

# PHALM: Building a Knowledge Graph from Scratch by Prompting Humans and a Language Model

Anonymous ACL submission

## Abstract

Despite the remarkable progress in natural language understanding with pretrained Transformers, neural language models often do not have commonsense knowledge. Toward commonsense-aware models, there have been attempts to obtain knowledge, ranging from automatic acquisition to crowdsourcing. However, it is difficult to obtain a high-quality knowledge base at a low cost, especially from scratch. In this paper, we propose PHALM, a method of building a knowledge graph from scratch, by prompting both crowdworkers and a large language model. We used this method to build a Japanese event knowledge graph and trained Japanese neural commonsense models. Experimental results revealed the acceptability of the built graph and inferences generated by the trained models. We also report the difference in prompting humans and a language model.

## 1 Introduction

Since pretrained models (Radford and Narasimhan, 2018; Devlin et al., 2019; Yang et al., 2019) based on Transformer (Vaswani et al., 2017) appeared, natural language understanding has made remarkable progress. In some benchmarks, the performance of natural language understanding models has already exceeded that of humans. These models are applied to various downstream tasks ranging from translation and question answering to narrative understanding and dialogue response generation. In recent years, the number of parameters in such models has continued to increase (Radford et al., 2019; Brown et al., 2020), and so has their performance.

When we understand or reason, we usually rely on commonsense knowledge. Computers also need such knowledge to answer open-domain questions and to understand narratives and dialogues, for example. However, pretrained models often do not have commonsense knowledge.

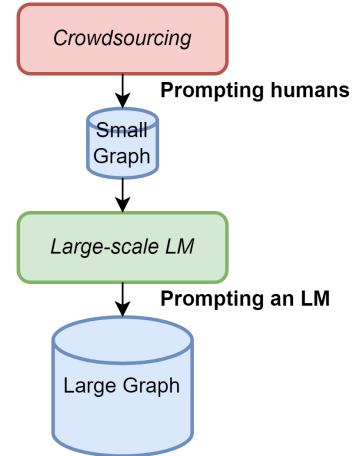


Figure 1: An overview of our method. We build a knowledge graph step by step from scratch, by prompting both humans and a language model.

There are many knowledge bases for commonsense inference. Some are built by crowdsourcing (Speer et al., 2017; Sap et al., 2019; Hwang et al., 2021), but acquiring a large-scale knowledge base is high-cost. Others are built by automatic acquisition (Zhang et al., 2019, 2020), but it is difficult to acquire high-quality commonsense knowledge. Recently, there have been some methods using large language models (LLMs) for building knowledge bases (Yuan et al., 2021; West et al., 2022; Liu et al., 2022). They often extend existing datasets, but do not build new datasets from scratch.

In this paper, we propose PHALM<sup>1</sup>, a method to build a knowledge graph from scratch with both crowdsourcing and an LLM. Asking humans to describe knowledge using crowdsourcing and generating knowledge using a language model are essentially the same (as it were, the latter is an analogy of the former), and both can be considered to be *prompting*. Therefore, we consider prompting for

<sup>1</sup>PHALM stands for **P**rompting **H**umans **A**nd a **L**anguage **M**odel.

062 both humans and a language model and gradually  
063 acquire a knowledge graph from a small scale to a  
064 large scale. Specifically, we acquire a small-scale  
065 knowledge graph by asking crowdworkers to de-  
066 scribe knowledge and use them as a few shots for  
067 an LLM to generate a large-scale knowledge graph.  
068 At each phase, we guarantee the quality of graphs  
069 by applying appropriate filtering.

070 We built a Japanese knowledge graph on events,  
071 considering prompts for both humans and a lan-  
072 guage model. With Yahoo! Crowdsourcing<sup>2</sup> and  
073 HyperCLOVA JP, a Japanese variant of the LLMs  
074 built by Kim et al. (2021), we obtained a knowl-  
075 edge graph that is not a simple translation, but  
076 unique to the culture. Then, we compared infer-  
077 ences collected by crowdsourcing and generated  
078 by the LLM. In addition to acquisition, we trained  
079 a Japanese neural commonsense model based on  
080 the built knowledge graph. With the model, we  
081 verified the acceptability of output inferences for  
082 unseen events. The resulting knowledge graph and  
083 the commonsense model created in this paper will  
084 be released to the public.<sup>3</sup>

## 085 2 Related Work

### 086 2.1 Commonsense Knowledge Datasets

087 There are several knowledge bases about common-  
088 sense, from what appears in the text to what is tacit  
089 but not written in the text. ConceptNet (Speer et al.,  
090 2017), for example, is a knowledge graph that con-  
091 nects words and phrases by relations. GenericsKB  
092 (Bhakthavatsalam et al., 2020) is a corpus describ-  
093 ing knowledge of entities in natural language rather  
094 than in graph.

