

C³KG: A Chinese Commonsense Conversation Knowledge Graph

Anonymous ACL submission

Abstract

Existing commonsense knowledge bases often organize tuples in an isolated manner, which is deficient for commonsense conversational models to plan the next steps. To fill the gap, we curate a large-scale multi-turn human-written conversation corpus, and create the first Chinese commonsense conversation knowledge graph which incorporates both social commonsense knowledge and dialog flow information. To show the potential of our graph, we develop a graph-conversation matching approach, and benchmark two graph-grounded conversational tasks. All the resources in this work will be released to foster future research.

1 Introduction

Commonsense knowledge describes facts and related judgments in our everyday world, which is essential for machine when interacting with humans. These years have witnessed a growing number of literature incorporating commonsense knowledge into various downstream tasks (Bauer et al., 2018; Chen et al., 2019; Lin et al., 2019; Guan et al., 2019; Ji et al., 2020).

Recently, Sap et al. (2019) curate ATOMIC, a large-scale commonsense knowledge base, which covers event-centered social aspects of inferential knowledge tuples. For example, there exist tuples like $\{PersonX \text{ adopts a cat}, xEffect, happy\}$ and $\{PersonX \text{ adopts a cat}, xWant, company\}$. Here, $xEffect$ and $xWant$ are two of nine relations defined in ATOMIC to infer people’s mental states for a given event, e.g., $PersonX \text{ adopts a cat}$. As such, it is promising to detect ATOMIC events mentioned in conversations, and utilize the inferred knowledge when developing social chatbots.

In spite of the potential, it has two major difficulties. For instance, when a friend in *distress* tells us that he recently adopted a cat, we humans will easily suspect that he might has allergies to the cat. However, such reasoning is difficult for

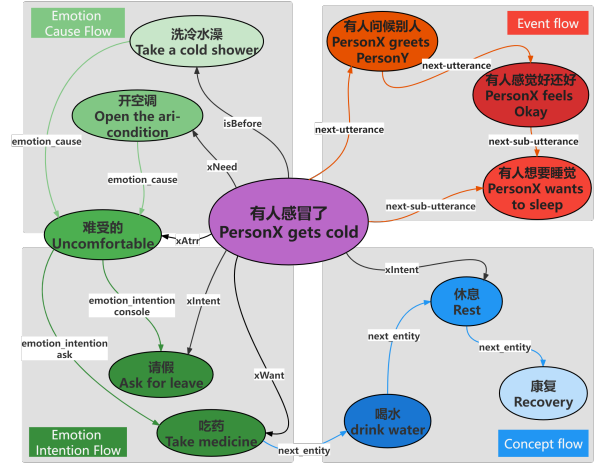


Figure 1: A tiny subset of C³KG, with four unique types of dialog flow relations.

chatbots. Given the event-relation pair $\{PersonX \text{ adopts a cat}, xEffect, __\}$, ATOMIC contains multiple tails like $\{finds \text{ out he has allergies}\}$ and the tail $\{becomes \text{ less lonely}\}$. To this end, the first difficulty comes from the existence of multiple tails, which will confuse the chatbots when inferring the cause behind the negative emotion. Secondly, the knowledge tuples in ATOMIC are isolated. It is thus more difficult for the chatbots to reason which tail(s) of knowledge should be used to produce coherent responses. For example, if the tuple $\{PersonX \text{ adopts a cat}, isAfter, finds \text{ a cat at the animal shelter}\}$ is detected from the dialogue history, then the tuple $\{PersonX \text{ adopts a cat}, xNeed, go \text{ to an animal rescue center}\}$ should not be considered anymore for future conversations. We argue that these issues hamper the application of ATOMIC to multi-turn dialogue modeling where the conversational agents need not only know the current state but also plan the future **dialog flow**.

To remedy these issues, we define 4 novel dialog flow relations, i.e., event flow, concept flow, emotion-cause flow, emotion-intent flow, as de-

064 pictured in Figure 1. To build up the relations, we
065 collect a large-scale multi-turn conversations in
066 everyday scenarios, and manually annotate the con-
067 versations with emotional information. Based on
068 the annotations, we are able to extract conversation-
069 related events in ATOMIC and connect them using
070 different dialog flows. In this way, we augment
071 ATOMIC with conversation-specific knowledge,
072 which facilitates chatbots to pick out useful comm-
073 monsense knowledge, and relieves their confusion
074 on noisy knowledge that are incoherent with dia-
075 log flows. We believe our graph is favorable for
076 commonsense conversation modeling.

077 To highlight: (1) We curate a new Chinese cor-
078 pus, containing multi-turn human-written conver-
079 sations on dailylife topics and rich, high-quality
080 annotations on the level of sub-utterance; (2)
081 We create and will release the first large-scale
082 Chinese commonsense conversation knowledge
083 graph, C^3KG , which contain 4 types of unique
084 dialog-flow edges to store the distilled conversation
085 knowledge from the multi-turn conversation cor-
086 pus; (3) We devise a graph-conversation matching
087 approach, and benchmark 2 typical tasks grounded
088 on commonsense conversation graph.

089 2 Related Work

090 2.1 Commonsense Knowledge Bases

091 ConceptNet (Speer et al., 2017a) is a popular com-
092 monsense knowledge base is, which has a Chi-
093 nese version with a relatively small set of knowl-
094 edge (Kuo et al., 2009). Another large-scale com-
095 monsense knowledge graph TransOMCS (Zhang
096 et al., 2020) is built automatically by converting
097 syntactic parses of Web sentences into structured
098 knowledge. However, the majority of relations
099 in these knowledge bases are taxonomic relations
100 such as *isA* and *Synonym* (Davis and Marcus,
101 2015), which inevitably limits their capabilities.
102 Differently, we rely on ATOMIC (Sap et al., 2019).
103 Despite the lack of Chinese version, ATOMIC cov-
104 ers unique mental knowledge. We thus translate
105 it into Chinese and build dialog flow relations on
106 it. Other Chinese knowledge bases include but
107 not limited to CN-DBPedia (Xu et al., 2017) and
108 zhishi.me (Niu et al., 2011).

