the demonstration. ASER is composed of eventualities and concepts. In line with the settings, we conduct case studies over extracted sub-graphs of eventuality and concept graphs. After that, we investigate the rule and meta-path-based inferences on ASER. For the rule-based inference, we leverage AMIE+ [59], a rule mining system on ontological knowledge bases (KBs), to discover closed and connected Horn rules on ASER. For the meta-path-based inference, we obtain the frequent meta-paths using statistical methods and perform case studies by instantiating the meta-paths in both eventuality and concept graphs.

¹⁹For the experiment on the ASER after the conceptualization, we only sample the conceptualized eventualities and ignore the original ones.

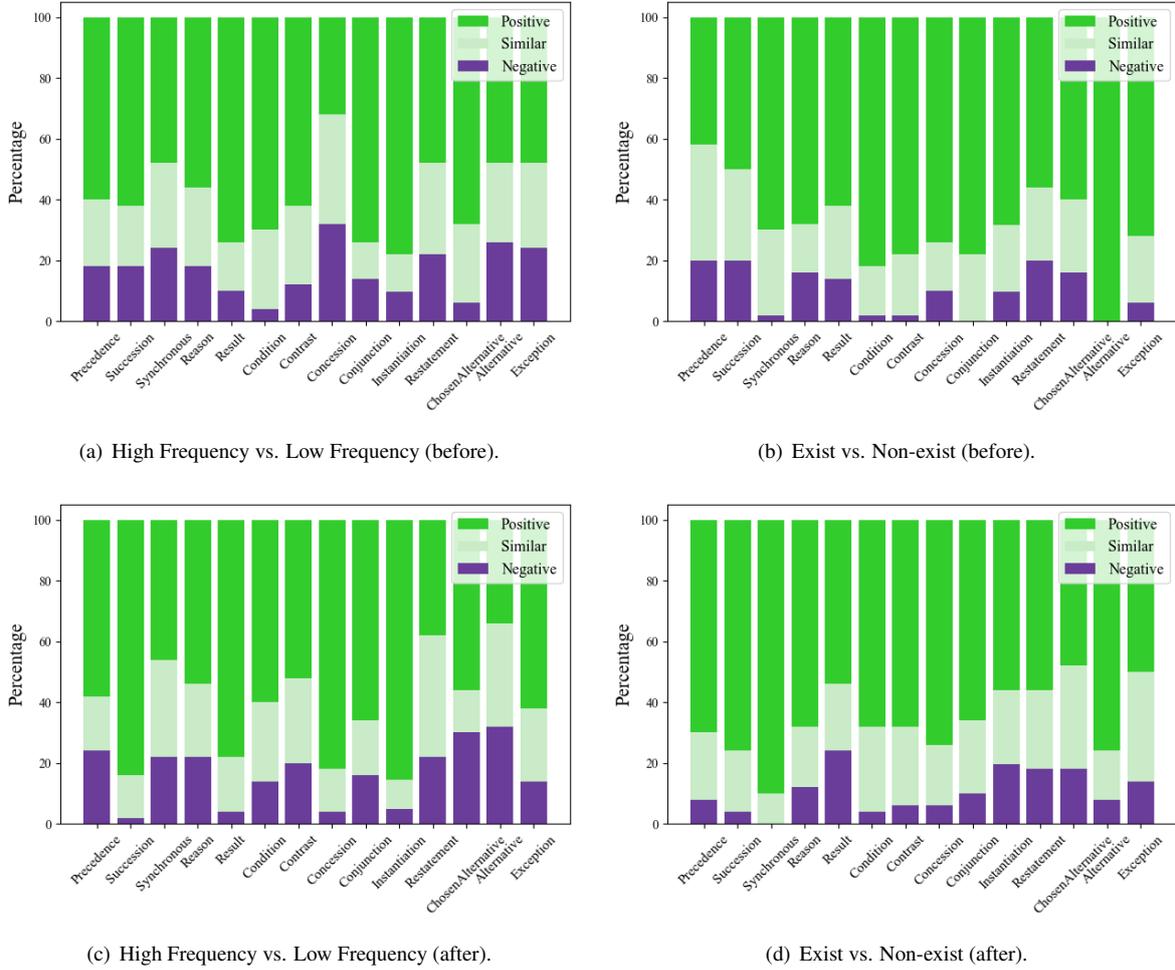


Figure 12: Human annotation of the higher selectional preference in ASER. Experiments on the eventualities before and after the conceptualization are denoted with (before) and (after), respectively. The green color indicates the number of edge pairs that the more frequent edge makes more sense, and the purple color indicates the number of edge pairs the less frequent edge makes more sense.

8.1. Eventuality Retrieval

The eventuality retrieval inference is defined as follows. Given a head eventuality²⁰ E_h and a relation list $\mathcal{L} = (T_1, T_2, \dots, T_k)$, find related eventualities and their associated probabilities such that for each eventuality E_t we can find a path, which contains all the relations in \mathcal{L} in order from E_h to E_t .

8.1.1. One-hop Inference

For the one-hop inference, we assume the target relation is T_1 . We then define the probability of any potential tail node E_t as:

$$\Pr(E_t|E_h, T_1) = \frac{w_{\langle E_h, T_1, E_t \rangle}^{(r)}}{\sum_{E'_t, s.t., \langle E_h, T_1, E'_t \rangle \in \mathcal{R}} w_{\langle E_h, T_1, E'_t \rangle}^{(r)}}, \quad (7)$$

where $w_{\langle E_h, T_1, E_t \rangle}^{(r)}$ is the relation weight, which is defined in Definition 3. If no node is connected with E_h via T_1 , $\Pr(E'|E_h, T_1)$ will be 0 for any $E' \in \mathcal{E}$.

²⁰ASER also supports the prediction of head eventualities given tail eventualities and relations. We omit it in this section for a clear presentation.

Head	Relation	Tail	Probability
You drink alcohol	Synchronous	You drown	0.50
I drink coffee	Result	I calm down	0.33
You are an employee	Contrast	You get fired	0.50
I am programmer	Result	I have free time	1.00
You go to restaurant	Precedence	You get sick	0.50
I am frightened	Reason	Dog barks	0.80
I order chicken	Concession	I am a vegan	1.00
It is my birthday	Result	We go to zoo	0.20
It is a cat	Condition	It is a tiger	0.67
The surgery goes well	Result	There is no complication	0.50

Table 8: Cases of one-hop eventuality inference in the eventuality graph. In the tables of the case study of eventualities, the words in eventualities are stored as lemmas in KBs. However, to clarify the examples, we manually correct the grammar mistakes.

Head	Relation	Tail	Probability
<i>Company</i> be <i>Stakeholder-Group</i>	Condition	<i>PersonX</i> be successful	0.53
<i>PersonX</i> have <i>Issue</i>	Reason	<i>PersonX</i> be proud	0.52
<i>PersonX</i> get <i>Symptom</i>	Synchronous	<i>PersonX</i> be <i>Vulnerable-Group</i>	0.50
<i>PersonX</i> be <i>Emotion</i>	Succession	<i>PersonX</i> marry	0.51
<i>AnimalX</i> bark	Result	<i>AnimalX</i> kill <i>AnimalY</i>	0.33
<i>PersonX</i> be <i>Predator</i>	Result	<i>PersonX</i> tease <i>PersonY</i>	0.25
<i>PersonX</i> do <i>Academic-Misconduct</i>	Contrast	<i>PersonX</i> tell <i>Institute</i>	0.52
<i>PersonX</i> play <i>Sport</i>	Reason	<i>PersonX</i> love <i>Activity</i>	0.27
<i>PersonX</i> hurt <i>Insect</i>	Condition	<i>PersonX</i> help <i>Insect</i>	0.83
<i>PersonX</i> have <i>Social-Medium</i>	Result	<i>PersonX</i> post it	0.72

Table 9: Cases of one-hop eventuality inference in the concept graph. The concepts are marked as *italic* texts.

Several interesting inference examples are observed. In Table 8, we list the reasonable examples of one-hop eventuality inference in the eventuality graph. We also list some of them as follows for discussion:

- ⟨“I drink coffee,” Reason, “I enjoy the flavor”⟩
- ⟨“You go to restaurant,” Precedence, “You got sick”⟩
- ⟨“It is a cat,” Condition, “It is a tiger”⟩

It is observed that “I enjoy the (coffee) flavor” is likely to be the reason for “I drink coffee.” It is also common that if you eat in an unhygienic restaurant, you would probably get sick after you go to the restaurant. Given the fact that the tiger is the largest cat species, it is reasonable to say that if “it is a tiger,” “it is a cat.”

The following examples in Table 9 show the results of one-hop eventuality inference in concept graph.

- ⟨“*Company* be *Stakeholder-Group*,” Condition, “*PersonX* be successful”⟩
- ⟨“*PersonX* hurt *Insect*,” Condition, “*PersonX* help *Insect*”⟩
- ⟨“*PersonX* be *Emotion*,” Succession, “*PersonX* marry”⟩

For instance, if someone is successful, his/her company is likely a big corporation of stakeholders. The second one shows a situation that if an unprofessional person helps insects out of good wills, he/she probably hurts them in reverse. We could also infer from the last case that people tend to be emotional when they get married.

Head	Relation1	Middle	Relation2	Tail	Probability
I go to school	Reason	[I admire]	Synchronous	I am grown up	0.50
I go to bed	Conjunction	[I sleep early]	Result	I am healthy	0.86
We have dinner	Conjunction	[Food is very good]	Contrast	Service is not	0.95
You go to restaurant	Condition	[They do something right]	Reason	There is a line-up	0.50
We have lunch	Conjunction	[We really hit it off]	Contrast	She has a boyfriend at time	0.50
You drink alcohol	Contrast	[You are fine]	Contrast	You have no work	0.75
I am a vegan	Result	[I do not eat fish]	Contrast	We are hungry	0.73
I go to bar	Precedence	[Our table is ready]	Result	We take seats	0.35
I go to restaurant	Reason	[I have a coupon]	Contrast	It is expired	0.36
I go to gym	Precedence	[I go on a date]	Contrast	We have nothing in common	0.25

Table 10: Cases of two-hop eventuality inference in the eventuality graph. In the table, we provide a typical example of middle nodes (embraced by brackets) to create a scenario for better understanding.

Head	Relation1	Middle	Relation2	Tail	Probability
<i>PersonX</i> wait for <i>PersonY</i>	Precedence	[<i>PersonX</i> be tired]	Result	<i>PersonX</i> go to sleep	0.50
<i>PersonX</i> hate <i>Animal</i>	Contrast	[<i>PersonX</i> be harmless]	Contrast	<i>PersonX</i> be <i>Symptom</i>	0.40
<i>PersonX</i> be cranky	Synchronous	[<i>PersonX</i> be hungry]	Result	<i>PersonX</i> order <i>Meat</i>	0.23
<i>PersonX</i> be <i>Artist</i>	Contrast	[<i>PersonX</i> play <i>Sport</i>]	Reason	<i>PersonX</i> be strong	0.33
<i>PersonX</i> regret	Condition	[<i>PersonX</i> despise <i>PersonY</i>]	Reason	<i>PersonY</i> be <i>Performer</i>	0.20
<i>PersonX</i> pull gun	Reason	[<i>PersonX</i> startle]	Synchronous	<i>Domestic-Animal</i> bark	0.50
<i>Predator</i> take down <i>Animal</i>	Reason	[It be <i>Predator</i>]	Synchronous	<i>PersonX</i> shoot	0.32
<i>PersonX</i> be <i>Academic-Title</i>	Result	[<i>PersonX</i> be right]	Contrast	<i>PersonY</i> doubt it	0.28
<i>PersonX</i> hear it	Synchronous	[<i>PersonY</i> play <i>Musical-Instrument</i>]	Synchronous	<i>PersonY</i> be blue	0.65
<i>PersonX</i> be <i>Artist</i>	Condition	[<i>PersonX</i> strike <i>PersonY</i>]	Synchronous	<i>PersonY</i> interview <i>PersonX</i>	0.40

Table 11: Cases of two-hop eventuality inference in the concept graph. In the table, we provide a typical example of middle nodes (embraced by brackets) to create a scenario for better understanding. The concepts are marked as *italic* texts.

8.1.2. Two-hop Inference

On top of Eq. (7), it is easy for us to define the probability of E_t on two-hop setting. Assume the two relations are T_1 and T_2 in order. We can define the probability as follows:

$$\Pr(E_t|E_h, T_1, T_2) = \sum_{E_m \in \mathcal{E}_m} \Pr(E_m|E_h, T_1) \Pr(E_t|E_m, T_2), \quad (8)$$

where \mathcal{E}_m is the set of intermediate node E_m such that (E_h, T_1, E_m) and $(E_m, T_2, E_t) \in \mathcal{R}$.

We list the intuitive examples of two-hop eventuality inference in eventuality and concept graph in Table 10 and Table 11. To better understand the two relations between the head node and the tail node, a typical middle node embraced by brackets is provided. In the eventuality graph, three examples in Table 10 are given for further explanation.

- ⟨“I go to bed,” Conjunction, [“I sleep early”], Result, “I am healthy”⟩
- ⟨“We have lunch,” Conjunction, [“We really hit it off”], Contrast, “She has a boyfriend at time”⟩
- ⟨“I go to restaurant,” Reason, [“I have a coupon”], Contrast, “It is expired”⟩

The first example illustrates that “I go bed” and “I sleep early” tend to result in “I am healthy.” The second one describes a common social situation that I have lunch with a girl and we really hit it off. But she has a boyfriend at that time. Also, it is inferred that the reason why I go to that restaurant is that I have a coupon. However, I find out that the coupon is expired.

Leveraging the same method, we perform two-hop eventuality inference in the concept graph and the results are presented in Table 11.

- ⟨“*PersonX* wait for *PersonY*,” Precedence, [“*PersonX* be tired”], Result, “*PersonX* go to sleep”⟩

Relation	Head	Tail	Probability
Result	You drink alcohol	You have to pee	1.00
Result	I drink coffee	I order a cappuccino	0.50
Alternative	You are a employee	You will be fired	0.50
Contrast	I eat meat	I am not a steak lover	1.00
Precedence	You go to sleep	you wake up	1.00
Contrast	I go to school	I drop out	0.50
Reason	I am not picky	I go to restaurant	0.43
Contrast	I love to cook	I go to restaurant	0.57
Precedence	He waves his hat	The train stops	1.00
Concession	I go to gym	I am tired	0.83

Table 12: Cases of one-hop relation inference in the eventuality graph.

Relation	Head	Tail	Probability
Condition	<i>Company be Stakeholder-Group</i>	<i>PersonX do Local-Ad</i>	0.10
Contrast	<i>PersonX call Agency</i>	<i>It take Duration</i>	0.18
Reason	<i>PersonX be Public-Figure</i>	<i>PersonX be professional</i>	0.30
Synchronous	<i>Animal bite</i>	<i>Animal be frightened</i>	0.23
Precedence	<i>PersonX be Vulnerable-Group</i>	<i>PersonX quit Activity</i>	0.88
Synchronous	<i>It be Domestic-Animal</i>	<i>It be Mammal</i>	0.77
Contrast	<i>Bird catch Animal</i>	<i>Animal get cheese</i>	0.67
Synchronous	<i>PersonX whistle</i>	<i>Animal bark</i>	1.00
Condition	<i>PersonX give lecture</i>	<i>PersonX be Academic-Title</i>	0.21
Result	<i>PersonX play Sport</i>	<i>PersonX be fit</i>	0.87

Table 13: Cases of one-hop relation inference in the concept graph. The concepts are marked as *italic* texts.