095 In some datasets, commonsense knowledge is  
096 collected in the form of question answering. Roem-  
097 mele et al. (2011) acquire plausible causes and ef-  
098 fects for premises as two-choice questions. Zellers  
099 et al. (2018) provide SWAG, acquiring inferences  
100 about a situation from video captions as four-choice  
101 questions. KUCI (Omura et al., 2020) is a dataset  
102 for commonsense inference in Japanese, which is  
103 obtained by combining automatic extraction and  
104 crowdsourcing. Talmor et al. (2019) build Com-  
105 monsenseQA, which treats commonsense on Con-  
106 ceptNet’s entities as question answering.

<sup>2</sup><https://crowdsourcing.yahoo.co.jp/>

<sup>3</sup>if accepted

### 107 2.2 Knowledge Graphs on Events

108 Regarding commonsense knowledge bases, there  
109 are several graphs that focus on events. ATOMIC  
110 (Sap et al., 2019) describes the relationship be-  
111 tween events, mental states (Rashkin et al., 2018),  
112 and personas. Hwang et al. (2021) merge ATOMIC  
113 and ConceptNet, proposing ATOMIC-2020.

114 There are also studies for leveraging con-  
115 text. GLUCOSE (Mostafazadeh et al., 2020)  
116 is a commonsense inference knowledge graph  
117 for short stories, built by annotating ROCStories  
118 (Mostafazadeh et al., 2016). CIDER (Ghosal et al.,  
119 2021) and CICERO (Ghosal et al., 2022) are the  
120 graphs for dialogues, where DailyDialog (Li et al.,  
121 2017) and other dialogue corpora are annotated  
122 with inferences.

123 ASER (Zhang et al., 2019) is an event knowledge  
124 graph, automatically extracted from text corpora by  
125 focusing on discourse. With ASER, TransOMCS  
126 (Zhang et al., 2020) aims at bootstrapped knowl-  
127 edge graph acquisition by pattern matching and  
128 ranking.

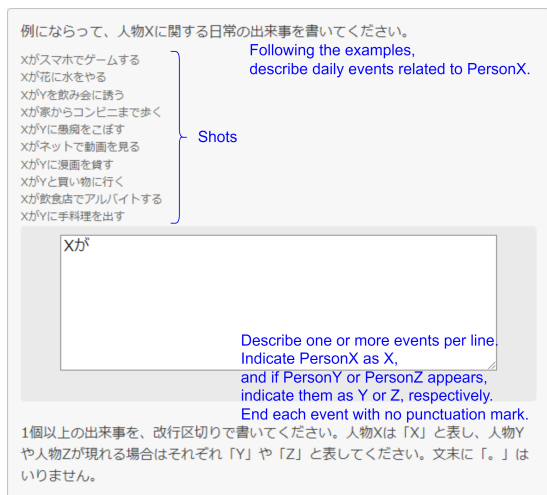
129 While ConceptNet and ATOMIC are acquired  
130 by crowdsourcing, ASER and TransOMCS are au-  
131 tomatically built. On one hand, a large-scale graph  
132 can be built easily in an automatic way, but it is  
133 difficult to obtain knowledge not appearing in the  
134 text. On the other hand, crowdsourcing can gather  
135 high-quality data, but it is expensive in terms of  
136 both money and time.

137 There is a method that uses crowdsourcing and  
138 neural language models together to build an event  
139 knowledge graph (West et al., 2022). Although it  
140 is possible to acquire a large-scale and high-quality  
141 graph, they assume that an initial graph, ATOMIC  
142 in this case, has already been available.