109 2.2 Extracting Knowledge from Conversation

110 To extract structured knowledge from conversa-
111 tions, previous works detect named entities from
112 each utterance in conversational datasets (Xu et al.,

2020c; Zou et al., 2021a; Ghosal et al., 2021) and build up the relationship based on their se-
quential order and Pointwise Mutual Information (PMI) (Church and Hanks, 1990). There also exists
some works use automatic extraction tools, such as OpenIE, to construct conversational knowledge
bases of certain domains (Ahmad et al., 2020). Although plausible, these knowledge graphs are built
on the granularities of word or phrase, which makes them hard to match the overall semantics of dia-
logue sentences. In this paper, we build a Chinese commonsense conversation knowledge graph
based on both multi-turn conversational corpus and event-centered knowledge base. At the same time,
we propose to use Sentence-BERT (Reimers and Gurevych, 2019a), a transformer-based semantic
similarity model, to construct dialog flow edges in our knowledge graph.

131 2.3 Knowledge Grounded Dialogue Modeling

132 There are growing interests in incorporating com-
133 monsense knowledge into dialogue tasks. Both
134 Zhou et al. (2018) and Zhang et al. (2019) intro-
135 duce knowledge triplets from ConceptNet (Speer
136 et al., 2017b) into open-domain response genera-
137 tion. Recently, Li et al. (2021a) and Zhong et al.
138 (2021) exploit ConceptNet to enhance emotion rea-
139 soning for response generation, and others design
140 graph reasoning methods to plan the topic tran-
141 sition in the responses (Moon et al., 2019; Tang
142 et al., 2019; Xu et al., 2020a; Li et al., 2021c).
143 One distinct work is Ghosal et al. (2020), which
144 utilizes ATOMIC (Hwang et al., 2020) in emo-
145 tional dialogue modeling for emotion identification.
146 In this paper, we connect the heads and tails in
147 ATOMIC according to four types of dialog flows.
148 Because the resulted graph C^3KG contains both so-
149 cial knowledge from ATOMIC and dialogue knowl-
150 edge from our corpus, it is thus more suitable for
151 empathetic conversation modeling.

152 3 A Scenario-based Multi-turn 153 Conversation Corpus

154 Our aim is to extract common dialog flow infor-
155 mation from real conversations. In this way, it is
156 crucial to ensure the quality of the conversation cor-
157 pus and the reliability of the extraction method. In
158 the following, we firstly introduce the conversation
159 corpus $CConv$ we depend on.

160 Instead using the noisy Internet data, we col-
161 lect a multi-turn human-written Chinese conversa-

tion corpus based on crowdsourcing. Initially, 100 workers are hired, and they are randomly paired to talk in text under a given scenario. Each scenario is one sentence describing the suggested conversation context which often involves certain everyday events. Besides, the workers are also required to follow certain rules like “each utterance should longer than 6 Chinese characters”, which are critical to help ensure the quality of the collected conversation. At the beginning of the crowdsourcing, we check each collected conversation and re-train the workers. To ensure the quality, we keep only 62 well-trained workers and let them finish our task. Note that the workers are paid with 1 CNY per utterance (nearly 0.2 dollar per utterance). Finally, we obtain 32k sessions of high-quality two-party conversations (650k utterances in total) on 200 scenarios of 15 daily topics.

To facilitate future research, we then hire another 3 well-trained assistants to manually annotate the conversations with fine-grained emotional labels including speaker’s emotion type, emotion cause, and response intention type. Following Rashkin et al. (2019), we define emotion type with 5 general classes {joy, angry, sad, caring, other}. Emotion cause span is a continuous text spans implying the reason of certain emotion (Li et al., 2021b). Response intention type is essential for building empathetic chatbots, and we define 6 commonly-adopted intent classes of {ask, advise, describe, opinion, console, other} following Welivita and Pu (2020). A snippet of an conversation example is given in Figure 2. In Appendix, we present more information of the constructed corpus.

By utilizing the annotations, we are able to distill dialogue knowledge to enhance the conversation graph and graph-grounded conversation modeling.

4 Translation of ATOMIC

Because our conversation corpus is Chinese, we want to build a Chinese conversation knowledge graph. It is well known that to build a knowledge graph from scratch is laborious and time-consuming. Instead, we base on ATOMIC and design a pipeline method to translate it into Chinese, meanwhile ensuring the resulted knowledge graph is reliable and suitable for conversation grounding.

4.1 Brief Introduction of ATOMIC

Before describing our detailed processing steps, we firstly give a brief description of ATOMIC (Sap

et al., 2019). ATOMIC organizes commonsense knowledge in the form of triplet $\langle \text{head}, \text{relation}, \text{tail} \rangle$, where head often describes a daily event.

There are two unique properties making ATOMIC suitable and attractive for building empathetic chatbots. Firstly, ATOMIC collects knowledge about how people will react to a given event. This kind of knowledge is related to people’s mental states, which is beneficial for understanding implicit emotions. For example, given a head event *PersonX makes PersonY’s coffee*, ATOMIC contains knowledge that PersonY will be *grateful* along the relation `oReact`. Secondly, ATOMIC organizes knowledge using several inferential relations and naturally supports *if-then* reasoning, which is crucial generating coherent responses. Totally, there are 9 relations defined in ATOMIC. The details can be found in Appendix.

4.2 Replacement of Certain Tokens

We begin with translating high-frequency patterns in the original triplets. As compared to the pre-defined set of relations, it is more difficult to handle the heads and tails. In ATOMIC, for example, there exist **185,046** heads and tails containing tokens like “*PersonX*” and “*PersonY*”. These personal pronouns stand for the givers and the receives for a certain event, and can be regarded as the speech parties in a conversation. Also, some ATOMIC heads like {*PersonX gets _____ as a pet*}, have a blank which can be filled with various tokens.

These aforementioned patterns bring ambiguity to the triplet semantics, and will confuse the translation model. To address, we devise a series of replacement rules to keep the original semantics while translation. For example, for the ATOMIC head *PersonX votes for personY*, we convert it to be “Someone votes for someone else” and send it to our translation model.

4.3 Joint Translation of Head and Tail

Nevertheless, the majority of the heads and tails in ATOMIC are short phrases, while machine translation models are often context-based. The multi-sense characteristics of language will further deteriorate the translation quality if we separately feed each single head and tail to a translation model.