- $\langle \text{“PersonX be cranky,” Synchronous, [“PersonX be hungry”], Result, “PersonX order Meat”} \rangle$
- $\langle \text{“PersonX be Artist,” Condition, [“PersonX strike PersonY”], Synchronous, “PersonY interview PersonX”} \rangle$

An interesting example shows that someone is waiting for his/her friend for such a long time that he/she is tired and decides to go to sleep. We also observe that the result of someone being cranky and hungry is most likely to be that he/she orders meats. The last one shows that the artifacts of an artist *PersonX* strikes *PersonY* and it happens at the same time as *PersonY* interviews with *PersonX*.

8.2. Relation Retrieval

The relation retrieval inference is defined as follows. Given two nodes E_h and E_t , find all relation lists and their probabilities such that for each relation list $\mathcal{L} = (T_1, T_2, \dots, T_k)$, we can find a path from E_h to E_t , which contains all the relations in \mathcal{L} in order.

8.2.1. One-hop Inference

Assuming that the path length is one, we define the probability of one relation $R = \langle E_h, T, E_t \rangle$ given E_h and E_t as:

$$\Pr(R|E_h, E_t) = \Pr(T|E_h, E_t) = \frac{w_{\langle E_h, T, E_t \rangle}^{(r)}}{\sum_{T' \in \mathcal{T}} w_{\langle E_h, T', E_t \rangle}^{(r)}}, \quad (9)$$

where \mathcal{T} is the relation type set.

In Table 12 and 13, we perform one-hop relation inference in eventuality and concept graph separately. The relations between head nodes and tail nodes are retrieved to present the commonsense in daily life. In Table 12, the eventuality relations are mined to show frequent patterns in eventuality graph.

- \langle Result, “I drink coffee,” “I order cappuccino” \rangle
- \langle Contrast, “I love to cook,” “I go to restaurant” \rangle
- \langle Concession, “I go to gym,” “I am tired” \rangle

For example, “You drink alcohol” usually leads to “You have to pee.” Another intriguing case illustrates that if someone loves to cook, he/she tends not to go to the restaurant regularly. The last one describes a diligent and determined person who decides to go to the gym, although he/she is tired.

In the concept graph, the same process is used and the results are stored in Table 13.

- \langle Contrast, “Bird catch Animal,” “Animal get cheese” \rangle
- \langle Condition, “PersonX give lecture,” “PersonX be Academic-Title” \rangle
- \langle Result, “PersonX play Sport,” “PersonX be fit” \rangle

We learn from the first example that if birds do not catch these animals (e.g., rats), they would probably get cheese. The second one shows that if someone has an academic title (e.g., professor), he/she will deliver a lecture. The third example tells that the result of *PersonX* plays sport is he/she is fit.

8.2.2. Two-hop Inference

Similarly, given two nodes E_h and E_t , we define the probability of a two-hop connection (T_1, T_2) between them as follows:

$$\begin{aligned} \Pr(T_1, T_2 | E_h, E_t) &= \sum_{E_m \in \mathcal{E}_m} \Pr(T_1, T_2, E_m | E_h, E_t) \\ &= \sum_{E_m \in \mathcal{E}_m} \Pr(T_1 | E_h) \Pr(E_m | T_1, E_h) \Pr(T_2 | E_m, E_t), \end{aligned} \quad (10)$$

where $\Pr(T | E_h)$ is the probability of a relation type T given a head eventuality E_h , which is defined as follows:

$$\Pr(T | E_h) = \frac{\sum_{E_t, s.t., (E_h, T, E_t) \in \mathcal{R}} W_{\langle E_h, T, E_t \rangle}^{(r)}}{\sum_{T' \in \mathcal{T}} \sum_{E_t, s.t., (E_h, T', E_t) \in \mathcal{T}} W_{\langle E_h, T', E_t \rangle}^{(r)}}. \quad (11)$$

The two-hop relations are inferred from the eventuality and concept graph in Table 14 and 15. Some reasonable results are listed below. In line with two-hop eventuality inference, we give a typical middle node embraced by brackets to show the circumstance more clearly.

- \langle Synchronous, Conjunction, “We have breakfast,” [“Our room is ready”], “The front desk staff is friendly” \rangle
- \langle Synchronous, Reason, “I sit on chair,” [“I get my hair washed”], “Stylist tells me” \rangle
- \langle Reason, Result, “I go to supermarket,” [“I have a coupon”], “The price is great” \rangle

In the eventuality graph, we find that some tourists visit a hotel. The guests have breakfast when their room is cleaned and ready. Meanwhile, they find the staff at the front desk is friendly and nice. The second example shows a common thing at the haircut salons. “I sit on chair” to get my hair washed because the stylist tells me to do so before haircut. We also find that someone goes to a supermarket because he/she has a coupon to lower the prices of groceries.

In Table 15, the two-hop relations among head, middle, and tail nodes are extracted to show some insights behind the relation inference.

- \langle Precedence, Precedence, “PersonX wait for PersonY,” [“PersonY send Information”], “PersonZ drag PersonX away” \rangle

Relation1	Relation2	Head	Middle	Tail	Probability
Conjunction	Synchronous	I go to bed	[I would sleep]	I heal	0.25
Result	Reason	I go to bed	[I fall right to sleep]	I am drunk	0.30
Synchronous	Conjunction	We have breakfast	[Our room is ready]	The front desk staffs are friendly	0.12
Contrast	Precedence	You are an employee	[You get fired]	The contract ends	0.11
Reason	Conjunction	I drink coffee	[I have severe ADHD]	I do not get any help	0.33
Synchronous	Reason	I sit on chair	[I get my hair wash]	Stylist tell me	0.25
Contrast	Reason	I am a vegan	[I eat meat]	It tastes good	0.33
Reason	Result	I go to supermarket	[I have a coupon]	The price is great	0.23
Reason	Contrast	I go to restaurant	[The service is great]	The food is mediocre	0.49
Contrast	Contrast	The surgery goes well	[She is in a coma]	She is stabilized	0.53

Table 14: Cases of two-hop relation inference in the eventuality graph. In the table, we provide a typical example of middle nodes (embraced by brackets) to create a scenario for better understanding.

Relation1	Relation2	Head	Middle	Tail	Probability
Precedence	Precedence	<i>PersonX be Stakeholder</i>	[<i>PersonX tell PersonX</i>]	<i>PersonX sign Document</i>	0.85
Precedence	Precedence	<i>PersonX wait for PersonX</i>	[<i>PersonX send Information</i>]	<i>PersonX drag PersonX away</i>	0.40
Result	Precedence	<i>PersonX be Vulnerable-Population</i>	[<i>PersonX be homeless</i>]	<i>Organization help PersonX</i>	0.73
Precedence	Result	<i>Predator catch Herbivore</i>	[<i>Predator eat Meat</i>]	<i>Mammal live</i>	0.79
Synchronous	Result	<i>PersonX be thirsty</i>	[<i>PersonX be hungry</i>]	<i>PersonX order Meat</i>	0.56
Contrast	Condition	<i>PersonX be Musician</i>	[<i>PersonX play Sport</i>]	<i>PersonX be tall</i>	0.42
Result	Synchronous	<i>PersonX excite</i>	[<i>PersonX imagine Emotion</i>]	<i>PersonX answer Electronic-Device</i>	0.58
Reason	Contrast	<i>PersonX eat Animal-Product</i>	[<i>PersonX enjoy it</i>]	<i>It be spicy</i>	0.25
Synchronous	Condition	<i>PersonX hear it</i>	[<i>PersonX play Musical-Instrument</i>]	<i>PersonX be Extracurricular-Activity</i>	0.16
Alternative	Result	<i>PersonX eat Meat</i>	[<i>PersonX be vegetarian</i>]	<i>PersonX starve</i>	0.50

Table 15: Cases of two-hop relation inference in the concept graph. In the table, we provide a typical example of middle nodes (embraced by brackets) to create a scenario for better understanding. The concepts are marked as *italic* texts.

- $\langle \text{Result, Synchronous, "PersonX excite," ["PersonX imagine Emotion"], "PersonX answer Electronic-Device"} \rangle$
- $\langle \text{Synchronous, Condition, "PersonX hear it," ["PersonX play Musical-Instrument"], "PersonX be Extracurricular-Activity"} \rangle$

In the first example, *PersonX* waits for *PersonY* before *PersonY* sends information (e.g., a letter) to inform *PersonX* not wait for him. *PersonX* is reluctant to leave until his/her friends drag him/her away. In the latter example, the result of a person being excited is that he/she imagines the situation and answers the phone. The third case shows that *PersonX* hears that *PersonY* plays an instrument since *PersonY* is having extracurricular activities.

8.3. Rule Mining

AMIE+ [59] aims at mining close and connected Horn Rules in the form of

$$\langle E_a, T_1, E_b \rangle \wedge \langle E_b, T_2, E_c \rangle \Rightarrow \langle E_a, T_3, E_b \rangle$$

where E_a, E_b, E_c and T_1, T_2, T_3 are eventuality variables and relation variables, respectively. As demonstrated through experiments on different KBs, AMIE+ provides an effective approach for the investigation and inference over KB from a logic-rule perspective. We therefore apply AMIE+ on ASER to probe whether it preserves logical properties among multi-hop eventualities and relations. To make the notation consistent with AMIE+, we denote a fact triple with variables at head and/or tail eventuality positions, such as $\langle E_a, T_1, E_b \rangle$, as an *atom*. A rule of interests comprise of a body of a set of atoms B_1, B_2, \dots, B_n ($n = 2$ in this case), and a head of a single atom $\langle E_h, T, E_t \rangle$. We abbreviate the rule as $\vec{B} \Rightarrow \langle E_h, T, E_t \rangle$

The Algorithm of AMIE+ adopts an iterative procedure: starting with a queue of all possible items (rule of size 1), it dequeues a rule in each iteration, and outputs/grows/prunes the rule with certain criterion (listed below), and enqueues the grown new rules back. Throughout the iteration, AMIE+ sets the following significance criterion:

Rule	$\langle E_b \xrightarrow{\text{Concession}} E_f \rangle \wedge \langle E_a \xrightarrow{\text{Result}} E_f \rangle \Rightarrow \langle E_a \xrightarrow{\text{Contrast}} E_b \rangle$
Instances	$\langle \text{I do not know} \rightarrow \text{I guess} \rangle \wedge \langle \text{I believe} \rightarrow \text{I guess} \rangle \Rightarrow \langle \text{I believe} \rightarrow \text{I do not know} \rangle$ $\langle \text{I am not sure} \rightarrow \text{I guess} \rangle \wedge \langle \text{I hope so} \rightarrow \text{I guess} \rangle \Rightarrow \langle \text{I hope so} \rightarrow \text{I am not sure} \rangle$ $\langle \text{I understand} \rightarrow \text{I can not speak} \rangle \wedge \langle \text{I am not a lawyer} \rightarrow \text{I can not speak} \rangle \Rightarrow \langle \text{I am not a lawyer} \rightarrow \text{I understand} \rangle$
Rule	$\langle E_f \xrightarrow{\text{Contrast}} E_b \rangle \wedge \langle E_a \xrightarrow{\text{Instantiation}} E_f \rangle \Rightarrow \langle E_a \xrightarrow{\text{Contrast}} E_b \rangle$
Instances	$\langle \text{I remember} \rightarrow \text{I could not find it} \rangle \wedge \langle \text{I get} \rightarrow \text{I remember} \rangle \Rightarrow \langle \text{I get} \rightarrow \text{I could not find it} \rangle$ $\langle \text{I would say} \rightarrow \text{I might be wrong} \rangle \wedge \langle \text{I hope} \rightarrow \text{I would say} \rangle \Rightarrow \langle \text{I hope} \rightarrow \text{I might be wrong} \rangle$ $\langle \text{It have been suggested} \rightarrow \text{This is unlikely} \rangle \wedge \langle \text{It is possible} \rightarrow \text{It have been suggested} \rangle \Rightarrow \langle \text{It is possible} \rightarrow \text{This is unlikely} \rangle$
Rule	$\langle E_e \xrightarrow{\text{ChosenAlternative}} E_b \rangle \wedge \langle E_a \xrightarrow{\text{ChosenAlternative}} E_e \rangle \Rightarrow \langle E_a \xrightarrow{\text{ChosenAlternative}} E_b \rangle$
Instances	$\langle \text{I will not go} \rightarrow \text{You come here} \rangle \wedge \langle \text{I want to see} \rightarrow \text{I will not go} \rangle \Rightarrow \langle \text{I want to see} \rightarrow \text{You come here} \rangle$ $\langle \text{I want} \rightarrow \text{It is} \rangle \wedge \langle \text{I wish} \rightarrow \text{I want} \rangle \Rightarrow \langle \text{I wish} \rightarrow \text{It is} \rangle$ $\langle \text{I want} \rightarrow \text{I get} \rangle \wedge \langle \text{I do not get that} \rightarrow \text{I want} \rangle \Rightarrow \langle \text{I do not get that} \rightarrow \text{I get} \rangle$
Rule	$\langle E_a \xrightarrow{\text{Reason}} E_e \rangle \wedge \langle E_e \xrightarrow{\text{Restatement}} E_b \rangle \Rightarrow \langle E_a \xrightarrow{\text{Reason}} E_b \rangle$
Instances	$\langle \text{I have ever see} \rightarrow \text{I know} \rangle \wedge \langle \text{I know} \rightarrow \text{They are} \rangle \Rightarrow \langle \text{I have ever see} \rightarrow \text{They are} \rangle$ $\langle \text{I am curious} \rightarrow \text{I think} \rangle \wedge \langle \text{I think} \rightarrow \text{It seems} \rangle \Rightarrow \langle \text{I am curious} \rightarrow \text{It seems} \rangle$ $\langle \text{It is not} \rightarrow \text{You are lying} \rangle \wedge \langle \text{You are lying} \rightarrow \text{I do not believe you} \rangle \Rightarrow \langle \text{It is not} \rightarrow \text{I do not believe you} \rangle$
Rule	$\langle E_a \xrightarrow{\text{Concession}} E_f \rangle \wedge \langle E_b \xrightarrow{\text{Reason}} E_f \rangle \Rightarrow \langle E_a \xrightarrow{\text{Contrast}} E_b \rangle$
Instances	$\langle \text{I have no clue} \rightarrow \text{I hope} \rangle \wedge \langle \text{It be} \rightarrow \text{I hope} \rangle \Rightarrow \langle \text{I have no clue} \rightarrow \text{It is} \rangle$ $\langle \text{I reckon} \rightarrow \text{I do not know} \rangle \wedge \langle \text{I can not talk about it} \rightarrow \text{I do not know} \rangle \Rightarrow \langle \text{I reckon} \rightarrow \text{I can not talk about it} \rangle$ $\langle \text{You do not understand it} \rightarrow \text{You are admitted} \rangle \wedge \langle \text{That is} \rightarrow \text{You are admitted} \rangle \Rightarrow \langle \text{You do not understand it} \rightarrow \text{That is} \rangle$
Rule	$\langle E_b \xrightarrow{\text{Alternative}} E_f \rangle \wedge \langle E_a \xrightarrow{\text{Result}} E_f \rangle \Rightarrow \langle E_a \xrightarrow{\text{Contrast}} E_b \rangle$
Instances	$\langle \text{I am going} \rightarrow \text{I am not going} \rangle \wedge \langle \text{I do not care} \rightarrow \text{I am not going} \rangle \Rightarrow \langle \text{I do not care} \rightarrow \text{I am going} \rangle$ $\langle \text{You do} \rightarrow \text{I do} \rangle \wedge \langle \text{I suppose} \rightarrow \text{I do} \rangle \Rightarrow \langle \text{I suppose} \rightarrow \text{You do} \rangle$ $\langle \text{I reckon} \rightarrow \text{I guess} \rangle \wedge \langle \text{I wonder} \rightarrow \text{I guess} \rangle \Rightarrow \langle \text{I wonder} \rightarrow \text{I reckon} \rangle$
Rule	$\langle E_a \xrightarrow{\text{Reason}} E_f \rangle \wedge \langle E_b \xrightarrow{\text{Succession}} E_f \rangle \Rightarrow \langle E_a \xrightarrow{\text{Reason}} E_b \rangle$
Instances	$\langle \text{I ask} \rightarrow \text{I am not sure} \rangle \wedge \langle \text{I do not know} \rightarrow \text{I am not sure} \rangle \Rightarrow \langle \text{I ask} \rightarrow \text{I do not know} \rangle$ $\langle \text{We are lucky} \rightarrow \text{We notice} \rangle \wedge \langle \text{We order} \rightarrow \text{We notice} \rangle \Rightarrow \langle \text{We are lucky} \rightarrow \text{We order} \rangle$ $\langle \text{I remember it} \rightarrow \text{I see it} \rangle \wedge \langle \text{I realize} \rightarrow \text{I see it} \rangle \Rightarrow \langle \text{I remember it} \rightarrow \text{I realize} \rangle$
Rule	$\langle E_a \xrightarrow{\text{Concession}} E_f \rangle \wedge \langle E_b \xrightarrow{\text{Precedence}} E_f \rangle \Rightarrow \langle E_a \xrightarrow{\text{Contrast}} E_b \rangle$
Instances	$\langle \text{I am unconscious} \rightarrow \text{I wake up} \rangle \wedge \langle \text{I see} \rightarrow \text{I wake up} \rangle \Rightarrow \langle \text{I am unconscious} \rightarrow \text{I see} \rangle$ $\langle \text{I swear} \rightarrow \text{I guess} \rangle \wedge \langle \text{I do not know} \rightarrow \text{I guess} \rangle \Rightarrow \langle \text{I swear} \rightarrow \text{I do not know} \rangle$ $\langle \text{I can not believe} \rightarrow \text{It is great} \rangle \wedge \langle \text{I think} \rightarrow \text{It is great} \rangle \Rightarrow \langle \text{I can not believe} \rightarrow \text{I think} \rangle$
Rule	$\langle E_a \xrightarrow{\text{Alternative}} E_e \rangle \wedge \langle E_e \xrightarrow{\text{Exception}} E_b \rangle \Rightarrow \langle E_a \xrightarrow{\text{Exception}} E_b \rangle$
Instances	$\langle \text{It is not} \rightarrow \text{It is wrong} \rangle \wedge \langle \text{It is wrong} \rightarrow \text{It is} \rangle \Rightarrow \langle \text{It is not} \rightarrow \text{It is} \rangle$ $\langle \text{I really want} \rightarrow \text{I think} \rangle \wedge \langle \text{I think} \rightarrow \text{I know} \rangle \Rightarrow \langle \text{I really want} \rightarrow \text{I know} \rangle$ $\langle \text{It is not} \rightarrow \text{I suppose} \rangle \wedge \langle \text{I suppose} \rightarrow \text{You know} \rangle \Rightarrow \langle \text{It is not} \rightarrow \text{You know} \rangle$
Rule	$\langle E_a \xrightarrow{\text{ChosenAlternative}} E_f \rangle \wedge \langle E_b \xrightarrow{\text{ChosenAlternative}} E_f \rangle \Rightarrow \langle E_a \xrightarrow{\text{Restatement}} E_b \rangle$
Instances	$\langle \text{I am hoping} \rightarrow \text{We get} \rangle \wedge \langle \text{I think} \rightarrow \text{We get} \rangle \Rightarrow \langle \text{I am hoping} \rightarrow \text{I think} \rangle$ $\langle \text{I suppose} \rightarrow \text{He is} \rangle \wedge \langle \text{I think} \rightarrow \text{He is} \rangle \Rightarrow \langle \text{I suppose} \rightarrow \text{I think} \rangle$ $\langle \text{I am glad} \rightarrow \text{I think} \rangle \wedge \langle \text{The food is good} \rightarrow \text{I think} \rangle \Rightarrow \langle \text{I am glad} \rightarrow \text{The food is good} \rangle$