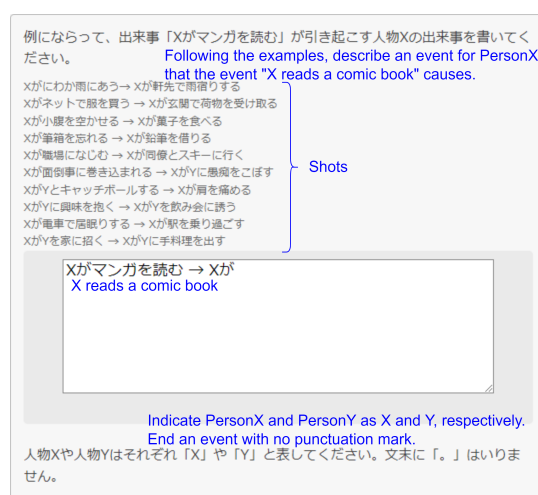
### 143 2.3 Neural Commonsense Models

144 There have been studies on storing knowledge in  
145 a neural form rather than a symbolic form. In par-  
146 ticular, methods of considering neural language  
147 models as knowledge bases (Petroni et al., 2019;  
148 Alkhamissi et al., 2022) have been developed.  
149 Bosselut et al. (2019) train COMET by finetuning  
150 pretrained Transformers on ATOMIC and Concept-  
151 Net, aiming at inference on unseen events and con-  
152 cepts. Gabriel et al. (2021) point out that COMET  
153 ignores discourse, introducing recurrent memory  
154 for paragraph-level information.

155 West et al. (2022) propose symbolic knowledge  
156 distillation where specific knowledge in a general



(a) For events



(b) For inferences (xEffect)

Figure 2: Examples of crowdsourcing interfaces. Crowdworkers are asked to describe events and inferences.

language model is distilled into a specific language model via a symbolic form. They expand ATOMIC using GPT-3 (Brown et al., 2020), filter the outputs using RoBERTa (Liu et al., 2019), and finetune GPT-2 (Radford et al., 2019) on the filtered ones.

### 3 Prompting Humans and a Language Model

We propose a method to build a knowledge graph for commonsense inference from scratch, with both crowdsourcing and a language model. In our method, we first construct a small-scale knowledge graph by crowdsourcing. Using the small-scale graph for prompts, we then extract commonsense knowledge from a language model. The flow of our method is shown in Figure 1. Building a knowledge graph from scratch only by crowdsourcing is expensive in terms of both money and time. Hence, the combination of crowdsourcing and a language model is expected to reduce the cost, especially in terms of time.

In other words, our method consists of the following two phases: (1) collecting a small-scale graph by crowdsourcing and (2) generating a large-scale graph by a language model. While crowdsourcing elicits commonsense from people, shots are used to extract knowledge from a language model. At this point, these phases are intrinsically the same, being considered as *prompting*. In the two phases, namely, we prompt people and a language model, respectively.

We build a commonsense inference knowledge graph in Japanese, with the concept of Section 3.

We focus on an event knowledge graph such as ATOMIC (Sap et al., 2019) and ASER (Zhang et al., 2019). Handling commonsense on events and mental states would facilitate understanding of narratives and dialogues. We use Yahoo! Crowdsourcing in the first phase and HyperCLOVA JP (Kim et al., 2021), an LLM in Japanese, in the second phase.

#### 3.1 Acquisition by Crowdsourcing

We first acquire a small-scale high-quality knowledge graph by crowdsourcing. With Yahoo! Crowdsourcing, specifically, we ask crowdworkers to write events and inferences. In a task, we provide them with 10 shots as a prompt for each event and inference. Note that for inferences, the prompts differ for each relation as mentioned later. We obtain a graph by filtering the collected inferences syntactically and semantically.

**Events** We ask crowdworkers to write daily events related to at least one person (PersonX). An example of the crowdsourcing task interface is shown in Figure 2a. The task provides instructions and 10 examples, and each crowdworker is asked to write at least one event. After all tasks are completed, we remove duplicate events. As a result, 257 events were acquired from 200 crowdworkers. We manually verified that all of the acquired events have a sufficient quality.

**Inferences** For the events collected above, we ask crowdworkers to write inferences about what happens and how a person feels before and after the events. In this paper, the relations for inference

	Inst #	Val #	Val %	IAA
Event	257	-	-	-
xNeed	504	402	79.76	39.85
xEffect	621	554	89.21	25.00
xIntent	603	519	86.07	36.11
xReact	639	550	86.07	31.82

Table 1: The statistics on events and inferences acquired by crowdsourcing.

are based on ATOMIC.<sup>4</sup> The following four are adopted as our target relations.