To remedy the issues, we instead translate the head and tail in each triplet together. Given a triplet $\langle h, r, t \rangle$, we connect the head h with its t using a heuristic connecting word r' w.r.t. the relation r , and obtain one long sentence l . After translating

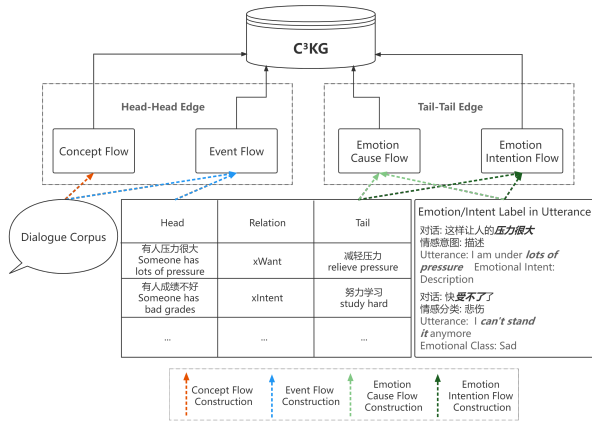


Figure 2: Construction Process of C³KG.

the long text, we split the translation result with the connecting word and turn it into h_{tr} and t_{tr} :

$$\begin{aligned}
 l &= \text{CONNECT}(h, r', t) \\
 l'_{tr} &= \text{TRANSLATION}(l) \quad (1) \\
 h_{tr}, r_{tr}, t_{tr} &= \text{SPLIT}(m'_{tr}, r'_{tr})
 \end{aligned}$$

where the resulted $\langle h_{tr}, r_{tr}, t_{tr} \rangle$ is the translated triplets. And CONNECT, SPLIT denote the corresponding operation. TRANSLATION stands for our translation model. By this means, we expect the connected l provides more contextual information for better semantic translation. The comparison results between separate translation and joint translation will be given in Section 6.3.1.

Note that auxiliary translation methods can be used. In this work, we use Xiaomi commercial Translation service.¹ For simplicity, we denote the translated ATOMIC as **ATOMIC-zh**.

5 Conversation Knowledge Graph Construction

5.1 Overview of C³KG

To supply dialog flow information for commonsense reasoning, we create a Chinese Commonsense Conversation Knowledge Graph, **C³KG**, whose statistics are summarized in below.

We then introduce our method of constructing a conversational knowledge graph based on ATOMIC-zh and our multi-turn conversation corpus. In general, we extract events from each conversations and match with the head in ATOMIC-zh. The core is how to build new dialog flow relations, which is depicted in Figure 2, and will be detailed present in the following section.

¹<http://fanyi.mioffice.cn>

#Relations	ATOMIC Relations	63,6656
	Event Flows	57,1196
	Concept Flows	7,7587
	Emotion-Cause Flows	187
#Triplets	Emotion-Intent Flows	196
		1285,822

Table 1: Statistics of C³KG.

5.2 Event Extraction

Knowledge in ATOMIC-zh is event-based and most of them are declarative sentences with some entities omitted. However, utterances in the open-domain dialogue dataset contain a lot of colloquial expressions and sub-sentences with more complex structures. To cope with the complexity, we develop a dependency parsing-based event detection pipeline to extract salient events in each utterance. The overview of our algorithm is described in Algorithm 1.

Pre-processing. We first split each utterance with punctuation, and operate on the level of sub-utterances. To reduce noise, we then filter short sub-utterances with transitive and dumb semantics like “好的” (OK), “就是这样” (That’s it). After that, we perform Dependency Syntactic Parsing and POS tagging using ltp4², and extract event mentions based on two kinds of structural patterns, verb-driven and adjective-driven clauses.

Verb-driven. Verb-driven clauses have a verb connecting to the root node in the dependency tree. After filtering some noisy words, we obtain verb-driven event mentions. For example, we extract the mention “催促提供物资的商家” (urged the merchants who provide supplies) from utterance “我和上司已经在催促提供物资的商家了” (My boss and I have already urged the merchants who provide supplies). In this utterance, we filter subject of utterance“我和上司” (My boss and I), adverbial“已经” (have already) and modal particle“了” (yet) at the end of the utterance.

Adjective-driven. Besides, adjective-driven clauses often have meaningful entities in sub-utterances. Similarly, we extract adjective-driven event mentions based on the adjective-driven clauses by keeping the modifier of its key adjective and filtering out other words. For example, we extract the mention “学习节奏快” (The pace of learning is fast) from the utterance “但学习节奏也太快了吧” (But the pace of learning is too fast).

²<https://github.com/HIT-SCIR/ltp>

Algorithm 1 Event Extraction from Utterance

Input: An utterance U **Output:** A set of event mentions M

- 1: Split U with punctuation, and get a series of sub-utterance SU , filter SU based on length
 - 2: **for** each $su \in SU$ **do**
 - 3: Obtain the dependency tree dep and POS tagging result pos of su
 - 4: Find the had node which connects directly to the $ROOT$ node in the dependency tree
 - 5: **if** POS tag of the had node $\in [v, a]$ **then**
 - 6: Append had to HAD
 - 7: **end if**
 - 8: **if** The number of verbs connected directly to had more than 1 **then**
 - 9: Recursively search verbs in the sub-tree of had and replace had in HAD with the founded verbs
 - 10: **end if**
 - 11: **for** $had \in HAD$ **do**
 - 12: **if** POS of node had is v **then**
 - 13: Keep words in su that appear after had and words connect directly to had and relation is ‘ADV’, connect them and append to M
 - 14: **else**
 - 15: Remain words in su that connect directly to had and relation is ‘SBV’, connect them and append to M
 - 16: **end if**
 - 17: **end for**
 - 18: **end for**
 - 19: Return M
-

332 In this utterance, we filter the initial conjunction
333 “但是” (but), adverbial “也” (no meaning) and “太”
334 (too) and modal particle “了” (yet) and “吧” (no
335 meaning) at the end of the utterance.