Table 16: Cases of AMIE+ rule mining in the eventuality graph. For the simplicity of formatting, we represent $\langle E_h, T, E_t \rangle$ triples as $\langle E_h \xrightarrow{T} E_t \rangle$.

Head coverage:

$$hc(\vec{B} \Rightarrow \langle E_h, T, E_t \rangle) := \frac{\text{supp}(\vec{B} \Rightarrow \langle E_h, T, E_t \rangle)}{\text{size}(T)},$$

where

$$\text{supp}(\vec{B} \Rightarrow \langle E_h, T, E_t \rangle) := \#(E_h, E_t) : \exists z_1, \dots, z_m : \vec{B} \wedge \langle E_h, T, E_t \rangle :$$

denotes the support, i.e., the number of correct prediction yielded with the rule in the current KB, and $\text{size}(T)$ denotes the number of facts with T as relations.

Standard confidence:

$$\text{conf}(\vec{B} \Rightarrow \langle E_h, T, E_t \rangle) := \frac{\text{supp}(\vec{B} \Rightarrow \langle E_h, T, E_t \rangle)}{\#(E_h, E_t) : \exists z_1, \dots, z_m : \vec{B}}$$

where $\#(E_h, E_t) : \exists z_1, \dots, z_m : \vec{B}$ denotes all possible predictions of the rule.

PCA confidence:

$$\text{conf}_{\text{pca}}(\vec{B} \Rightarrow \langle E_h, T, E_t \rangle) := \frac{\text{supp}(\vec{B} \Rightarrow \langle E_h, T, E_t \rangle)}{\#(E_h, E_t) : \exists z_1, \dots, z_m : \vec{B} \wedge \langle E_h, T, E_t \rangle},$$

where $\#(E_h, E_t) : \exists z_1, \dots, z_m : \vec{B} \wedge \langle E_h, T, E_t \rangle$ denotes the number of pairs of (E_h, E_t) predicted by corresponding relation body \vec{B} but with an existing pair of $\langle E_h, T, E_t \rangle$ in KB.

AMIE+ focuses on RDF Knowledge Bases, where an RDF KB could be represented as a set of facts in the form $\langle \text{Subject}, \text{Relation}, \text{Object} \rangle$. To pair ASER Graph with AMIE+, we extract all $\langle E_1, \text{Relation}, E_2 \rangle$ triples from the relation table as our set of facts. To preserve the frequency information, we duplicate each extracted fact for f times where f is the corresponding triple frequency in the relation table. During our experiments, we set the threshold of minimal PCA confidence $\text{minPCA}=0.1$ and minimal head coverage $\text{minHC}=0.01$, and run AMIE+ on both the eventuality graph and concept graph of ASER.

Some of the mined rules and instantiated cases are shown in Table 16 and 17, respectively. For the eventuality graph, the rule “ $\langle E_b \xrightarrow{\text{Concession}} E_f \rangle \wedge \langle E_a \xrightarrow{\text{Result}} E_f \rangle \Rightarrow \langle E_a \xrightarrow{\text{Contrast}} E_b \rangle$ ” demonstrates that if the result of some event is opposite to what should happen, then the reason that induces this “opposite” result is likely to contrast the original event in some core properties. The last instance “ $\langle \text{I understand} \rightarrow \text{I can not speak} \rangle \wedge \langle \text{I am not a lawyer} \rightarrow \text{I can not speak} \rangle \Rightarrow \langle \text{I am not a lawyer} \rightarrow \text{I understand} \rangle$ ” illustrates this rule with a lawsuit scenario where the “opposite” result “I can not speak” is induced by the reason “I am not a lawyer,” which contrasts the original event “I understand.” And the rule “ $\langle E_a \xrightarrow{\text{ChosenAlternative}} E_f \rangle \wedge E_b \langle \xrightarrow{\text{ChosenAlternative}} E_f \rangle \Rightarrow \langle E_a \xrightarrow{\text{Restatement}} E_b \rangle$ ” captures that if the subject consistently prefers two events as alternative to a third event, then it is possible that the first two events have very similar semantic meaning. This is illustrated with its last instance “ $\langle \text{I am glad} \rightarrow \text{I think} \rangle \wedge \langle \text{The food is good} \rightarrow \text{I think} \rangle \Rightarrow \langle \text{I am glad} \rightarrow \text{The food is good} \rangle$ ” where both “I am glad” and “The food is good” express similar meanings in which the subject consistently prefers to “I think.”

For the concept graph, the rule “ $\langle E_e \xrightarrow{\text{Instantiation}} E_a \rangle \wedge \langle E_e \xrightarrow{\text{Instantiation}} E_b \rangle \Rightarrow \langle E_a \xrightarrow{\text{Conjunction}} E_b \rangle$ ” shows that the concepts that both describe a shared concept in details are likely to both happen. Its first instance “ $\langle \text{PersonX realize} \rightarrow \text{PersonX point out} \rangle \wedge \langle \text{PersonX realize} \rightarrow \text{PersonX have Information} \rangle \Rightarrow \langle \text{PersonX point out} \rightarrow \text{PersonX have Information} \rangle$ ” demonstrates this through a information-capture process, where the two concepts “PersonX point out” and “PersonX have Information” that both describe the concept “PersonX realize” happen together. The rule “ $\langle E_e \xrightarrow{\text{Exception}} E_b \rangle \wedge \langle E_e \xrightarrow{\text{Succession}} E_a \rangle \Rightarrow \langle E_a \xrightarrow{\text{Contrast}} E_b \rangle$ ” shows the association between Exception and Contrast that is bridged via Succession. This is demonstrated with its first instance “ $\langle \text{Item be ready} \rightarrow \text{PersonX wait} \rangle \wedge \langle \text{Item be ready} \rightarrow \text{PersonX check} \rangle \Rightarrow \langle \text{PersonX check} \rightarrow \text{PersonX be wait} \rangle$ ” that shows a scenario of a customer checking and waiting for products.

8.4. Meta-path Mining

ASER is a complex heterogeneous graph that encodes the commonsense knowledge. ASER is composed of two types of nodes (i.e., extracted eventuality and conceptualized eventuality) and 15 types of edges (e.g., Reason and Precedence). We leverage meta-path [60] mining which studies the semantic meanings behind paths to tackle the heterogeneity of ASER. A meta-path is a path that consists of a sequence of different relations defined among various node types. Formally, a meta-path P is defined as a path $E_1 \xrightarrow{T_1} E_2 \xrightarrow{T_2} \dots \xrightarrow{T_{l-1}} E_l$, in which $T = T_1 \circ T_2 \circ \dots \circ T_l$ is the composite relation between N_1 and N_l . Take an example from Table 18, the meta-path “ $E_1 \xrightarrow{\text{Conceptualization}} C_1 \xrightarrow{\text{ConceptInstantiation}} E_2$ ” defines a composite relation in which the two eventuality E_1 and E_2 are conceptualized to the same concept C_1 .

To automatically select the most frequent and influential meta-paths, we first perform a random walk on the hybrid graph. Specifically, 50,000 seed nodes are chosen independently and uniformly from the nodes of ASER. Starting

Rule	$\langle E_e \xrightarrow{\text{Restatement}} E_a \rangle \wedge \langle E_e \xrightarrow{\text{Restatement}} E_b \rangle \Rightarrow \langle E_a \xrightarrow{\text{Conjunction}} E_b \rangle$
Instances	$\langle \text{PersonX laugh} \rightarrow \text{PersonX smile} \rangle \wedge \langle \text{PersonX laugh} \rightarrow \text{PersonX open Facial-Feature} \rangle \Rightarrow \langle \text{PersonX smile} \rightarrow \text{PersonX open Facial-Feature} \rangle$ $\langle \text{PersonX love it} \rightarrow \text{It be good} \rangle \wedge \langle \text{PersonX love it} \rightarrow \text{It be tasty} \rangle \Rightarrow \langle \text{It be good} \rightarrow \text{It be tasty} \rangle$ $\langle \text{PersonX wish} \rightarrow \text{PersonX need} \rangle \wedge \langle \text{PersonX wish} \rightarrow \text{PersonX need} \rangle \Rightarrow \langle \text{PersonX need} \rightarrow \text{PersonX need} \rangle$
Rule	$\langle E_e \xrightarrow{\text{Instantiation}} E_a \rangle \wedge \langle E_e \xrightarrow{\text{Instantiation}} E_b \rangle \Rightarrow \langle E_a \xrightarrow{\text{Conjunction}} E_b \rangle$
Instances	$\langle \text{PersonX realize} \rightarrow \text{PersonX point out} \rangle \wedge \langle \text{PersonX realize} \rightarrow \text{PersonX have Information} \rangle \Rightarrow \langle \text{PersonX point out} \rightarrow \text{PersonX have Information} \rangle$ $\langle \text{PersonX have} \rightarrow \text{PersonX get} \rangle \wedge \langle \text{PersonX have} \rightarrow \text{PersonX own} \rangle \Rightarrow \langle \text{PersonX get} \rightarrow \text{PersonX own} \rangle$ $\langle \text{PersonX know} \rightarrow \text{PersonX be sure} \rangle \wedge \langle \text{PersonX know} \rightarrow \text{PersonX remember} \rangle \Rightarrow \langle \text{PersonX be sure} \rightarrow \text{PersonX remember} \rangle$
Rule	$\langle E_e \xrightarrow{\text{Concession}} E_b \rangle \wedge \langle E_e \xrightarrow{\text{Restatement}} E_a \rangle \Rightarrow \langle E_a \xrightarrow{\text{Contrast}} E_b \rangle$
Instances	$\langle \text{PersonX order Dish} \rightarrow \text{PersonX be hungry} \rangle \wedge \langle \text{PersonX order Dish} \rightarrow \text{PersonX order} \rangle \Rightarrow \langle \text{PersonX order} \rightarrow \text{PersonX be hungry} \rangle$ $\langle \text{PersonX wish} \rightarrow \text{PersonX doubt} \rangle \wedge \langle \text{PersonX wish} \rightarrow \text{PersonX need} \rangle \Rightarrow \langle \text{PersonX doubt} \rightarrow \text{PersonX need} \rangle$ $\langle \text{PersonX love it} \rightarrow \text{PersonX hate it} \rangle \wedge \langle \text{PersonX love it} \rightarrow \text{It be good} \rangle \Rightarrow \langle \text{PersonX hate it} \rightarrow \text{It be good} \rangle$
Rule	$\langle E_e \xrightarrow{\text{Exception}} E_b \rangle \wedge \langle E_e \xrightarrow{\text{Succession}} E_a \rangle \Rightarrow \langle E_a \xrightarrow{\text{Contrast}} E_b \rangle$
Instances	$\langle \text{Item be ready} \rightarrow \text{PersonX wait} \rangle \wedge \langle \text{Item be ready} \rightarrow \text{PersonX check} \rangle \Rightarrow \langle \text{PersonX check} \rightarrow \text{PersonX be wait} \rangle$ $\langle \text{PersonX say} \rightarrow \text{PersonX be sorry} \rangle \wedge \langle \text{PersonX say} \rightarrow \text{PersonX be surprised} \rangle \Rightarrow \langle \text{PersonX be sorry} \rightarrow \text{PersonX be surprised} \rangle$ $\langle \text{It be} \rightarrow \text{PersonX guess} \rangle \wedge \langle \text{It be} \rightarrow \text{It be factor} \rangle \Rightarrow \langle \text{PersonX guess} \rightarrow \text{It be factor} \rangle$
Rule	$\langle E_a \xrightarrow{\text{Restatement}} E_f \rangle \wedge \langle E_b \xrightarrow{\text{Restatement}} E_f \rangle \Rightarrow \langle E_a \xrightarrow{\text{Synchronous}} E_b \rangle$
Instances	$\langle \text{PersonX love it} \rightarrow \text{It be good} \rangle \wedge \langle \text{PersonX feel} \rightarrow \text{It be good} \rangle \Rightarrow \langle \text{PersonX love it} \rightarrow \text{PersonX feel} \rangle$ $\langle \text{It be cool} \rightarrow \text{It be good} \rangle \wedge \langle \text{PersonX think} \rightarrow \text{It be okay} \rangle \Rightarrow \langle \text{It be cool} \rightarrow \text{It be okay} \rangle$ $\langle \text{PersonX like it} \rightarrow \text{It be good} \rangle \wedge \langle \text{PersonX be happy} \rightarrow \text{It be good} \rangle \Rightarrow \langle \text{PersonX like it} \rightarrow \text{PersonX be happy} \rangle$

Table 17: Cases of AMIE+ rule mining in the concept graph. For the simplicity of formatting, we represent $\langle E_h, T, E_t \rangle$ triples as $\langle E_h \xrightarrow{T} E_t \rangle$.

from each seed node, a random walk is used to generate 50 multi-hop paths of different nodes and relations. The nodes in a path are represented by their types rather than their contents. For example, a path “I drink coffee $\xrightarrow{\text{Result}}$ I stay up late” is converted into a meta-path “ $E_1 \xrightarrow{\text{Result}} E_2$.” After collecting meta-paths, we search for the frequent patterns of 2-hop and 3-hop meta-paths. The appearance of meta paths is counted. The 2-hop and 3-hop meta paths are later ranked by their frequencies. The frequent meta paths are selected for the further case study. We list the intriguing instances from these meta-paths in Table 18.