- What would have happened before (xNeed)
- What would happen after (xEffect)
- What PersonX would have felt before (xIntent)
- What PersonX would feel after (xReact)

While xNeed and xEffect are inferences about events, xIntent and xReact are inferences about mental states.

Three crowdworkers are hired per event. Given an instruction and 10 examples, each crowdworker is asked to write one inference. An example of the crowdsourcing task interface is shown in Figure 2b. We remove duplicate inferences as in the case of events, and then apply syntactic filtering<sup>5</sup> using the Japanese syntactic parser KNP<sup>6</sup>.

The statistics of the acquired events and inferences are shown in the Inst # column of Table 1. The whole process costed 16,844 JPY (approximately 123 USD) by hiring 547 crowdworkers. Examples of acquired inferences are shown in Table 2.

### 3.2 Evaluation and Filtering

To examine the qualities of the inferences acquired by crowdsourcing, we crowdsource their evaluation. We ask three crowdworkers whether the inferences are acceptable or not and judge their acceptability by majority voting. The evaluation is

<sup>4</sup>The relations are not exactly the same as those of ATOMIC. xIntent in this paper covers xIntent and xWant in ATOMIC, and tails for our xIntent and xReact may contain not mental states but events. The reason for the difference is that English and Japanese have different linguistic characteristics, i.e., it is difficult to collect knowledge in the same structure as the original.

<sup>5</sup>KNP determines if the subject is PersonX, if the tense is present, and if the event is a single sentence.

<sup>6</sup><https://nlp.ist.i.kyoto-u.ac.jp/?KNP>

1. Xがスマホでゲームする (X plays a game on X's phone)
2. Xが花に水をやる (X waters flowers)
3. XがYを飲み会に誘う (X invites Y to a drinking party)
- ...
11. XがYに謝る (X apologizes to Y)

(a) For events

1. Xがにわか雨にあう。結果として、Xが軒先で雨宿りする。  
(X gets caught in a shower. As a result, X takes shelter from the rain under the eaves.)
2. Xがネットで服を買う。結果として、Xが荷物を受け取る。  
(X buys clothes on the Internet. As a result, X receives a package.)
3. Xが小腹を空かせる。結果として、Xが菓子を食べる。  
(X gets hungry. As a result, X eats a snack.)
- ...
11. Xが筆箱を忘れる。結果として、Xが鉛筆を借りる。  
(X forgets to bring X's pencil case. As a result, X borrows a pencil.)

(b) For inferences (xEffect)

Figure 3: Prompts for generating events and inferences from an LLM. The underlined parts are generated.

crowdsourced independently for each relation. The inferences judged to be unacceptable by majority voting are filtered out.

The inferences collected in Section 3.1 are evaluated and filtered as above. The statistics are listed in the middle two columns of Table 1. As a result, we employed 465 crowdworkers and spent 8,679 JPY (approximately 63 USD). We also calculated Fleiss's  $\kappa$  as an inner-annotator agreement in the evaluation, which is shown in the rightmost column of Table 1.

There are several tendencies in the inferences filtered out, i.e., judged to be unacceptable. In some inferences, the order is reversed, as in the triple ⟨PersonX sleeps twice, xEffect, PersonX thinks that they are off work today⟩. Others are not plausible, as in ⟨PersonX surfs the Internet, xNeed, PersonX gets to the ocean⟩.

### 3.3 Generation from an LLM

From a small-scale high-quality knowledge graph acquired in Sections 3.1 and 3.2, we generate a large-scale knowledge graph with an LLM. We use the Koya 39B model of HyperCLOVA JP as a language model. Both events and inferences are generated by providing 10 shots. The shots are randomly chosen from the small-scale graph for each generation.