336 **Recursive Applying.** The resulted event mentions
337 may still contain multiple verbs and several seman-
338 tic units. In this case, we apply a secondary de-
339 composition. For example, we will split the event
340 mention “以为进了大学就可以放松放松” (could
341 relax after entering university) into two events “进
342 了大学” (entering university) and “就可以放松放
343 松” (could relax). To do so, we count the number
344 of verbs connected to the root word in the mention
345 as well as the depth of the sub-trees led by those
346 verbs. Based on the results, we determine whether
347 the mention needs a secondary decomposition us-

ing a threshold. If needed, we recursively search
348 verbs in the original dependency tree and replace
349 the key verb with the verbs we found.
350

5.3 Event Linking as Matching

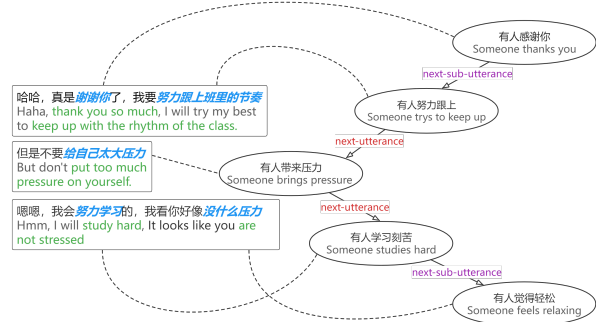


Figure 3: An Example of Head-Head Edge Construction for Event Flows.

In order to discover common dialog flows among
352 the knowledge base, the event mentions in the con-
353 versations are then linked to ATOMIC heads using
354 matching techniques.
355

356 Typically, we adopt Sentence-BERT, a power-
357 ful semantic matching model, which is based on
358 Siamese and Triplet Network and pre-trained on
359 sentence pairs in different relationships (Reimers
360 and Gurevych, 2019b). It encodes two given sen-
361 tences separately and calculates the similarity be-
362 tween their representations, and thus performing
363 efficiently in large-scale many-to-many matching.

364 To enhance the matching performance, we fine-
365 tune Sentence-BERT on our corpus. Specifically,
366 we randomly select 8,000 $\langle m, h \rangle$ mention-head
367 pairs matched by pre-trained Sentence-BERT, and
368 manually label a matching score in $\{0,1\}$ for fine-
369 tuning. Note the reason why we adopt discrete
370 $\{0,1\}$ instead of continuous $[0, 1]$ scores is that us-
371 ing the former effectively mitigates the domain gap.
372 It will induce the matching model to label 0 for
373 those $\langle m, h \rangle$ share similar characters in surface
374 but different meanings in semantics. After fine-
375 tuning, we calculate the cosine similarity scores
376 and choose the head with the highest score as the
377 matching result given an event mention.

5.4 Edge Construction

378 Now we have 32k sessions of multi-turn conver-
379 sations and link their event mentions to ATOMIC
380 heads. The remaining is how to utilize them and
381 build commonsense conversation knowledge graph.
382

In this work, we propose three kinds of edges to reflect different types of dialog flows.

5.4.1 Head-Head Edge Construction

Event Flow. Naturally, a dialogue is hierarchical in that it consists of a sequence of utterances produced by two interlocutors, where each utterance is composed of one or several sub-utterances. If two event mentions are detected together within in a conversation, the co-occurrence can be regarded as a dialog flow example. Following the flow, it is then intuitive to connect the ATOMIC heads linked by the mentions, as illustrated in Figure 3. By connecting intra-utterance and inter- utterance mentions, we acquire the event flows of `next-sub-utterance` and `next-utterance`.

Concept Flow. ATOMIC also has entity-level heads in addition to the phrase-level events. To utilize them, we perform entity linking by detecting word entities with POS tag belonging to {verb, noun, adjective} in the original conversations, and match them with the entity-level AOMIC heads to construct concept flow edges similarly. These concept flows are helpful for planning and transiting the contents in topic-aware conversation (Yao et al., 2018; Moon et al., 2019; Xu et al., 2020b; Zou et al., 2021b).

Because we are interested in the most common dialog flows, we only keep those highly-frequent connections, and create a head-to-head dialog flow between the ATOMIC head entities and events.

5.4.2 Tail-Tail Edge Construction

Besides, we also consider another essential type of dialog flow, i.e., emotion-based empathy flow. In this paper, we utilize the emotional labels on our corpus (in Section 3) to construct two kinds of emotion-based edges connecting tails in our knowledge graph. Intuitively, `emotion-cause` dialog flow reflects the reasons for a specific emotion, which is useful for fine-grained emotion understanding. And `emotion-intent` empathy flow indicates what response intentions are proper to use when the other one is in a specific emotion, which is critical for response empathy.

Pre-processing. To construct emotion-based edges, we category the tails into 3 classes according to their connecting relations, as listed in Table 2. The first class of tails are linked by relations $xAttr$ or $xReact$, which reflects people’s psychological reaction towards a certain event (head). For instance, {*PersonX runs out of steam*, $xAttr$, *tired*}

$Tail_{emotion}$	$xAttr, xReact$
$Tail_{before}$	$isAfter, xNeed$
$Tail_{after}$	$isBefore, xWant, xIntent, xEffect, oEffect$

Table 2: Relation Categories For Emotion-based Edge Construction.

indicates that someone is lacking energy. We denote the first class as $Tail_{emotion}$. The second class $Tail_{before}$ states the events commonly happen before the heads, e.g., {*PersonX runs out of steam*, $isAfter$, *PersonX exercises in gym*}. On the contrary, the last class $Tail_{after}$ contain the events following the head events like {*PersonX runs out of steam*, $xWant$, *to get some energy*}.

By analyzing these relations and tails, we find heuristics to build emotion-based dialog flows. By connecting the head and tails in class $Tail_{emotion}$, we are able to create causal emotional inference like {*PersonX exercises in gym*, `emotion-cause`, *tired*}. Through cross linking the tails in class $Tail_{emotion}$ and $Tail_{after}$, we are able to develop the inferential edges like {*tired*, `emotion-intent`, *to get some energy*}.

Filtering. Based on the heuristics, we apply Sentence-BERT³ to match each tail in class $Tail_{emotion}$ to one of 4 emotion labels defined in our dataset, i.e., {joy, sad, angry, caring}, and set a threshold of 0.7 to determine the emotion class of the tails. The tails sharing the same emotion class with the original utterance are kept to build emotion-based dialog flows.

Emotion Cause Flow. Then, we apply keyword-based exact matching between the tails in $Tail_{before}$ with dialogue context. For $Tail_{before}$, if there is an keyword exactly matched with some keywords in the previous utterances, we create an `emotion-cause` edge flowed from the tail of $Tail_{before}$ to those filtered tails in $Tail_{emotion}$, indicating that the event of $Tail_{before}$ may cause person to feel the emotion of the tail in $Tail_{emotion}$.