For 2-hop meta paths, the results are very similar to the ones of eventuality/relation retrieval inference. For example, “ $E_1 \xrightarrow{\text{Reason}} E_2 \xrightarrow{\text{Result}} E_3$ ” describes paths following cause and effect relations. A typical instance in daily life is “I am in pain \rightarrow I am alone \rightarrow I sit at bar,” describing a scenario in which a man suffers from loneliness and goes to the bar to numb the pains. In addition to relations among eventualities, the interaction between concepts and eventualities are also discovered by the meta-paths. In the cases of “ $E_1 \xrightarrow{\text{Conceptualization}} C_1 \xrightarrow{\text{ConceptInstantiation}} E_2$,” two semantically distinct eventualities are unified in the concept-level. For example, “He is psychiatrist” and “I am attorney” follows the same pattern, “PersonX be specialist.”

For 3-hop meta paths, the reasoning paths are longer and illustrates the daily life in more details. For example, in the meta-path “ $E_1 \xrightarrow{\text{Result}} E_2 \xrightarrow{\text{Contrast}} E_3 \xrightarrow{\text{Conjunction}} E_4$,” an instantiated example, “I have you number \rightarrow I call you \rightarrow I have a meeting \rightarrow I have a presentation,” shows that a person wants to call his friends with the phone number. However, he has to do a presentation in the coming meeting and decides to call his friend later. As for the hybrid meta-path with extracted and conceptualized eventualities, “ $E_1 \xrightarrow{\text{Conjunction}} E_2 \xrightarrow{\text{Conceptualization}} C_1 \xrightarrow{\text{ConceptInstantiation}} E_3$,” we find out that someone is sweating because of the hot weather while someone is unfortunately in a coma. Both of them are unified under the concept “PersonX be Symptom.”

9. From ASER to Commonsense Knowledge

In this section, we investigate the connection between ASER and existing commonsense knowledge bases. Specifically, we check the coverage and similarities between the selectional preference knowledge in ASER and the human-defined commonsense knowledge in ConceptNet [1] and ATOMIC [26].

#Hop	meta-path	Instances
2	$E_1 \xrightarrow{\text{Conjunction}} E_2 \xrightarrow{\text{Contrast}} E_3$	I go to bed → I go to sleep → I wake up I have breakfast → I have milk → I feel sick I take bus → I go to work → I go home
	$E_1 \xrightarrow{\text{Precedence}} E_2 \xrightarrow{\text{Precedence}} E_3$	You go to sleep → You wake up → You hit the ground You drink alcohol → You go to toilet → You have to pee You go to restaurant → You are sick → You go to hospital
	$E_1 \xrightarrow{\text{Conceptualization}} C_1 \xrightarrow{\text{ConceptInstantiation}} E_2$	He is psychiatrist → <i>PersonX</i> is <i>Specialist</i> → I am attorney I want milk → <i>PersonX</i> want <i>Animal-Product</i> → He wants burgers You make reservation → <i>PersonX</i> make <i>Service</i> → He makes statement
	$E_1 \xrightarrow{\text{Conjunction}} E_2 \xrightarrow{\text{Conjunction}} E_3$	I go to gym → I have to wait → I go home I am vegan → My wife is vegan → I used to eat meat It is a cat → It is fine → It is beautiful
	$E_1 \xrightarrow{\text{Reason}} E_2 \xrightarrow{\text{Result}} E_3$	I go to bar → I have many friends → I have parties I go to school → We could afford → I get my first job I am in pain → I am alone → I sit at bar
3	$E_1 \xrightarrow{\text{Precedence}} E_2 \xrightarrow{\text{Conjunction}} E_3 \xrightarrow{\text{Precedence}} E_4$	The rain comes down → The engine whistles → The train starts → The train moves on The moon arises → The weather is pleasant → The snow ceases → The night is still She sleeps → The phone rings → We gets home → She hangs up the phone
	$E_1 \xrightarrow{\text{Conjunction}} E_2 \xrightarrow{\text{Conceptualization}} C_1 \xrightarrow{\text{ConceptInstantiation}} E_3$	I play piano → I am musician → <i>PersonX</i> be <i>Artist</i> → He is actor I am chill → It is a snake → It be <i>Predator</i> → It is a bear It is hot → I am sweating → <i>PersonX</i> be <i>Symptom</i> → She is in a coma
	$E_1 \xrightarrow{\text{Condition}} E_2 \xrightarrow{\text{Reason}} E_3 \xrightarrow{\text{Conjunction}} E_4$	Everyone knows him → He comes off the bench → He makes his debut for club → He scores his first goal I am healthy → I sleep → I am exhausted → I am cold We get the check → We order dessert → I am still hungry → We eat everything
	$E_1 \xrightarrow{\text{Result}} E_2 \xrightarrow{\text{Contrast}} E_3 \xrightarrow{\text{Conjunction}} E_4$	I am tired → I go to bed → The sun is shining → The wind blows There is a storm coming → The rain falls → The sky is clear → The air is warm I have you number → I call you → I have a meeting → I have a presentation
	$E_1 \xrightarrow{\text{Contrast}} E_2 \xrightarrow{\text{Reason}} E_3 \xrightarrow{\text{Reason}} E_4$	I am a vegan → I eat meat → I enjoy it → It tastes good The painting is controversial → It is a masterpiece → It belongs to museum → It is valuable I get over it quickly → I go to mall → I buy clothes → I have a job interview

Table 18: Instances of meta paths generated by random walk. E represents eventuality while C represents concept. The concepts in the instances are marked as *italic* texts. For example, “I go to bed → I go to sleep → I sleep” is an instance of the meta-path “ $E_1 \xrightarrow{\text{Conjunction}} E_2 \xrightarrow{\text{Contrast}} E_3$.”

9.1. Relationship with ConceptNet

After around 20 years development, ConceptNet 5.0 [61] now contains 21 million edges over 8 million nodes, built from the original ConceptNet [1]. The core of ConceptNet, which is inherited from the Open Mind CommonSense (OMCS) project [1], only contains 600K pieces of high-quality commonsense knowledge in the format of tuples, e.g., (‘song’, *UsedFor*, ‘sing’). However, there is a huge gap between the small scale of existing commonsense knowledge resources and the broad demand of downstream applications, motivating us to acquire more commonsense knowledge cheaply and feasibly.

Based on the observation that selectional preference can naturally reflect commonsense knowledge about word choice in various contexts [4], we proposed TransOMCS [11] to transfer the selectional preference knowledge in ASER to ConceptNet-like commonsense tuples. Specifically, we adopt the English subset of ConceptNet 5 [61] as seed commonsense knowledge, and only relations covered by the original OMCS project [1] are selected. Different from OpenIE [62] and Hearst patterns [63], where human-defined patterns are leveraged to extract relations, we develop a pipeline to discover dependency patterns automatically. As shown in Figure 13, for each commonsense relation r in ConceptNet, we first try to find patterns over dependency and discourse relations in ASER automatically from the overlap of ConceptNet assertions and ASER sub-graphs. The percentage of ConceptNet knowledge that can be matched in ASER is presented in Figure 14. After that, a pattern selection scoring function is designed to select highly plausible patterns. We present the most plausible dependency patterns for each relation as an illustration in Table 20. Based on the patterns, we can traverse the whole ASER to acquire a large-scale commonsense knowledge graph in the format of ConceptNet.

A running example of TransOMCS is shown in Figure 13. For the OMCS-like assertion ⟨“Good grades,” Causes, “Graduate”⟩, we can extract the corresponding dependency relation from the ASER edge ⟨“he gets good grades,” Result, “he graduates college”⟩. Such dependency pattern is in turn used for other ASER edges to extract novel knowledge. As a result, we successfully acquire 18 million ConceptNet-like commonsense assertions with high novelty and accuracy. From the case study in Table 21 we can see that ConceptNet-like commonsense knowledge is indeed contained in ASER. For example, with the help of ASER, we can know that students are often at school, artists often create art, and the wall is part of the house.

We also conducted qualitative analysis regarding the accuracy, novelty, and quantity of the acquired commonsense

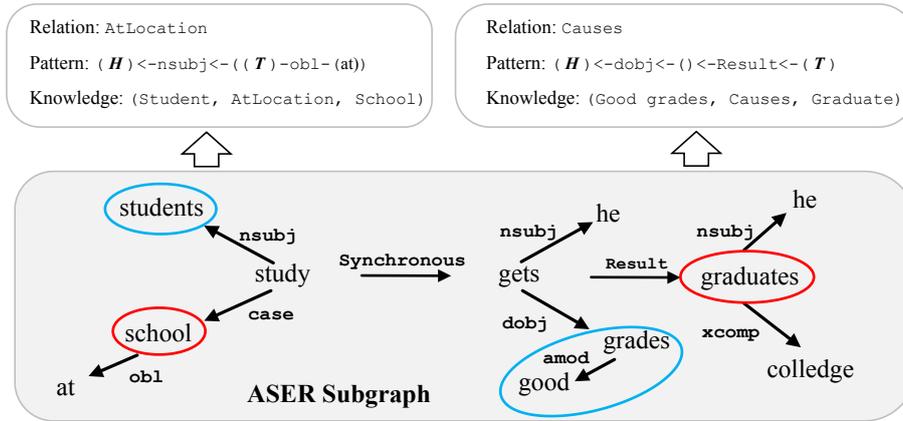


Figure 13: Example of transferring selectional preference knowledge in ASER to commonsense knowledge. By exploring different dependency patterns, we can acquire different forms of candidate knowledge. Extracted head and tail concepts are indicated with blue and red circles respectively, and are denoted as H and T placeholders in the patterns.

Model	# Vocab	# Tuple	Novel _t	Novel _c	ACC _n	ACC _o
COMET _{Original} (Greedy decoding)	715	1,200	33.96%	5.27%	58%	90%
COMET _{Original} (Beam search - 10 beams)	2,232	12,000	64.95%	27.15 %	35 %	44%
COMET _{Extended} (Greedy decoding)	3,912	24,000	99.98%	55.56%	34%	47%
COMET _{Extended} (Beam search - 10 beams)	8,108	240,000	99.98%	78.59%	23%	27%
LAMA _{Original} (Top 1)	328	1,200	-	-	-	49%
LAMA _{Original} (Top 10)	1,649	12,000	-	-	-	20%
LAMA _{Extended} (Top 1)	1,443	24,000	-	-	-	29%
LAMA _{Extended} (Top 10)	5,465	240,000	-	-	-	10%
TransOMCS _{Original} (No Ranking)	33,238	533,449	99.53%	89.20%	72%	74%
TransOMCS (Top 1%)	37,517	184,816	95.71%	75.65%	86%	87%
TransOMCS (Top 10%)	56,411	1,848,160	99.55%	92.17%	69%	74%
TransOMCS (Top 30%)	68,428	5,544,482	99.83%	95.22%	67%	69%
TransOMCS (Top 50%)	83,823	9,240,803	99.89%	96.32%	60%	62%
TransOMCS (No Ranking)	100,659	18,481,607	99.94%	98.30%	54%	56%
OMCS in ConceptNet 5.0	36,954	207,427	-	-	-	92%

Table 19: Main evaluation results of TransOMCS compared with COMET and LAMA.

knowledge in TransOMCS. For accuracy, human annotators from Amazon Mechanical Turk are invited to evaluate whether the 100 randomly sampled generated commonsense triple are plausible or not. If at least four annotators out of five agree that the triple is plausible, then it is considered plausible. ACC_n is the accuracy of the novel triples and ACC_o is the overall accuracy of all triples. For novelty, the proportion of all generated tuples that are novel (Novel_t) and that have a novel object/tail (Novel_c) are used to measure novelty of the generated knowledge. We also include the quantity of generated triples for reference.

For baseline models, we compare COMET [64] and LAMA [10] with TransOMCS. When decoding, head-relation pairs in OMCS are fed into the language models to acquire tail outputs, where this setting is denoted as COMET_{Original} and LAMA_{Original}. Due to the small size of OMCS (around 1.2K), we also include 24K additional head-relation pairs from the concepts extracted by TransOMCS as additional inputs, where the corresponding models are denoted as COMET_{Extended} and LAMA_{Extended}.

The evaluation results of TransOMCS are shown in Table 19 (Results are from Table 2 in the original paper of TransOMCS [11]). In the meantime of producing commonsense knowledge with two more orders of magnitude in terms of quantity, TransOMCS can produce commonsense tails with more novelty and accuracy. For COMET, as it’s a pure machine learning based approach which can fit the training data too well to generate novel tails. For quality, when the test data is similar to the training set, COMET provides the best quality. For example, in the COMET_{Original} setting under greedy decoding, it achieves 90% overall accuracy. The quality of LAMA is whereas less satisfying as a matter of the unsupervised setting and over simple prompts. Compared with them, TransOMCS (top 1%) can generate

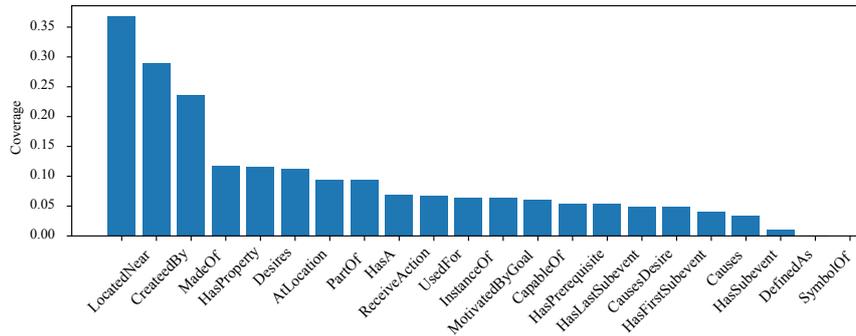


Figure 14: The matching statistics of ConceptNet assertions in ASER grouped by each relation. The coverage indicates the proportion of ConceptNet assertions where both heads and tails can be matched to an ASER unit, i.e., a discourse edge or an eventuality.

Relation	Dependency Pattern	Relation	Dependency Pattern
AtLocation	()->compound->()	HasSubevent	(-dobj-)->neg->()
CapableOf	()<-nsubj<-()	HasFirstSubevent	(-prep-)<-Succession<-()
Causes	(-compound-)<-pobj<-of<-prep<-()	HasLastSubevent	()<-acomp<-be<-Reason<-()
CausesDesire	()<-pobj<-to<-prep<-()	InstanceOf	()<-acomp<-be->nsubj->()
CreatedBy	()<-dobj<-make->nsubj->()	LocatedNear	()<-nsubj<-be->prep->on->pobj->()
DefinedAs	()<-nsubj<-be->attr->(-amod-)	MadeOf	()<-compound<-()
Desires	()<-nsubj<-()	MotivatedByGoal	()<-xcomp<-()
HasA	()<-nsubj<-have->dobj->()	PartOf	()->compound->()
HasPrerequisite	()->dobj->()	ReceivesAction	()<-dobj<-()
HasProperty	()<-nsubj<-be->acomp->()	UsedFor	()<-pobj<-(-prep-)

Table 20: Examples of extracted dependency patterns in TransOMCS. We select the pattern ranked as most plausible for each relation. () are placeholders for words, and attributes like nsubj are names of the dependency edges.

commonsense knowledge with comparable quality as COMET.

9.2. Relationship with ATOMIC

9.2.1. Overlaps

Besides ConceptNet, another substantial commonsense knowledge base is ATOMIC [26], a large-scale human-annotated commonsense knowledge graph that provides inferential knowledge about daily events. Like ASER, the ATOMIC nodes are events described in free-form text, while not parsed to be canonical. There are nine *if-then* relationships defined across ATOMIC, measuring the daily causes and effects for certain base events. To tackle the limitations in terms of novelty and coverage of current *if-then* commonsense acquisition methods, we proposed a novel framework DISCOS (from DIScourse to COMmonSense) [65, 66], which transfers selectional preference knowledge in ASER to complex commonsense knowledge in ATOMIC. As a result, we acquire 3.4 Million *if-then* commonsense knowledge in the format of ATOMIC. An illustration of the process in DISCOS is presented in Figure 15.

Specifically, we first conduct an alignment from ATOMIC to ASER. In ATOMIC, the personal pronouns are represented with wildcards like “PersonX” and “PersonY,” and in ASER, the subjects of events are concrete personal pronouns like “she” and “he.” Moreover, as all of the tail events in ATOMIC are written by human annotators, the form of ATOMIC tails can be arbitrary and sometimes subjects are omitted. Based on those observations, we develop some string substitution rules to align the nodes in ATOMIC and ASER, as illustrated in Table 22. After conducting the string substitution operations, we use the parser in ASER to parse the acquired text into standard ASER format.

Table 23 presents the coverage statistics between ATOMIC and ASER. We first conduct the string match to check the coverage of ATOMIC nodes in ASER, and find that the average percentage of ATOMIC nodes found in ASER is 62.9%. For edges, we present the percentage of ATOMIC edges whose head and tail are both covered by ASER, which is 35.91% on average. On top of the matched edges, we check the shortest path length between the matched head and tail in ASER and report the average among all edges in the *Avg. Shortest Path Length* column. The range of

Head	Relation	Tail	Head	Relation	Tail
student	AtLocation	school	talk	HasProperty	cheap
curator	AtLocation	museum	future	HasProperty	uncertain
leader	AtLocation	group	be sure	HasSubevent	ask
glue	CapableOf	dry	be hungry	HasSubevent	eat
anyone	CapableOf	think	intrude into	HasFirstSubevent	shoot
door	CapableOf	open	go at	HasFirstSubevent	work
love	Causes	be friendly	closer	HasLastSubevent	go
attract	Causes	be vulgar	world	MadeOf	country
want	Causes	be closer	whole	MadeOf	part
music	CausesDesire	listen	run	MotivatedByGoal	afraid
friend	CausesDesire	talk	eat	MotivatedByGoal	hungry
choice	CausesDesire	entitle	sleep	MotivatedByGoal	tired
art	CreatedBy	artist	wall	PartOf	house
playoff	CreatedBy	team	child	PartOf	family
money	CreatedBy	bank	bone	PartOf	fish
earth	DefinedAs	world	crime	ReceivesAction	commit
god	DefinedAs	truth	game	ReceivesAction	play
door	DefinedAs	entrance	video	ReceivesAction	watch
idea	Desires	come	table	UsedFor	sit at
word	HasA	meaning	radio	UsedFor	listen to
house	HasA	wall	pool	UsedFor	swim in
bathroom	HasA	sink	nose	LocatedNear	eye
save	HasPrerequisite	do part	heat	LocatedNear	fire
enter	HasPrerequisite	ask i	beaver	LocatedNear	dam

Table 21: TransOMCS generated ConceptNet-like commonsense knowledge tuples.

		Mapping rules
	Head	Replace <i>PersonX</i> and <i>PersonY</i> with concrete singular personal pronouns, i.e., I/he/she/man/women/person
Tail	xWant/oWant/ xIntent/xNeed	Add a personal pronoun in front of the tail and remove the initial “to”
	xEffect/oEffect	Add a personal pronoun in front of the tail
	xReact/oReact	Add a personal pronoun and “be” in front of the tail
	xAttr	Add a personal pronoun and “be” in front of the tail

Table 22: Mapping rules from ATOMIC to ASER.

shortest path length starts from 1, where the shortest path length between two directly connected nodes is 1. We can conclude that, within a few hops of reasoning in ASER, a decent percentage of ATOMIC relations can be inferred. Some examples are presented in Table 24. For instance, the knowledge that if *PersonX bites PersonX’s tongue* then the person would want to *cry*, can be entailed from the Precedence discourse relation in ASER.

9.2.2. Mining ATOMIC-like Knowledge from ASER

As the heads and tails in ATOMIC are all arbitrary sentences, the aforementioned pattern mining approach used in TransOMCS is no longer suitable. To effectively convert ASER knowledge into the ATOMIC format, we propose to use a neural network based classifier instead of hard patterns. After we match ATOMIC and ASER, we will use the matched eventualities and associated sub-graph as the positive training examples. For each matched eventuality, we consider its one-hop or two-hop neighbors in ASER to be the candidate eventualities for populating commonsense knowledge of the corresponding ATOMIC relation, whose examples are shown in Table 26. With the help of a graph-based knowledge graph population model and the random negative example sampling, we successfully acquire large-scale commonsense knowledge in the format of ATOMIC. As demonstrated in Figure 15 and Table 27, both the original extracted eventualities and edges and those after the conceptualization can help us find rich commonsense about daily events. For example, before the conceptualization, we can find some knowledge like ‹‹She takes antibi-

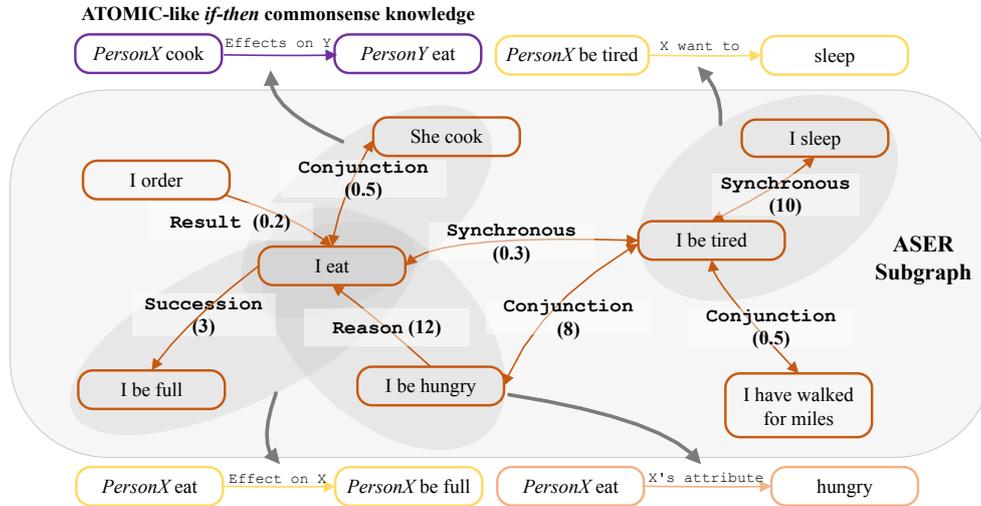


Figure 15: An illustration about exploring novel inferential commonsense knowledge about events. The center of the figure is a real subgraph of ASER. The grey ovals across ASER nodes are the relations that can be transferred to plausible *if-then* commonsense relations. For example, the <I eat, Succession, I eat> tuple in ASER can be intuitively written as the ATOMIC format, <PersonX eats, Effects on X, be full>).

Relation	Nodes	Edges	Avg. Shortest Path Length
oEffect	31.1%	25.36%	2.41
oReact	87.3%	51.53%	2.22
oWant	61.6%	36.95%	2.47
xAttr	95.8%	53.67%	2.38
xEffect	33.1%	21.81%	2.51
xIntent	33.8%	21.06%	2.56
xNeed	52.9%	24.91%	2.67
xReact	88.7%	52.66%	2.25
xWant	58.8%	30.60%	2.59
Average	62.9%	35.91%	2.44

Table 23: Mapping statistics of ATOMIC nodes and edges in ASER. The *Nodes* and *Edges* columns denote the percentage of ATOMIC nodes or edges that can be found in ASER. The *Avg. Shortest Path Length* column presents the average shortest path length of the matched ATOMIC edges in ASER.

otic,” Result, “She gets better”), which is rather specific. After the conceptualization, we can get a more abstract level commonsense that <“PersonX takes medicine,” Result, “PersonX gets better”>. Further experiments in [65] also show that compared with a pure supervised model, the knowledge populated with our approach is much more novel and diverse with the comparable high quality.

Similar but a bit different with that in TransOMCS, we evaluate the acquired commonsense knowledge by DISCOS using accuracy, novelty, and diversity. For accuracy, we ask annotators from Amazon Mechanical Turk to determine whether the commonsense tails generated by either COMET or DISCOS are plausible or not. We randomly sampled 50 heads for each relations, and acquire the top 10 results provided by the two models. For COMET, the top 10 results are acquired by Top 10 results using beam search with beam size 10. For DISCOS, the results are acquired by selecting the top 10 neighbors from ASER that are ranked the highest by BERTSAGE, a graph-aware model for populating commonsense knowledge on ASER. For novelty, we report the proportion of generated tails that are novel ($Novelty_{tail}$), and the proportion of novel tails in the set of all the unique generated tails ($Novelty_{unique}$). The second novelty metric is also expected to be high to avoid the situation when $Novelty_{tail}$ is high while the *novel* tails are all the same. We also check the diversity among the ten generated tails for each head-relation pair. The proportion of distinct unigrams (Dist-1) and bigrams (Dist-2) among the total number of generated unigrams and bigrams are used

Head	Tail	ATOMIC-Rel	ASER-Rel
<i>PersonX</i> bites <i>PersonX</i> 's tongue	<i>PersonX</i> cries	xWant	Precedence
<i>PersonX</i> feels hungry	<i>PersonX</i> eats	xWant	Conjunction
<i>PersonX</i> opens the envelope	<i>PersonX</i> read the letter	xWant	Co_Occurance
<i>PersonX</i> pays <i>PersonX</i> 's bill	<i>PersonX</i> leaves the restaurant	xWant	Co_Occurance
<i>PersonX</i> bleeds profusely	<i>PersonX</i> passes out	xEffect	Co_Occurance
<i>PersonX</i> goes to party	<i>PersonX</i> gets drunk	xEffect	Conjunction
<i>PersonX</i> plays well	<i>PersonX</i> wins	xEffect	Co_Occurance
<i>PersonX</i> wins the lottery	<i>PersonX</i> becomes rich	xEffect	Co_Occurance
<i>PersonX</i> would better go	<i>PersonX</i> is busy	xAttr	Condition
<i>PersonX</i> bites <i>PersonX</i> 's nail	<i>PersonX</i> is nervous	xAttr	Synchronous
<i>PersonX</i> eats <i>PersonX</i> 's breakfast	<i>PersonX</i> is hungry	xAttr	Condition
<i>PersonX</i> holds <i>PersonX</i> 's tongue	<i>PersonX</i> is quiet	xAttr	Co_Occurance
<i>PersonX</i> can not sleep	<i>PersonX</i> is stressed	xReact	Reason, Condition
<i>PersonX</i> is away from home	<i>PersonX</i> is lonely	xReact	Conjunction
<i>PersonX</i> is looking forward to it	<i>PersonX</i> is excited	xReact	Conjunction
<i>PersonX</i> tells <i>PersonY</i> everything	<i>PersonX</i> is trusted	xReact	Conjunction
<i>PersonX</i> accepts the challenge	<i>PersonX</i> wins	xIntent	Co_Occurance
<i>PersonX</i> bows <i>PersonX</i> 's head	<i>PersonX</i> prays	xIntent	Co_Occurance
<i>PersonX</i> removes <i>PersonX</i> 's hat	<i>PersonX</i> shows respect	xIntent	Co_Occurance
<i>PersonX</i> sits in car	<i>PersonX</i> waits for <i>PersonY</i>	xIntent	Co_Occurance
<i>PersonX</i> begins <i>PersonX</i> 's work	<i>PersonX</i> gets up	xNeed	Contrast, Conjunction
<i>PersonX</i> closes the door	<i>PersonX</i> has opened it	xNeed	Synchronous
<i>PersonX</i> gets a divorce	<i>PersonX</i> gets married	xNeed	Reason
<i>PersonX</i> makes amends	<i>PersonX</i> apologizes	xNeed	Conjunction
<i>PersonX</i> calls <i>PersonY</i> 's name	<i>PersonY</i> turns around	oEffect	Co_Occurance
<i>PersonX</i> receives a text	<i>PersonY</i> waits	oEffect	Synchronous
<i>PersonX</i> takes <i>PersonY</i> a picture	<i>PersonY</i> smiles	oEffect	Co_Occurance
<i>PersonX</i> tries to tell <i>PersonY</i>	<i>PersonY</i> refuses to listen	oEffect	Contrast
<i>PersonX</i> gets pregnant	<i>PersonY</i> wants to marry <i>PersonX</i>	oWant	Precedence
<i>PersonX</i> has not seen <i>PersonY</i> in years	<i>PersonY</i> wants to see <i>PersonX</i>	oWant	Co_Occurance
<i>PersonX</i> puts <i>PersonX</i> 's arm around <i>PersonY</i>	<i>PersonY</i> pushes <i>PersonX</i> away	oWant	Conjunction
<i>PersonX</i> steals <i>PersonY</i> 's wallet	<i>PersonY</i> calls the police	oWant	Conjunction
<i>PersonX</i> complains to the manager	<i>PersonY</i> is sorry	oReact	Co_Occurance
<i>PersonX</i> did an excellent job	<i>PersonY</i> is happy	oReact	Conjunction
<i>PersonX</i> gives <i>PersonY</i> money	<i>PersonY</i> is grateful	oReact	Conjunction

Table 24: Overlaps of ASER and ATOMIC.

here.