Emotion Intention Flow. For tails in class $Tail_{after}$, we create an `emotion-intent` flow from those filtered tails in $Tail_{emotion}$ to the tails in $Tail_{after}$. Notably, we also assign one of five intent labels to each `emotion-intent` edge, i.e., {ask, advise, describe, opinion, console} (Section 3).

Figure 4 depicts the process of constructing the labeled `emotion-intent` edge. We start by

³This model is not fine-tuned on our dataset.

Method	Fluency	Logic
Separate translation	0.825	0.71
Joint translation	0.92	0.88

Table 4: Evaluation of Translation Quality.

on both Fluency and Logic aspects clearly demonstrate the superiority of joint translation method. In terms of logical coherence, we find many sample cases are labeled with 0 logical score due to the incompleteness of their heads, which somehow confuses the semantics and obstacles logical connection to the tails. For example, {有人把他父亲, xAttr, 告密者} ({PersonX gets PersonX's father, xAttr, a tattletale}) seems ridiculous. However, if we add 出卖 (betrayed) in the end of the heads, then it becomes complete and logical. Nonetheless, such seemingly illogical knowledge might still be informative for downstream tasks with fuzzy matching techniques. Hence, we retain this kind of incomplete heads.

6.3.2 Edge Evaluation

At the heart of C³KG is the novel dialog flow relations we develop in this work. To validate the quality of these relations, we utilize another open-domain multi-turn Chinese dialogue dataset, MOD (Fei et al., 2021)⁴. In specific, we extract event mentions from MOD utterances and match them to our graph using the methods as in Section 5.2. Then we evaluate the connectivity and distance of the matched results, w.r.t. both next_utterance and next_sub_utterance relations. This aims to assess the aggregation degree of related content in our knowledge graph.

KG	next_utterance		next_sub_utterance	
	Connectivity	Dist.	Connectivity	Dist.
C ³ KG	45.50	2.08	12.96	2.13

Table 5: Edge Evaluation Result on MOD dataset.

Table 5 shows our edge evaluation result on MOD. While the connectivity of the matching node between utterance is quite high, we find that the connectivity of matching node within the same utterance is relative fair. This result gives us inspiration in the future study to enlarge window size to find more latent event-level transfer within the utterance.

⁴<https://github.com/lizekang/DSTC10-MOD>

Method	Emotion	Intent
Base	90.7	65.3
Knowledge	91.4	67.3
History	90.5	64.7
Knowledge+History	91.2	65.4

Table 6: Baselines for Graph-grounded Tasks.

7 Proposed Tasks

To show the potential, we propose two graph-grounded conversational tasks, i.e., emotion classification and intent prediction, and train benchmark models using our labeled corpus CConv.

Task 1: Emotion Classification requires to produce an emotion label conditions on the conversations. Following common practice, we choose the BERT model, and sample the xAttr, xReact tails from our matching head as extra input.

Task 2: Intent Prediction requires to predict a proper type of response intent for the conversations. We choose BERT model, and sample the oReact, oEffect tails from our matching heads. As simple baselines, we introduce history and graph knowledge through concatenation with an input format as U_{i-2} [SEP] U_{i-1} [SEP] U_i [SEP] oReact tail [SEP] oEffect tail.

The accuracies of baseline methods are reported in Table 6. *Base* denotes only using the utterance to do prediction. *Knowledge* and *History* denote whether to add knowledge we sampled and dialogue history to the model. While adding knowledge improves the emotion classification and intent prediction performances, it seems problematic to directly concatenating history dialogues, which may bring noises. The moderate scores also indicate that there is still a room to improve for graph-grounded conversation understanding.

8 Discussions of Future Work

In this work, we provide a systematic approach from event mention detection, event linking to conversation graph construction which consists of 4 distinguished types of dialog flows. For each step, there exist possible refinements. For example, we plan to include other event-based resources to improve graph-conversation matching accuracy as well as the graph knowledge coverage.

We also plan to continue the annotations to supply more dialog flow information especially those empathy ones, and evaluate more dialog flow relations on other datasets. Ethical statements are given in Appendix.