Table 25 shows the performance of DISCOS compared with COMET. DISCOS can achieve comparable or even better human annotated accuracy on oEffect, oReact, oWant, xIntent, and xNeed) among the nine relations. These relations are either of a smaller amount of annotations in ATOMIC or require more temporal knowledge (i.e., xIntent, and xNeed are the causes of the head event instead of effects). In addition, DISCOS can significantly outperform COMET in terms of novelty. The reason behind this is similar to that in TransOMCS, which is that COMET fits the training data too well and can suffer from selection bias [67]. Due to the limitation of using beam search to generate multiple tails given a head-relation pair, COMET also performs poorly on both diversity metrics than COMET. As DISCOS adopts an information extraction plus classification framework instead of a text generation framework, it does not suffer from that problem.

Relation	Model	Novelty _{tail}	Novelty _{unique}	Dist-1	Dist-2	Accuracy
oEffect	COMET@10	16.8	40.4	60.3	76.3	59.8
	DISCOS@10	62.9	76.2	66.7	89.3	68.3
oReact	COMET@10	0.4	4.9	35.5	13.5	69.6
	DISCOS@10	22.5	50.4	33.5	35.9	67.1
oWant	COMET@10	9.8	32.4	46.6	84.1	69.0
	DISCOS@10	55.8	75.4	69.0	93.8	69.9
xAttr	COMET@10	0.1	0.8	8.3	4.2	77.7
	DISCOS@10	12.0	30.4	26.0	27.4	66.7
xEffect	COMET@10	8.0	24.1	58.4	81.8	75.4
	DISCOS@10	54.5	71.1	67.2	90.4	60.9
xIntent	COMET@10	12.7	31.2	42.9	75.7	86.2
	DISCOS@10	51.7	74.1	61.5	87.3	87.8
xNeed	COMET@10	18.6	41.0	41.4	75.7	80.7
	DISCOS@10	44.2	66.2	63.6	88.4	84.9
xReact	COMET@10	0.4	4.7	27.1	12.1	75.6
	DISCOS@10	9.1	42.8	29.3	32.9	68.4
xWant	COMET@10	12.3	30.5	42.2	78.7	78.9
	DISCOS@10	38.1	62.0	65.3	91.5	73.4

Table 25: Main evaluation results of DISCOS and COMET.

ATOMIC Head	ATOMIC Tail	ATOMIC-Rel	Add. Neigh. by ASER
<i>PersonX</i> bites <i>PersonX</i> 's tongue	<i>PersonX</i> cries	xWant	<i>PersonY</i> strikes <i>PersonX</i> carefully with back
<i>PersonX</i> bows <i>PersonX</i> 's head	<i>PersonX</i> prays	xWant	<i>PersonX</i> cover <i>PersonX</i> 's face with hands
<i>PersonX</i> catch <i>PersonY</i> 's eye	<i>PersonX</i> makes an impression	xWant	<i>PersonY</i> is interested
<i>PersonX</i> becomes angry	<i>PersonX</i> yells	xEffect	<i>PersonX</i> asks for explanation
<i>PersonX</i> goes to the party	<i>PersonX</i> gets drunk	xEffect	<i>PersonX</i> 's stomach hurts
<i>PersonX</i> wins the lottery	<i>PersonX</i> becomes rich	xEffect	<i>PersonX</i> would quit <i>PersonX</i> 's job
<i>PersonX</i> can not sleep	<i>PersonX</i> is stressed	xReact	<i>PersonX</i> had a bad day at work
<i>PersonX</i> is away from home	<i>PersonX</i> is lonely	xReact	<i>PersonX</i> tries to talk to people
<i>PersonX</i> is looking forward to it	<i>PersonX</i> is excite	xReact	<i>PersonX</i> is working hard to get there
<i>PersonX</i> accepts the challenge	<i>PersonX</i> win	xIntent	<i>PersonY</i> plays
<i>PersonX</i> bows <i>PersonX</i> 's head	<i>PersonX</i> prays	xIntent	<i>PersonX</i> is silent
<i>PersonX</i> sits in car	<i>PersonX</i> waits for <i>PersonY</i>	xIntent	the police gets <i>PersonX</i> out
<i>PersonX</i> gets a divorce	<i>PersonX</i> gets married	xNeed	<i>PersonX</i> 's spouse cheats on <i>PersonX</i>
<i>PersonX</i> makes amends	<i>PersonX</i> apologizes	xNeed	<i>PersonX</i> did wrong

Table 26: Additional commonsense neighbors that ASER can provide, which can be learned by a knowledge graph population model.

10. Applications on Downstream Tasks

After the release of the ASER database²¹, many efforts have been devoted to applying the ASER knowledge for downstream tasks. In this section, we briefly introduce representative works of applying the ASER knowledge for downstream tasks and their key observations. More technical details can be found in the original papers.

10.1. Converting ASER into the Format of Human-crafted Commonsense Knowledge Graph

As discussed in Section 9.1, we explored how to convert ASER knowledge into the format of ConceptNet [1]. To test whether the converted ASER knowledge can help downstream tasks, we conduct experiments on two downstream tasks: commonsense reading comprehension [68] and dialogue generation [69]. Besides the original ConceptNet knowledge base, we also compare with other commonsense knowledge retrieval methods (i.e., COMET [64] and LAMA [10]).

The experimental results are shown in Table 28. For the reading comprehension task, adding the ASER knowledge contributes 0.37 overall accuracy, compared to 0.21 contribution of OMCS. Meanwhile, the contributions of COMET and LAMA are minor for this task. For the dialogue generation task, ASER knowledge also shows remarkable improvement in the quality of generated responses. At the same time, adding other knowledge resources to OMCS does not provide any meaningful improvements to the performance. The reason behind this could be that

²¹<https://github.com/HKUST-KnowComp/ASER>

	Head	Tail
Extracted	she take antibiotic	she get better
Conceptualized	<i>PersonX take Medicine</i>	<i>PersonX get better</i>
Extracted	he pay he bill	money be not plentiful with he
Conceptualized	<i>PersonX pay Short-Dated-Asset</i>	money be not plentiful with <i>PersonX</i>
Extracted	i win the lottery	i become rich
Conceptualized	<i>PersonX win Form-of-Gambling</i>	<i>PersonX become rich</i>
Extracted	he spill coffe	i ask for refill
Conceptualized	<i>PersonX spill Beverage</i>	<i>PersonY ask for refill</i>

Table 27: Examples of *if-then* commonsense knowledge in ASER. The knowledge before and after with Conceptualization are indicated with “Extracted” and “Conceptualized.”

Commonsense Knowledge Resource	Reading Comprehension		Dialog Generation	
	Accuracy (%)	Δ (%)	BLEU	Δ
Base model (no external knowledge resource)	82.90	-	0.54	-
+OMCS	83.11	+0.21	0.72	+0.18
+COMET	83.12	+0.22	0.61	+0.07
+LAMA	83.13	+0.23	0.56	+0.02
+ASER knowledge	83.27	+0.37	1.85	+1.31

Table 28: Effect of different knowledge resources on commonsense reading comprehension [68] and dialogue generation [69] tasks.

COMET and LAMA provide limited high quality novel commonsense knowledge. For example, the original OMCS on average contributes 1.46 supporting tuples²² and ASER knowledge contributes another 3.36 supporting tuples. As a comparison, COMET and LAMA only provide 0.01, 0.49 additional tuples respectively.

10.2. Combining ASER Knowledge with Language Models

Besides converting ASER into commonsense triplets, another work [70] tries to combine the structured knowledge and pre-trained language models. Motivated by the observation that while language models have already captured rich knowledge, they often only perform well when the semantic unit is a single token while poorly when the semantic unit is more complex (e.g., a multi-token named entity or an eventuality [71]). For example, if we follow LAMA [10] to analyze the knowledge contained in BERT-large [8] with a token prediction task, we can find out that BERT can understand that birds can fly, and a car is used for transportation, but it fails to understand the relation between “Jim yells at Bob” and relevant eventualities. An important reason behind this is that current language models heavily rely on token-level masked language models (MLMs) as the loss function, which can effectively represent and memorize token co-occurrence statistics²³ but struggle at perceiving multi-token concepts. To address this issue, [70] proposed to first verbalize the sub-graphs in ASER into sentences and then further fine-tune the pre-trained language models. A specific loss is added during the training phase to help the models to learn the complex eventuality knowledge in ASER.

To test whether the knowledge in ASER can help improve language models’ commonsense reasoning ability, [70] conducted experiments on three popular commonsense reasoning tasks: (1) ROCStories [73], which is widely used for story comprehension tasks such as Story Cloze Test; (2) MATRES [74], that focuses on the temporal commonsense between events; (3) COPA [75] that works on the causal commonsense. Experimental results show that the ASER knowledge can significantly improve the performance of pre-trained language models on these downstream tasks. It also supports our assumption that due to the limitation of the training loss, language models still need the support of structured knowledge to understand those complex commonsense knowledge.

²²Here by supporting tuple, we mean that the head and tail concept appear in the post and response respectively.

²³Sinha et al., [72] also explains the success of LMs due to distributional information. These models pre-trained over sentences with shuffled word order still achieve high accuracy.

10.3. Leveraging the Knowledge in ASER for Script Learning

ASER has been found useful for the task of Script Learning [76, 77]. The task of Script Learning aims to predict plausible subsequent events given an event chain describing previous states [78]. For example, a script depicting someone going to the restaurant may contain “*PersonX* goes to the restaurant,” “*PersonX* reads the menu,” and “*PersonX* orders food.” Script learning aims to predict the following events given the known event chain, for example in the previous case the next step can be “*PersonX* eats food.” Understanding scripts can be of vital importance on tasks such as storytelling, dialogue generation, and event understanding.

Lv et al. [76] use Elastic Search to match the events from event chains to ASER nodes, and select relevant supporting knowledge from their neighbors in ASER. The retrieved knowledge from ASER is then encoded with RoBERTa [9] and aggregated using an attention mechanism. The knowledge representation is then concatenated with the representation of the event chain as the final representation. Such a knowledge-aware model can boost the performance of RoBERTa-Large by over 2 points in terms of accuracy on the Multi-Choice Narrative Cloze (MCNC) dataset [79]. Furthermore, instead of only focusing on related subgraphs from ASER of a certain event chain, which may not be enough to equip the model with general script reasoning ability, Zhou et al. [77] proposed to pre-train a discriminative knowledge model on ASER, where the task is to classify the relationship given head and tail in a (h, r, t) triple. The head and tail are encoded separately with pre-trained language models and an interactive concatenation is applied to model their inner relationship. The finetuned encoder is then used as the encoder for events in the event chains. A chain-contextualized Bi-LSTM is then applied to deal with event chains. This model can learn rich relational patterns in the ASER graph for a script in a more supportive way than including local sub-structures only. Experimental results show that it can further boost the performance of Lv et al. [77] by 5 points.

11. Related Works

In this section, we introduce related works about commonsense knowledge acquisition, linguistic relation based information extraction systems, and conceptualization.

11.1. Commonsense Knowledge Acquisition

The acquisition of commonsense knowledge can be categorized into three main categories, crowdsourcing [7, 80, 1, 26, 27, 28], automatic construction from large-scale corpora [81, 82, 83, 84], and more recently, mining from pre-trained language models [10, 85, 86, 87]. Details are as follows.

Crowdsourcing Commonsense Knowledge Bases: Commonsense knowledge, primarily possessed by ordinary people, was first formalized and collected from human beings ourselves [80] with specific guidance towards specific domains. The CYC project asked knowledge engineers to write assertions and formalize the text to logical formats to support logical reasoning. ConceptNet [1] is originated from the Open-Mind CommonSense (OMCS) [88] project, human annotations are applied to acquire over 400K commonsense assertions among world entities. The latest version of ConceptNet 5 [61] now involves the English version of previous ConceptNets, as well as millions of facts from other taxonomy like WordNet and DBPedia. For each entity in ConceptNet, it can be linked to WordNet, Wiktionary, OpenCyc, and DBPedia. Moreover, ConceptNet is now a multi-lingual knowledge base that can also build connections between 83 languages. While ConceptNet focuses on commonsense relations among entities or noun phrases, ATOMIC [26] is proposed to investigate rich *if-then* relationships among daily social events. Nine social interaction related relations are developed and human annotators are asked to write the corresponding causes or effects of a certain base event. ATOMIC₂₀ [27] is further proposed to unify the triples from ConceptNet and ATOMIC, together with some newly developed relations. GLUCOSE [28] is a commonsense knowledge base constructed based on ROC Story [73]. The commonsense causal relations in GLUCOSE are based on cognitive psychology theories that humans primarily focus on events, their timeline, locations of entities, causes and motivations of the event, and emotional trajectory of the character, when focusing on a piece of narrative.

Commonsense Knowledge by Information Extraction: Though in general, commonsense knowledge is not explicitly expressed, there is still a non-negligible amount of commonsense knowledge of certain types that can be mined using information extraction tools, such as salient properties of objects [81, 83], verb-oriented selectional preference commonsense [84], and general statements [89, 90]. WebChild [81] uses semi-supervised label propagation over

constructed graphs from web contents, where the seed commonsense knowledge is derived from WordNet. Quasimodo [83] derives commonsense knowledge from search-engine query logs and QA forums. Syntactical patterns are designed to capture salient properties of objects, for example, detecting questions starting with *Why* and some specific auxiliary verbs of a certain object. Verb-Oriented Commonsense Knowledge [84] explores plausible subjects and objects of certain verbs. A large-scale probabilistic taxonomy, Probbase [43], is used to conceptualize subject and object in a verb phrase to get a general s-v-o phrase. An entropy-based filter is applied to determine the appropriate level of conceptualization and a language model is used to score the quality of the provided s-v-o triples. To capture knowledge that goes beyond (h, r, t) triples, some knowledge bases storing general statements are proposed to be more flexible in representing commonsense knowledge. GenericsKB [89] is constructed from large corpora using BERT-based scoring as a filter and including contextual metadata as supporting information. Such kind of knowledge is more flexible and can help some downstream tasks such as question-answering.

Commonsense Knowledge in Pre-trained Language Models: With the number of parameters in pre-trained language models [8, 9, 91, 92, 93, 94, 86] increasing exponentially, researchers are exploring ways to mine commonsense knowledge directly from pre-trained language models, in view of their strong representation ability on large-scale corpora and compositional generalization ability. Such exploration includes both supervised approaches [64, 27] and unsupervised approaches [10, 85, 95, 86]. For supervised learning based approaches, pre-trained language models such as BART [93] and GPT-2 [92] are finetuned on large-scale commonsense knowledge bases on a conditional generation task, where the head and relation in the commonsense triple are given as input and the tail serves as the expected output. Those models finetuned on ConceptNet, ATOMIC, and ATOMIC₂₀²⁰ can generate commonsense tails with high precision, though may not be generalized enough to generate novel knowledge that is required for commonsense knowledge acquisition. For unsupervised approaches, prompts are designed to probe commonsense knowledge directly from large pre-trained models. LAMA [10] and Davison et al. [85] designed simple hand-written prompts to conduct factual probing in ConceptNet from BERT [8]. Automatically generated prompts such as best paraphrase-based prompts [95], the best sequence of tokens maximizing the gold label likelihood [96], directly optimized embeddings instead of prompts in the form of text [97] are used to feed into pre-trained language models to generate outputs, whereas the language model remains untrained. ATOMIC^{10x} [87] leverages GPT-3, with 100x larger the scale than models such as GPT-2-XL, where some seeds from ATOMIC are used as prompts to acquire commonsense knowledge directly from GPT-3. Human evaluations demonstrate that such an automatically constructed commonsense knowledge base can outperform human annotation in terms of correctness and diversity.

11.2. Conceptualization

Conceptualization in Cognitive Science: People posit the importance of a specific element of human commonsense, conceptualization. As observed by psychologists, “concepts are the glue that holds our mental world together” [98]. Human beings are able to make reasonable inferences by utilizing the IsA relationship between real-world concepts and instances. For example, without knowing what a “floppy disk” is, given that it is a “memory device,” people may infer that it may store data and be readable by a computer. In K-lines theory [99], people conceptualize the world as a pyramid, and map a K-node (a mental state) to this pyramid, which has a lower-band limit and a higher band limit to ensure right common and non-conflicting properties. When we want to remember something, we create a K-line for it; when later it is activated, the K-line induces a subset of those mental agencies resembling states that created the K-line. A lower K-line could affect the instantiation of a more abstract higher level K-line so that K-nodes help us to make abstraction, logical, and procedural reasoning. For example, we could create a K-line for Tesla by mapping and connecting “company,” “big company,” “IT company,” “AI company,” “high-tech company,” and “automobile company.” As properties are usual non-conflicting, combining the concrete accumulation of particular instances with the rejection of strongly dissonant properties automatically leads to a rather abstract unification.

Conceptualization in Computer Science: In the computer science community, researchers also explored how to leverage the conceptualization to help machines understand the world. Probbase [43] is a large-scale probabilistic taxonomy to store such “IsA” relations between instances and concepts, where 2.7 million concepts are automatically harnessed from 1.68 billion documents. It has been found useful for several natural language understanding tasks [16, 100]. Besides, pattern-based word co-occurrence statistics [101, 102] and distributed embedding models [103, 104] can help detect the hypernymy relation to enrich the conceptualization knowledge base. However, conceptualization needs to address the typicality and ambiguity. Various computational approaches have been analyzed for deriving basic-level categorization as a trade-off [105]. To address this issue, contextualized conceptualization

was proposed. Previous works have explored how to leverage topic modeling [106] and external knowledge [100] to better conceptualize the concepts based on the local context. A recent work also explored how to capture the connection between nouns and associated verbs for the better conceptualization [84]. Last but not least, some attempts on pre-trained models for context-dependent conceptualization also indicated the counter intuitiveness and conceptual inconsistency [107].

12. Conclusions and Future Works

In this paper, we focus on the commonsense knowledge acquisition problem. Throughout the years, the community has devoted enormous efforts to acquiring commonsense knowledge with either human annotation or information extraction techniques. However, these works are either not scalable or can only handle a specific kind of pre-defined commonsense knowledge. To explore a more fundamental understanding of the commonsense knowledge about daily events and states, we follow previous research on the lower bound of semantic theory [53], partial information [5], and K-lines theory [99], and propose to represent commonsense knowledge with higher-order selectional preference over eventualities. Specifically, we first leverage the distribution of daily eventualities and their relations in raw corpus to simulate the plausibility of different semantic combinations, and then leverage the conceptualization module to conceptualize the observed knowledge into an abstract level. Following this methodology, we develop a large-scale eventuality-centric commonsense knowledge graph ASER, which is a large-scale eventuality knowledge graph that contains 438 million eventualities and 648 million edges. Considering the large scale of commonsense, we propose an unsupervised pipeline to extract rich commonsense knowledge about events from the raw corpus instead of human annotation. To effectively represent humans' preference about daily events, we design ASER to be weighed, and larger weight indicates that the eventuality or edge is more likely to happen. We conduct human evaluations, case studies, and extrinsic evaluations to evaluate the quality of ASER. As one of the main extraction methodologies of our approach is that we prefer accuracy over recall because we can easily scan more data, even though our current extraction pipeline may sacrifice the recall, it guarantees the high quality of the extracted knowledge. Further experiments also demonstrate that the knowledge in ASER can be effectively converted into human-crafted commonsense knowledge in other commonsense knowledge bases such as ConceptNet [1] and ATOMIC [26] and then help downstream tasks such as reading comprehension [68], dialogue generation [69], story completion [73], temporal relation prediction [74], and causal relation prediction [75].

As a long-standing artificial intelligence problem, commonsense reasoning is still challenging for current natural language understanding models. In this work, even though we shed some light on how to represent the commonsense knowledge from the angle of partial information, there is still a long way to go to fully solve the commonsense reasoning problem. Specifically, our current research has the following limitations that need to be addressed in the future:

1. **Evaluation:** The first challenge we are still facing is the lack of a good evaluation system. Unlike other tasks, most current commonsense reasoning tasks (e.g., Winograd Schema Challenge [6]) are not directly evaluating models' commonsense reasoning abilities. Instead, they are a kind of approximation. Take WSC as an example, many research has discovered that current models can bypass the essential commonsense reasoning and solve the questions with other information [108]. It is quite often that we are just solving a "dataset" without solving the underlining "task" we truly want to solve.
2. **Storage and Computation Efficiency:** As aforementioned, ASER has 438 million eventualities and 648 million edges. Such a large scale guarantees the coverage of ASER, but it also brings a huge burden for storage and computation. Our current hardware architecture and inference algorithms still cannot support fast inference and response. We can try to address this issue from two angles: (1) Better hardware architecture; (2) Better knowledge graph organization.
3. **Contextualized Conceptualization:** Another critical challenge we are facing is how to correctly contextualize the observed concepts. As discussed by the K-lines theory [99] and recent research on conceptualization [84], it is important to conceptualize the observed objects into the correct concept level based on the local context. However, to the best of our knowledge, there is still no reliable contextualized conceptualization model that can handle the open-world scenario. In this work, we use the distribution of concepts over a big corpus instead of the local context to partially remedy this issue. For example, after observing "dogs can bark" and the probability of

an animal being a dog is 0.08, we will conclude that the plausibility of eventuality “animal bark” is 0.08, which indicates that an animal may not always be able to bark, but compared with other entities such as “house,” an animal is more likely to bark. A potential limitation of this method is the reporting bias issue, as studied in [13], the correlation between the natural distribution and human’s commonsense knowledge is slightly less than 0.8. How to handle the reporting bias issue and effectively conceptualize the observed entities based on the local context is a problem worth exploring in the future.

All codes, data, and APIs are published at the project page²⁴ to encourage further research on commonsense and event understanding.

Acknowledgements

This paper was supported by the GRF (16211520) and the RIF (R6020-19 and R6021-20) from RGC of Hong Kong, the NSFC Fund (U20B2053) from the NSFC of China, the MHKJFS (MHP/001/19) from ITC of Hong Kong with special thanks to HKMAAC and CUSBLT, and the Jiangsu Province Science and Technology Collaboration Fund (BZ2021065).

Contributions

The contributions of all authors are as follows.

- **Hongming Zhang:** Proposing the idea of using higher-order selectional preference over eventualities to represent commonsense knowledge, designing the ASER structure, designing the eventuality and edge patterns, designing the eventuality extraction algorithm, selecting data, conducting intrinsic evaluation, conducting extrinsic evaluations (except dialogue system), and writing the paper.
- **Xin Liu:** Designing and implementing data pre-processing, constituency parsing, clause analyzing, relation extracting with discourse parsing systems, the ASER database schema, and the construction pipeline, providing scripts for extraction and conceptualization and APIs for knowledge databases, managing data and code, and drafting the major of Section 5 and 6.
- **Haojie Pan:** Exploring and implementing the conceptualization with Probase, conducting the extrinsic evaluation on the dialogue system, designing the client-server model for the distributed ASER system, and preparing the online demo, and drafting sections relevant to conceptualization.
- **Haowen Ke:** Pre-processing raw data with CoreNLP to acquire lemmatized tokens, pos-tags, name entities, dependency tree, and constituency tree, analyzing the inference results in ASER, and drafting Section 8.
- **Jiefu Ou:** Implementing the rule-based inference over ASER with the AMIE+ system. Assisting Xin Liu for discourse relation extraction and assisting Haowen Ke for analyzing the inference results in ASER.
- **Tianqing Fang:** Analyzing the relation between ASER and other commonsense knowledge bases, and drafting Section 9.
- **Yangqiu Song:** Proposing the ideas of building an eventuality centric knowledge graph, using conceptualization for abstraction and instantiation, managing the ASER project, and revising the paper.

²⁴<https://github.com/HKUST-KnowComp/ASER>

References

- [1] H. Liu, P. Singh, ConceptNet—a practical commonsense reasoning tool-kit, *BT technology journal* 22 (4) (2004) 211–226.
- [2] J. Gordon, B. V. Durme, L. K. Schubert, Learning from the web: Extracting general world knowledge from noisy text, in: *Proceedings of the Collaboratively-Built Knowledge Sources and Artificial Intelligence Workshop at the 24th AAAI Conference on Artificial Intelligence*, Atlanta, USA, 2010.
- [3] J. Gordon, L. K. Schubert, Quantificational sharpening of commonsense knowledge, in: *AAAI Fall Symposium on Commonsense Knowledge*, Arlington, USA, 2010.
- [4] P. Resnik, Selectional preference and sense disambiguation, in: *Tagging Text with Lexical Semantics: Why, What, and How?*, 1997.
- [5] Y. Wilks, An intelligent analyzer and understander of english, *Communications of the ACM* 18 (5) (1975) 264–274.
- [6] H. J. Levesque, E. Davis, L. Morgenstern, The winograd schema challenge, in: *AAAI Spring Symposium on Logical Formalizations of Commonsense Reasoning*, Stanford, USA, 2011, p. 47.
- [7] D. B. Lenat, R. V. Guha, *Building large knowledge-based systems; representation and inference in the Cyc project*, Addison-Wesley Longman Publishing Co., Inc., 1989.
- [8] J. Devlin, M. Chang, K. Lee, K. Toutanova, BERT: Pre-training of deep bidirectional transformers for language understanding, in: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Minneapolis, USA, 2019, pp. 4171–4186.
- [9] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, V. Stoyanov, RoBERTa: A robustly optimized BERT pretraining approach, *CoRR abs/1907.11692*.
- [10] F. Petroni, T. Rocktäschel, S. Riedel, P. S. H. Lewis, A. Bakhtin, Y. Wu, A. H. Miller, Language models as knowledge bases?, in: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*, Hong Kong, China, 2019, pp. 2463–2473.
- [11] H. Zhang, D. Khashabi, Y. Song, D. Roth, TransOMCS: From linguistic graphs to commonsense knowledge, in: *Proceedings of the 29th International Joint Conference on Artificial Intelligence*, 2020, pp. 4004–4010.
- [12] E. Bach, The algebra of events, *Linguistics and philosophy* 9 (1) (1986) 5–16.
- [13] H. Zhang, H. Ding, Y. Song, SP-10K: A large-scale evaluation set for selectional preference acquisition, in: *Proceedings of the 57th Conference of the Association for Computational Linguistics*, Florence, Italy, 2019, pp. 722–731.
- [14] J. Wang, M. Lan, A refined end-to-end discourse parser, in: *Proceedings of the 19th Conference on Computational Natural Language Learning: Shared Task*, Beijing, China, 2015, pp. 17–24.
- [15] J. M. Zacks, B. Tversky, Event structure in perception and conception, *Psychological Bulletin* 127 (1) (2001) 3.
- [16] Y. Song, H. Wang, Z. Wang, H. Li, W. Chen, Short text conceptualization using a probabilistic knowledgebase, in: *Proceedings of the 22nd International Joint Conference on Artificial Intelligence*, Barcelona, Spain, 2011, pp. 2330–2336.
- [17] Y. Song, S. Wang, H. Wang, Open domain short text conceptualization: A generative + descriptive modeling approach, in: *Proceedings of the 24th International Joint Conference on Artificial Intelligence*, Buenos Aires, Argentina, 2015, pp. 3820–3826.
- [18] C. F. Baker, C. J. Fillmore, J. B. Lowe, The berkeley FrameNet project, in: *Proceedings of 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics*, Montréal, Canada, 1998, pp. 86–90.
- [19] U. NIST, et al., The ace 2003 evaluation plan, *US National Institute for Standards and Technology (2003) 2003–08*.
- [20] J. Aguilar, C. Beller, P. McNamee, B. Van Durme, S. Strassel, Z. Song, J. Ellis, A comparison of the events and relations across ace, ere, tac-kbp, and FrameNet annotation standards, in: *Proceedings of the 2nd Workshop on EVENTS: Definition, Detection, Coreference, and Representation*, EVENTS@ACL, Baltimore, USA, 2014, pp. 45–53.
- [21] M. Palmer, D. Gildea, P. Kingsbury, The proposition bank: An annotated corpus of semantic roles, *Computational Linguistics* 31 (1) (2005) 71–106.
- [22] A. L. Meyers, R. Reeves, C. Macleod, R. Szekely, V. Zielinska, B. Young, R. Grishman, The NomBank project: An interim report, in: *Proceedings of the Workshop Frontiers in Corpus Annotation@HLT-NAACL*, Boston, USA, 2004.
- [23] J. Pustejovsky, P. Hanks, R. Sauri, A. See, R. Gaizauskas, A. Setzer, D. Radev, B. Sundheim, D. Day, L. Ferro, et al., The timebank corpus, in: *Corpus Linguistics*, Vol. 2003, 2003, p. 40.
- [24] N. A. Smith, Y. Choi, M. Sap, H. Rashkin, E. Allaway, Event2Mind: Commonsense inference on events, intents, and reactions, in: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, Melbourne, Australia, 2018, pp. 463–473.
- [25] B. Dalvi, L. Huang, N. Tandon, W. tau Yih, P. Clark, Tracking state changes in procedural text: A challenge dataset and models for process paragraph comprehension, *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, New Orleans, USA (2018) 1595–1604.
- [26] M. Sap, R. LeBras, E. Allaway, C. Bhagavatula, N. Lourie, H. Rashkin, B. Roof, N. A. Smith, Y. Choi, ATOMIC: An atlas of machine commonsense for if-then reasoning, in: *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*, Honolulu, USA, 2019, pp. 3027–3035.
- [27] J. D. Hwang, C. Bhagavatula, R. L. Bras, J. Da, K. Sakaguchi, A. Bosselut, Y. Choi, (Comet-) Atomic 2020: On symbolic and neural commonsense knowledge graphs, in: *Proceedings of the 25th AAAI Conference on Artificial Intelligence*, Virtual Event, 2021, pp. 6384–6392.
- [28] N. Mostafazadeh, A. Kalyanpur, L. Moon, D. W. Buchanan, L. Berkowitz, O. Biran, J. Chu-Carroll, GLUCOSE: Generalized and contextualized story explanations, in: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, Virtual Event, 2020, pp. 4569–4586.
- [29] N. Tandon, G. de Melo, A. De, G. Weikum, Knowlywood: Mining activity knowledge from hollywood narratives, in: *Proceedings of the 24th ACM International Conference on Information and Knowledge Management*, Melbourne, Australia, 2015, pp. 223–232.
- [30] R. Prasad, N. Dinesh, A. Lee, E. Miltsakaki, L. Robaldo, A. K. Joshi, B. L. Webber, The penn discourse treebank 2.0, in: *Proceedings of the International Conference on Language Resources and Evaluation*, Marrakech, Morocco, 2008.
- [31] R. Jackendoff, *Semantic Structures*, MIT Press, 1992.