References

- Zishan Ahmad, Asif Ekbal, Shubhashis Sengupta, Anutosh Mitra, Roshni Rammani, and Pushpak Bhattacharyya. 2020. Active learning based relation classification for knowledge graph construction from conversation data. In *International Conference on Neural Information Processing*, pages 617–625. Springer.
- Lisa Bauer, Yicheng Wang, and Mohit Bansal. 2018. Commonsense for generative multi-hop question answering tasks. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4220–4230, Brussels, Belgium. Association for Computational Linguistics.
- Jiaao Chen, Jianshu Chen, and Zhou Yu. 2019. Incorporating structured commonsense knowledge in story completion. In *AAAI*.
- Kenneth Church and Patrick Hanks. 1990. Word association norms, mutual information, and lexicography. *Computational linguistics*, 16(1):22–29.
- E. Davis and G. Marcus. 2015. Commonsense reasoning and commonsense knowledge in artificial intelligence. *Communications of the ACM*, 58:92 – 103.
- Zhengcong Fei, Zekang Li, Jinchao Zhang, Yang Feng, and Jie Zhou. 2021. Towards expressive communication with internet memes: A new multimodal conversation dataset and benchmark. *arXiv preprint arXiv:2109.01839*.
- Deepanway Ghosal, Pengfei Hong, Siqi Shen, Navonil Majumder, Rada Mihalcea, and Soujanya Poria. 2021. Cider: Commonsense inference for dialogue explanation and reasoning. *arXiv preprint arXiv:2106.00510*.
- Deepanway Ghosal, Navonil Majumder, Alexander Gelbukh, Rada Mihalcea, and Soujanya Poria. 2020. Cosmic: Commonsense knowledge for emotion identification in conversations. *arXiv preprint arXiv:2010.02795*.
- Jian Guan, Yansen Wang, and Minlie Huang. 2019. Story ending generation with incremental encoding and commonsense knowledge. In *AAAI*.
- Jena D Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. 2020. Comet-atomic 2020: On symbolic and neural commonsense knowledge graphs. *arXiv preprint arXiv:2010.05953*.
- Haozhe Ji, Pei Ke, Shaohan Huang, Furu Wei, Xiaoyan Zhu, and Minlie Huang. 2020. Language generation with multi-hop reasoning on commonsense knowledge graph. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 725–736, Online. Association for Computational Linguistics.
- Yen-Ling Kuo, Jong-Chuan Lee, Kai yang Chiang, Rex Wang, Edward Shen, Cheng wei Chan, and Jane Yung jen Hsu. 2009. Community-based game design: experiments on social games for commonsense data collection. In *HCOMP '09*.
- Qintong Li, Piji Li, Zhumin Chen, and Zhaochun Ren. 2021a. Towards empathetic dialogue generation over multi-type knowledge. In *AAAI*.
- Yanran Li, Ke Li, Hongke Ning, Xiaoqiang Xia, Yalong Guo, Chen Wei, Jianwei Cui, and Bin Wang. 2021b. Towards an online empathetic chatbot with emotion causes. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2041–2045.
- Yanran Li, Wenjie Li, and Zhitao Wang. 2021c. Graph-structured context understanding for knowledge-grounded response generation. *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- Bill Yuchen Lin, Xinyue Chen, Jamin Chen, and Xiang Ren. 2019. Kagnet: Knowledge-aware graph networks for commonsense reasoning. *ArXiv*, abs/1909.02151.
- Seungwhan Moon, Pararth Shah, Anuj Kumar, and Rajen Subba. 2019. Opendialkg: Explainable conversational reasoning with attention-based walks over knowledge graphs. In *ACL*.
- Xing Niu, Xinruo Sun, Haofen Wang, Shunlin Rong, Guilin Qi, and Yong Yu. 2011. Zhishi.me - weaving chinese linking open data. In *SEMWEB*.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic open-domain conversation models: A new benchmark and dataset. In *ACL*.
- Nils Reimers and Iryna Gurevych. 2019a. Sentencebert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.
- Nils Reimers and Iryna Gurevych. 2019b. SentenceBERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A. Smith, and Yejin Choi. 2019. ATOMIC: an atlas of machine commonsense for if-then reasoning. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 3027–3035. AAAI Press.

730	R. Speer, Joshua Chin, and Catherine Havasi. 2017a.	Yicheng Zou, Zhihua Liu, Xingwu Hu, and	781
731	Conceptnet 5.5: An open multilingual graph of gen-	Qi Zhang. 2021a. Thinking clearly, talking fast:	782
732	eral knowledge. In <i>AAAI</i> .	Concept-guided non-autoregressive generation for	783
		open-domain dialogue systems. <i>arXiv preprint</i>	784
733	Robyn Speer, Joshua Chin, and Catherine Havasi. 2017b.	<i>arXiv:2109.04084</i> .	785
734	Conceptnet 5.5: An open multilingual graph of gen-		
735	eral knowledge. In <i>Thirty-first AAAI conference on</i>	Yicheng Zou, Zhihua Liu, Xingwu Hu, and Qi Zhang.	786
736	<i>artificial intelligence</i> .	2021b. Thinking clearly, talking fast: Concept-	787
		guided non-autoregressive generation for open-	788
737	Jianheng Tang, Tiancheng Zhao, Chenyan Xiong, Xi-	domain dialogue systems. In <i>Proceedings of the</i>	789
738	aodan Liang, Eric P. Xing, and Zhiting Hu. 2019.	<i>2021 Conference on Empirical Methods in Natural</i>	790
739	Target-guided open-domain conversation. <i>ArXiv</i> ,	<i>Language Processing</i> , pages 2215–2226, Online and	791
740	abs/1905.11553.	Punta Cana, Dominican Republic. Association for	792
		Computational Linguistics.	793
741	A. Welivita and Pearl Pu. 2020. A taxonomy of empa-		
742	thetic response intents in human social conversations.		
743	In <i>COLING</i> .		
744	Bo Xu, Yong Xu, Jiaqing Liang, Chenhao Xie, Bin		
745	Liang, Wanyun Cui, and Yanghua Xiao. 2017. Cn-		
746	dbpedia: A never-ending chinese knowledge extrac-		
747	tion system. In <i>IEA/AIE</i> .		
748	J. Xu, Haifeng Wang, Zheng-Yu Niu, Hua Wu, and		
749	Wanxiang Che. 2020a. Knowledge graph grounded		
750	goal planning for open-domain conversation genera-		
751	tion. In <i>AAAI</i> .		
752	J. Xu, Haifeng Wang, Zheng-Yu Niu, Hua Wu, Wanxi-		
753	ang Che, and Ting Liu. 2020b. Conversational graph		
754	grounded policy learning for open-domain conversa-		
755	tion generation. In <i>ACL</i> .		
756	Jun Xu, Haifeng Wang, Zheng-Yu Niu, Hua Wu, Wanxi-		
757	ang Che, and Ting Liu. 2020c. Conversational graph		
758	grounded policy learning for open-domain conversa-		
759	tion generation. In <i>Proceedings of the 58th Annual</i>		
760	<i>Meeting of the Association for Computational Lin-</i>		
761	<i>guistics</i> , pages 1835–1845.		
762	Lili Yao, Ruijian Xu, C. Li, Dongyan Zhao, and Rui		
763	Yan. 2018. Chat more if you like: Dynamic cue		
764	words planning to flow longer conversations. <i>ArXiv</i> ,		
765	abs/1811.07631.		
766	Hongming Zhang, Daniel Khashabi, Yangqiu Song, and		
767	Dan Roth. 2020. Transomcs: From linguistic graphs		
768	to commonsense knowledge. In <i>IJCAI</i> .		
769	Houyu Zhang, Zhenghao Liu, Chenyan Xiong, and		
770	Zhiyuan Liu. 2019. Grounded conversation genera-		
771	tion as guided traverses in commonsense knowledge		
772	graphs. <i>arXiv preprint arXiv:1911.02707</i> .		
773	Peixiang Zhong, Di Wang, Pengfei Li, Chen Zhang,		
774	Hao Wang, and Chunyan Miao. 2021. CARE:		
775	commonsense-aware emotional response generation		
776	with latent concepts. In <i>AAAI</i> .		
777	Hao Zhou, Tom Young, Minlie Huang, Haizhou Zhao,		
778	Jingfang Xu, and Xiaoyan Zhu. 2018. Commonsense		
779	knowledge aware conversation generation with graph		
780	attention. In <i>IJCAI</i> , pages 4623–4629.		