- [32] A. P. Mourelatos, Events, processes, and states, *Linguistics and philosophy* 2 (3) (1978) 415–434.
- [33] L. Ehrlinger, W. Wöß, Towards a definition of knowledge graphs, in: *Joint Proceedings of the Posters and Demos Track of the 12th International Conference on Semantic Systems and the 1st International Workshop on Semantic Change & Evolving Semantics co-located with the 12th International Conference on Semantic Systems*, Leipzig, Germany, Vol. 1695, 2016.
- [34] G. A. Miller, *WordNet: an electronic lexical database*, MIT Press, 1998.
- [35] K. D. Bollacker, C. Evans, P. K. Paritosh, T. Sturge, J. Taylor, Freebase: a collaboratively created graph database for structuring human knowledge, in: *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, Vancouver, Canada, 2008, pp. 1247–1250.
- [36] O. Etzioni, M. Cafarella, D. Downey, Webscale information extraction in knowitall (preliminary results), in: *Proceedings of the 13th international conference on World Wide Web*, New York, USA, 2004, pp. 100–110.
- [37] M. Banko, M. J. Cafarella, S. Soderland, M. Broadhead, O. Etzioni, Open information extraction from the web, in: *Proceedings of the 20th International Joint Conference on Artificial Intelligence*, Hyderabad, India, 2007, pp. 2670–2676.
- [38] F. M. Suchanek, G. Kasneci, G. Weikum, YAGO: a core of semantic knowledge, in: *Proceedings of the 16th International Conference on World Wide Web*, Banff, Canada, 2007, pp. 697–706.
- [39] J. Hoffart, F. M. Suchanek, K. Berberich, G. Weikum, YAGO2: A spatially and temporally enhanced knowledge base from wikipedia, *Artificial Intelligence* 194 (2013) 28–61.
- [40] R. Navigli, S. P. Ponzetto, BabelNet: The automatic construction, evaluation and application of a wide-coverage multilingual semantic network, *Artificial Intelligence* 193 (2012) 217–250.
- [41] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, Z. G. Ives, DBpedia: A nucleus for a web of open data, in: *Proceedings of 6th International Semantic Web Conference*, Busan, Korea, Vol. 4825, 2007, pp. 722–735.
- [42] A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. R. H. Jr., T. M. Mitchell, Toward an architecture for never-ending language learning, in: *Proceedings of the 24th AAAI Conference on Artificial Intelligence*, Atlanta, Georgia, USA, 2010.
- [43] W. Wu, H. Li, H. Wang, K. Q. Zhu, Probase: A probabilistic taxonomy for text understanding, in: *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*, Scottsdale, USA, 2012, pp. 481–492.
- [44] X. Dong, E. Gabrilovich, G. Heitz, W. Horn, N. Lao, K. Murphy, T. Strohmman, S. Sun, W. Zhang, Knowledge vault: A web-scale approach to probabilistic knowledge fusion, in: *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, New York, USA, 2014, pp. 601–610.
- [45] J. Berant, A. Chou, R. Frostig, P. Liang, Semantic parsing on freebase from question-answer pairs, in: *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, Seattle, USA, 2013, pp. 1533–1544.
- [46] G. Glavas, J. Snajder, M. Moens, P. Kordjamshidi, HiEve: A corpus for extracting event hierarchies from news stories, in: *Proceedings of the 9th International Conference on Language Resources and Evaluation*, Reykjavik, Iceland, 2014, pp. 3678–3683.
- [47] H. Wang, M. Chen, H. Zhang, D. Roth, Joint constrained learning for event-event relation extraction, in: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, Virtual Event, 2020, pp. 696–706.
- [48] R. Han, Q. Ning, N. Peng, Joint event and temporal relation extraction with shared representations and structured prediction, in: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*, Hong Kong, China, 2019, pp. 434–444.
- [49] X. Liu, J. Ou, Y. Song, X. Jiang, On the importance of word and sentence representation learning in implicit discourse relation classification, in: *Proceedings of the 29th International Joint Conference on Artificial Intelligence*, 2020, pp. 3830–3836.
- [50] N. Xue, H. T. Ng, S. Pradhan, R. Prasad, C. Bryant, A. Rutherford, The conll-2015 shared task on shallow discourse parsing, in: *Proceedings of the 19th Conference on Computational Natural Language Learning: Shared Task*, Beijing, China, 2015, pp. 1–16.
- [51] E. V. Siegel, K. R. McKeown, Learning methods to combine linguistic indicators: Improving aspectual classification and revealing linguistic insights, *Computational Linguistics* 26 (4) (2000) 595–627.
- [52] H. Zhang, M. Chen, H. Wang, Y. Song, D. Roth, Analogous process structure induction for sub-event sequence prediction, in: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, Virtual Event, 2020, pp. 1541–1550.
- [53] J. J. Katz, J. A. Fodor, The structure of a semantic theory, *Language* 39 (2) (1963) 170–210.
- [54] P. S. Resnik, Selection and information: A class-based approach to lexical relationships, *IRCS Technical Reports Series* (1993) 200.
- [55] M. Steedman, J. Baldridge, Combinatory categorial grammar, *Non-Transformational Syntax: Formal and explicit models of grammar* (2011) 181–224.
- [56] P. Kingsbury, M. Palmer, From treebank to propbank, in: *Proceedings of the Third International Conference on Language Resources and Evaluation*, Las Palmas, Spain, 2002, pp. 1989–1993.
- [57] E. Sandhaus, The new york times annotated corpus, *Linguistic Data Consortium*, Philadelphia 6 (12) (2008) e26752.
- [58] P. Lison, J. Tiedemann, OpenSubtitles2016: Extracting large parallel corpora from movie and TV subtitles, in: *Proceedings of the 10th International Conference on Language Resources and Evaluation*, Portorož, Slovenia, 2016.
- [59] L. Galárraga, C. Teflioudi, K. Hose, F. M. Suchanek, Fast rule mining in ontological knowledge bases with AMIE+, *The VLDB Journal* (2015) 707–730.
- [60] Y. Sun, J. Han, *Mining Heterogeneous Information Networks: Principles and Methodologies*, Synthesis Lectures on Data Mining and Knowledge Discovery, Morgan & Claypool Publishers, 2012.
- [61] R. Speer, C. Havasi, ConceptNet 5: A large semantic network for relational knowledge, in: *The People’s Web Meets NLP, Collaboratively Constructed Language Resources*, 2013, pp. 161–176.
- [62] G. Angelì, M. J. J. Premkumar, C. D. Manning, Leveraging linguistic structure for open domain information extraction, in: *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing*, Beijing, China, 2015, pp. 344–354.
- [63] M. A. Hearst, Automatic acquisition of hyponyms from large text corpora, in: *The 14th International Conference on Computational Linguistics*, Nantes, France, 1992, pp. 539–545.
- [64] A. Bosselut, H. Rashkin, M. Sap, C. Malaviya, A. Celikyilmaz, Y. Choi, COMET: Commonsense transformers for automatic knowledge

- graph construction, in: *Proceedings of the 57th Conference of the Association for Computational Linguistics*, Florence, Italy, 2019, pp. 4762–4779.
- [65] T. Fang, H. Zhang, W. Wang, Y. Song, B. He, DISCOS: Bridging the gap between discourse knowledge and commonsense knowledge, in: *The 2021 Web Conference*, Virtual Event, 2021, pp. 2648–2659.
- [66] T. Fang, W. Wang, S. Choi, S. Hao, H. Zhang, Y. Song, B. He, Benchmarking commonsense knowledge base population with an effective evaluation dataset, in: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Virtual Event, 2021, pp. 8949–8964.
- [67] B. Zadrozny, Learning and evaluating classifiers under sample selection bias, in: *Proceedings of the 21st International Conference on Machine Learning*, Banff, Canada, 2004, p. 114.
- [68] S. Ostermann, M. Roth, A. Modi, S. Thater, M. Pinkal, SemEval-2018 Task 11: Machine comprehension using commonsense knowledge, in: *Proceedings of The 12th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2018*, New Orleans, USA, 2018, pp. 747–757.
- [69] Y. Li, H. Su, X. Shen, W. Li, Z. Cao, S. Niu, DailyDialog: A manually labelled multi-turn dialogue dataset, in: *Proceedings of the 8th International Joint Conference on Natural Language Processing*, Taipei, Taiwan, 2017, pp. 986–995.
- [70] C. Yu, H. Zhang, Y. Song, W. Ng, CoCoLM: Complex commonsense enhanced language model, *CoRR abs/2012.15643*.
- [71] P. Verga, H. Sun, L. B. Soares, W. W. Cohen, Facts as experts: Adaptable and interpretable neural memory over symbolic knowledge, *CoRR abs/2007.00849*.
- [72] K. Sinha, R. Jia, D. Hupkes, J. Pineau, A. Williams, D. Kiela, Masked language modeling and the distributional hypothesis: Order word matters pre-training for little, in: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Virtual Event, 2021, pp. 2888–2913.
- [73] N. Mostafazadeh, N. Chambers, X. He, D. Parikh, D. Batra, L. Vanderwende, P. Kohli, J. F. Allen, A corpus and cloze evaluation for deeper understanding of commonsense stories, in: *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, San Diego, USA, 2016, pp. 839–849.
- [74] Q. Ning, H. Wu, D. Roth, A multi-axis annotation scheme for event temporal relations, in: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, Melbourne, Australia, 2018, pp. 1318–1328.
- [75] A. S. Gordon, Z. Kozareva, M. Roemmele, SemEval-2012 Task 7: Choice of plausible alternatives: An evaluation of commonsense causal reasoning, in: *Proceedings of the 6th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT*, Montréal, Canada, 2012, 2012, pp. 394–398.
- [76] S. Lv, F. Zhu, S. Hu, Integrating external event knowledge for script learning, in: *Proceedings of the 28th International Conference on Computational Linguistics*, Virtual Event, 2020, pp. 306–315.
- [77] Y. Zhou, X. Geng, T. Shen, J. Pei, W. Zhang, D. Jiang, Modeling event-pair relations in external knowledge graphs for script reasoning, in: *Findings of the 59th Annual Meeting of the Association for Computational Linguistics*, Virtual Event, 2021, pp. 4586–4596.
- [78] N. Chambers, D. Jurafsky, Unsupervised learning of narrative event chains, in: *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics*, Columbus, USA, 2008, pp. 789–797.
- [79] Z. Li, X. Ding, T. Liu, Constructing narrative event evolutionary graph for script event prediction, in: *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, Stockholm, Sweden, 2018, pp. 4201–4207.
- [80] D. B. Lenat, CYC: A large-scale investment in knowledge infrastructure, *Communications of the ACM* 38 (11) (1995) 33–38.
- [81] N. Tandon, G. de Melo, F. M. Suchanek, G. Weikum, WebChild: harvesting and organizing commonsense knowledge from the web, in: *Seventh ACM International Conference on Web Search and Data Mining*, New York, USA, 2014, pp. 523–532.
- [82] N. Tandon, G. de Melo, G. Weikum, WebChild 2.0: Fine-grained commonsense knowledge distillation, in: *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (System Demonstrations)*, Vancouver, Canada, 2017, pp. 115–120.
- [83] J. Romero, S. Razniewski, K. Pal, J. Z. Pan, A. Sakhadeo, G. Weikum, Commonsense properties from query logs and question answering forums, in: *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019*, Beijing, China, November 3-7, 2019, 2019, pp. 1411–1420.
- [84] J. Liu, Y. Zhou, D. Wu, C. Wang, H. Jiang, S. Zhang, B. Xu, Y. Xiao, Mining verb-oriented commonsense knowledge, in: *Proceedings of the 36th IEEE International Conference on Data Engineering*, Dallas, USA, 2020, pp. 1830–1833.
- [85] J. Davison, J. Feldman, A. M. Rush, Commonsense knowledge mining from pretrained models, in: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*, Hong Kong, China, 2019, pp. 1173–1178.
- [86] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, D. Amodei, Language models are few-shot learners, in: *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020*, December 6-12, 2020, virtual, 2020.
- [87] P. West, C. Bhagavatula, J. Hessel, J. D. Hwang, L. Jiang, R. L. Bras, X. Lu, S. Welleck, Y. Choi, Symbolic knowledge distillation: from general language models to commonsense models, *CoRR abs/2110.07178*.
- [88] P. Singh, T. Lin, E. T. Mueller, G. Lim, T. Perkins, W. L. Zhu, Open Mind Common Sense: Knowledge acquisition from the general public, in: *On the Move to Meaningful Internet Systems, Heidelberg, Berlin*, 2002, pp. 1223–1237.
- [89] S. Bhakthavatsalam, C. Anastasiades, P. Clark, GenericsKB: A knowledge base of generic statements, *CoRR abs/2005.00660*.
- [90] T.-P. Nguyen, S. Razniewski, G. Weikum, Advanced semantics for commonsense knowledge extraction, in: *The 2021 Web Conference*, Virtual Event, 2021, pp. 2636–2647.
- [91] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, Improving language understanding by generative pre-training.
- [92] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al., Language models are unsupervised multitask learners.
- [93] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, L. Zettlemoyer, BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension, in: *Proceedings of the 58th Annual Meeting of the Association*

- for Computational Linguistics, Virtual Event, 2020, pp. 7871–7880.
- [94] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, P. J. Liu, Exploring the limits of transfer learning with a unified text-to-text transformer, *Journal of Machine Learning Research* 21 (2020) 1–67.
- [95] Z. Jiang, F. F. Xu, J. Araki, G. Neubig, How can we know what language models know, *Transactions of the Association for Computational Linguistics* 8 (2020) 423–438.
- [96] T. Shin, Y. Razeghi, R. L. L. IV, E. Wallace, S. Singh, AutoPrompt: Eliciting knowledge from language models with automatically generated prompts, in: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, Virtual Event, 2020, pp. 4222–4235.
- [97] Z. Zhong, D. Friedman, D. Chen, Factual probing is [MASK]: Learning vs. learning to recall, in: *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Virtual Event, 2021, pp. 5017–5033.
- [98] G. Murphy, *The big book of concepts*, MIT press, 2004.
- [99] M. Minsky, K-lines: A theory of memory, *Cognitive Science* 4 (2) (1980) 117–133.
- [100] Z. Wang, K. Zhao, H. Wang, X. Meng, J. Wen, Query understanding through knowledge-based conceptualization, in: *Proceedings of the 24th International Joint Conference on Artificial Intelligence*, Buenos Aires, Argentina, 2015, pp. 3264–3270.
- [101] S. Roller, D. Kiela, M. Nickel, Hearst patterns revisited: Automatic hypernym detection from large text corpora, in: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, Melbourne, Australia, 2018, pp. 358–363.
- [102] M. Le, S. Roller, L. Papaxanthos, D. Kiela, M. Nickel, Inferring concept hierarchies from text corpora via hyperbolic embeddings, in: *Proceedings of the 57th Conference of the Association for Computational Linguistics*, Florence, Italy, 2019, pp. 3231–3241.
- [103] K. A. Nguyen, M. Köper, S. S. im Walde, N. T. Vu, Hierarchical embeddings for hypernymy detection and directionality, in: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, Copenhagen, Denmark, 2017, pp. 233–243.
- [104] H. Chang, Z. Wang, L. Vilnis, A. McCallum, Distributional inclusion vector embedding for unsupervised hypernymy detection, in: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, New Orleans, USA, 2018, pp. 485–495.
- [105] Z. Wang, H. Wang, J. Wen, Y. Xiao, An inference approach to basic level of categorization, in: *Proceedings of the 24th ACM International Conference on Information and Knowledge Management*, Melbourne, Australia, 2015, pp. 653–662.
- [106] D. Kim, H. Wang, A. H. Oh, Context-dependent conceptualization, in: *Proceedings of the 23rd International Joint Conference on Artificial Intelligence*, Beijing, China, 2013, pp. 2654–2661.
- [107] I. Porada, K. Suleman, A. Trischler, J. C. K. Cheung, Modeling event plausibility with consistent conceptual abstraction, in: *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Virtual Event, 2021, pp. 1732–1743.
- [108] Y. Elazar, H. Zhang, Y. Goldberg, D. Roth, Back to square one: Artifact detection, training and commonsense disentanglement in the winograd schema, in: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Virtual Event, 2021, pp. 10486–10500.