A Ethical Considerations

At last, we discuss the potential ethic impacts of this work. (1) **Transparency**: We will release the newly introduced corpus and the built conversation knowledge graph, as well as the benchmark approaches to facilitate future research. Similar datasets and knowledge bases include Empathetic Dialogues (Rashkin et al., 2019) and ATOMIC (Sap et al., 2019), which are often public available and have been used extensively. (2) **Privacy**: The corpus is crowdsourced under a set of specific rules to forbid the workers disclosure sensitive and personal identifiable information. (3) **Politeness**: Because our conversations are human-written and are related to healthy daily life scenarios, they are expected to be clean, legal, and polite. The crowdsourcing rules are designed to avoid emotionally triggering words as much as possible.

B Corpus: CConv

B.1 Example & Statistics

In our corpus CConv, conversations are conducted based on a scenario between two parties. Table 8 gives an example conversation. The statistics of CConv is also present in Table 7. Since there are 200 scenarios in total, and hence we have 160 diverse multi-turn conversations in average.

# sessions of dialogues	32,612
# utterances	650,147
# unique scenarios	200
# conversation topics	15
Avg. # words per utterance	7.8
Avg. # turns per dialogue	19.9

Table 7: The Statistics of the Corpus CConv.

B.2 Topics and Scenarios

To ensure the diversity of the conversations, we select 15 everyday topics. For each topic, we manually write tens of one-sentence scenario to guide the conversation context.

In total, we have 15 topics and 200 scenarios. To better understand, we show some example topics and scenarios in Table 9.

B.3 Annotation Criteria

To facilitate future research, we hire another 3 well-trained assistants to manually annotate the conversations with fine-grained emotional labels including speaker’s emotion type, emotion cause, and

response intention type. The annotation example is given along with the example in Table 8.

Emotion Class. Following Rashkin et al. (2019), we define emotion type with 5 general classes {joy, angry, sad, caring, other}.

Emotion Cause Span. Emotion cause span is a continuous text spans implying the reason of certain emotion (Li et al., 2021b).

Response Intent. Response intention type is essential for building empathetic chatbots, and we define 6 commonly-adopted intent classes of {ask, advise, describe, opinion, console, other} following Welivita and Pu (2020), which are described in Table 10.

C ATOMIC

In this work, we introduce ATOMIC (Sap et al., 2019) as the commonsense knowledge base due to its attractive properties of mental state inferences and *if-then* causal relations, as analyzed before.

ATOMIC (Sap et al., 2019) is a novel event-centered knowledge graph, consisting of 880K tuples of social commonsense knowledge. Distinguished from ConceptNet (Speer et al., 2017a), there are two unique properties making ATOMIC suitable and attractive for building empathic chatbots. Firstly, ATOMIC collects knowledge about how people will feel and react to a given event. This kind of knowledge is related to people’s mental states, which is beneficial for understanding implicit emotions. For example, given a head event *PersonX makes PersonY’s coffee*, ATOMIC contains knowledge that PersonY will be *grateful* along the relation `oReact`. Secondly, ATOMIC organizes knowledge using several inferential relations and naturally supports *if-then* reasoning, which is crucial generating coherent responses.

Here, we adopt the figures and demonstrations from the original ATOMIC paper (Sap et al., 2019) to present the 9 relations defined in ATOMIC and give some examples in Figure 6 and Table 11.

D Evaluation

D.1 Template-based Event Extraction Methods

To evaluate our matching methods proposed in this work, we randomly choose 100 utterances and compare with several approaches. In specific, we propose a baseline *POS* matching method, which employs POS tagging-based templates to extract events. The templates are given in Table 12.

Situation			
同事之间，一方身体不舒服，另一方表达关心 Acted as colleagues, one person is sick, and the other one cares about his/her health.			
Conversation			
Speaker	Utterance	Emotion	Intent
1	你今天来得比平时晚呀。 <u>是身体不舒服吗?</u> (You are later than normal days. Are you OK?)	caring	ask
2	呜呜，昨晚空调开的太，一大早起来头就特别疼。 (Yesterday the air conditioner was too cold that I had a headache this morning.)	sad	description
1	怪不得，那你吃过感冒药了吗? (I know. Have you taken the medicine?)	caring	ask
2	吃过了，现在已经好多了，就是有点想睡觉。 (Sure. I feel better now, just feel a little bit sleepy.)	other	description
...
2	今天的工作安排多么? (What are today's arrangements?)	other	other
1	我会帮你做的。 <u>你好好休息吧!</u> (I will help finish them. You'd better take a good rest.)	caring	advise
2	真是太感谢你了! (I really appreciate a lot for your help!)	joy	other

Table 8: Example Conversation with Annotations. Note that the underlined words stand for the emotion cause span. Words are shorten due to space limit.

Topic	Scenario
Study	两个学生之间，讨论课业压力大，总是做不完作业 (Between two students, discuss the overload homework) 考研失败，向朋友倾诉自己的伤心和烦恼 (Fail the entry exam of graduate study, express the distress to a friend)
Entertainment	讨论自己最喜欢的一部电影，以及为什么喜欢它 (Discuss one of your favorite films and why) 聊一聊自己曾经单曲循环过的歌曲，以及当时自己的感受 (Talk about a music or a song you have put on repeat all the night)
Love	情侣之间，因为生活作息不一致而吵架闹别扭 (Between a couple, quarrel with the lover due to inharmonious habits) 自己订婚了，激动地与好友分享喜讯 (Being engaged, share the good news to the best friend)

Table 9: Example Topics and Scenarios.

Intent Type	Definition	Example
ask	to know further details or clarify	<i>What happened?</i>
describe	present more details and explain the reasons	<i>I'm sad because I failed the exam.</i>
advise	give explicit solutions	<i>Try to exercise more.</i>
opinion	share own thoughts	<i>I don't like being disturbed after work.</i>
console	pacify others	<i>I hope you'd feel better.</i>
other	-	<i>Goodbye.</i>

Table 10: Annotation Criteria for Response Intent.

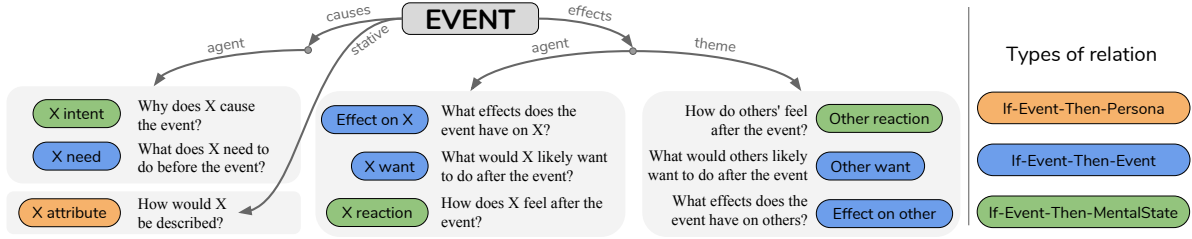


Figure 6: The taxonomy of *if-then* reasoning types. We consider nine *if-then* relations that have overlapping hierarchical structures as visualized above. One way to categorize the types is based on the type of content being predicted: (1) **If-Event-Then-Mental-State**, (2) **If-Event-Then-Event**, and (3) **If-Event-Then-Persona**. Another way is to categorize the types based on their causal relations: (1) “**causes**”, (2) “**effects**”, and (3) “**stative**”. Some of these categories can further divide depending on whether the reasoning focuses on the “agent” (X) or the “theme” (Other) of the event.

Event	Type of relations	Inference examples	Inference dim.
“PersonX pays PersonY a compliment”	If-Event-Then-Mental-State	PersonX wanted to be nice PersonX will feel good PersonY will feel flattered	xIntent xReact oReact
	If-Event-Then-Event	PersonX will want to chat with PersonY PersonY will smile PersonY will compliment PersonX back	xWant oEffect oWant
	If-Event-Then-Persona	PersonX is flattering PersonX is caring	xAttr xAttr
“PersonX makes PersonY’s coffee”	If-Event-Then-Mental-State	PersonX wanted to be helpful PersonY will be appreciative PersonY will be grateful	xIntent oReact oReact
	If-Event-Then-Event	PersonX needs to put the coffee in the filter PersonX gets thanked PersonX adds cream and sugar	xNeed xEffect xWant
	If-Event-Then-Persona	PersonX is helpful PersonX is deferential	xAttr xAttr
“PersonX calls the police”	If-Event-Then-Mental-State	PersonX wants to report a crime Others feel worried	xIntent oReact
	If-Event-Then-Event	PersonX needs to dial 911 PersonX wants to explain everything to the police PersonX starts to panic Others want to dispatch some officers	xNeed xWant xEffect oWant
	If-Event-Then-Persona	PersonX is lawful PersonX is responsible	xAttr xAttr

Table 11: Examples of **If-Event-Then-X** commonsense knowledge present in Sap et al. (2019). For inference dimensions, “x” and “o” pertain to PersonX and others, respectively (e.g., “xAttr”: attribute of PersonX, “oEffect”: effect on others).

882 D.2 More Examples of Constructed Scenario 883 Graphs

contain the full set of C^3KG relations.

891

884 In this section, we visualize more snippets of the
885 scenario graphs. They are “insomnia” in Figure 8,

886 Please kindly note that for clarity, we only visu-
887 alize a small set of relation and tails in each figure,
888 and try to give a comprehensive view of the rela-
889 tions by showing different relations in different
890 scenario graphs. In fact, every scenario graphs

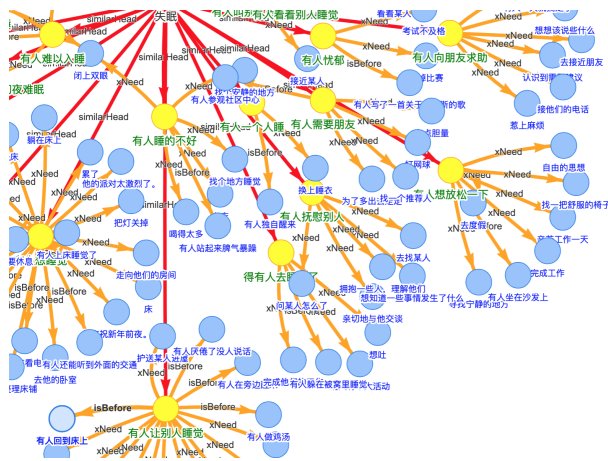


Figure 7: Scenario Graph of “Insomnia”.

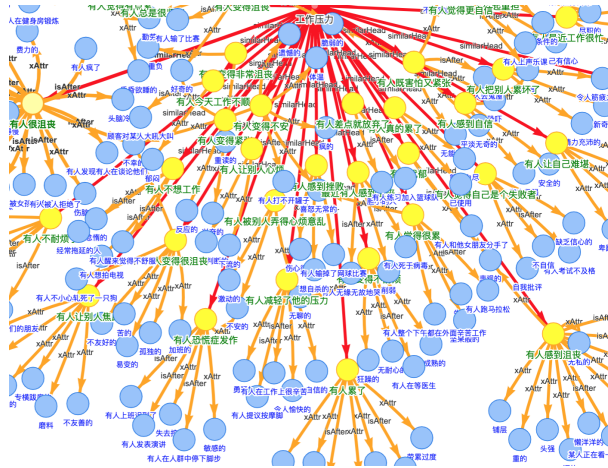


Figure 8: Scenario Graph of “Work Pressure”.

Original pattern	Replaced pattern
PersonX...PersonX...	Someone...himself...
PersonX...PersonY...	Someone...some one else...
PersonX...PersonX’s...	Someone...his...
PersonX...PersonY’s...	Someone...someone else’s
...something...

Table 13: Pattern replacement we use when translating ATOMIC

POS sequence	Example
v+v	想睡觉 (want to sleep)
v+n	做作业 (do homework)
v+i	感觉如释重负 (feel relieved)
v+u+z	跑得飞快 (run fast)
v+u+m	看了一下 (take a look)
v+c+v	讨论并通过 (discuss and approve)
v+c+i	尝试但一无所获 (try but find nothing)
a+v	热烈鼓掌 (applause warmly)

Table 12: POS templates we use in event extraction method POS.