**Events** New events are generated by HyperCLOVA JP, using the events acquired in Section 3.2 as shots. An example prompt for event generation



Head	Rel	Tail	Eval
Xが顔を洗う (X washes X's face)	xNeed	Xが水道で水を出す (X runs water from the tap)	✓
	xEffect	Xが歯を磨く (X brushes X's teeth) Xがタオルを準備する (X prepares a towel) Xが鏡に映った自分の顔に覚えのない傷を見つける (X finds an unrecognized scar on X's face in the mirror)	✓ ✓ ✓
	xIntent	Xが歯磨きをする (X brushes his teeth) スッキリしたい (Want to feel refreshed) 眠いのでしゃまっとしたい (Sleepy and Want to feel refreshed)	✓ ✓ ✓
	xReact	さっぱりして眠気覚ましになる (Feel refreshed and shake off X's sleepiness) きれいになる (Be clean) さっぱりした (Felt refreshed)	✓ ✓ ✓

Table 2: Examples of inferences acquired through crowdsourcing. Triples with ✓ in the eval column were judged to be acceptable by the evaluation in Section 3.2.

Rel	Template
xNeed	<i>h</i> ためには、 <i>t</i> 必要がある。(To <i>h</i> , need to <i>t</i> .)
xEffect	<i>h</i> . 結果として、 <i>t</i> . ( <i>h</i> . As a result, <i>t</i> .)
xIntent	<i>h</i> のは、 <i>t</i> と思ったから。( <i>h</i> because felt <i>t</i> .)
xReact	<i>h</i> と、 <i>t</i> と思う。( <i>h</i> then feel <i>t</i> .)

Table 3: The templates of shots for an LLM. *h* and *t* stand for head and tail, respectively. When generating, *t* is extracted.

	Inst #	Val %	IAA
Event	1,471	-	-
xNeed	9,403	80.81	36.07
xEffect	8,792	85.45	34.03
xIntent	10,155	86.06	43.42
xReact	10,941	90.30	21.51

Table 4: The statistics of events and inferences generated from an LLM. % Val and IAA are the evaluation results of 500 randomly selected inferences.

is shown in Figure 3a. We generate 10,000 events, remove duplicates, and apply the same syntactic filtering as in Section 3.1.

**Inferences** As in event generation, the inferences acquired in Sections 3.1 and 3.2 are used as shots. We generate 10 inferences for each event and remove duplicate triples. While we simply list the shots as a prompt in event generation, different prompts are used for each relation in inference generation. An example prompt for xEffect generation is shown in Figure 3b. Shots are given in natural language, and tails are extracted by pattern matching. Shot templates for each relation are shown in Table 3. Finally, the syntactic filtering is applied to obtain the graph.

The statistics of events and inferences generated by HyperCLOVA JP are shown in Table 4, and the results of the evaluation and the inter-annotator agreement are also shown in Table 4. For this evaluation, we sampled 500 inferences per relation. We hired 409 crowdworkers for a fee of 7,260 JPY (approximately 53 USD) in total. A comparison with Table 1 indicates that the quality is as good as those written by crowdworkers. Examples of generated inferences are shown in Table 5.

The generated knowledge graph in Japanese reflects the culture of Japan, such as (PersonX goes to the office, xNeed, *PersonX takes a train*). This fact indicates the importance of building from scratch for a specific language, rather than translating a similar dataset in a different language, which emphasizes the value of our method proposed in this paper.

## 4 Analysis on the Built Knowledge Graph

### 4.1 Effect of Filtering

In this paper, a small-scale knowledge graph is collected as in Sections 3.1 and 3.2, and a large-scale knowledge graph is generated as in Section 3.3. Here, we examine how effective the filtering in Section 3.2 is. As an experiment, we use filtered and unfiltered small-scale graphs as prompts to generate a large-scale graph. Then, we randomly select 500 generated triples for each relation and evaluate them by crowdsourcing as in Section 3.2. Note that the results for the filtered triples are the same as Section 3.3. For the triples without filtering, we crowdsourced again, paying 393 crowdworkers 7,260 JPY (approximately 53 USD).

The ratios of appropriate inferences with and without filtering are shown in Table 6. For all rela-

Head	Rel	Tail
Xがコンビニへ行く (X goes to a convenience store)	xNeed	Xが財布を持っている (X has X's wallet), Xが外出する (X goes out), Xが外出着に着替える (X changes into going-out clothes), Xが財布を持って出かける (X goes out with X's wallet), Xが外へ出る (X goes outside)
	xEffect	Xが買い物をする (X goes shopping), Xが雑誌を立ち読みする (X browses through magazines), XがATMでお金をおろす (X withdraws money from ATM), Xが弁当を買う (X buys lunch), Xがアイスを買う (X buys ice cream)
	xIntent	何か買いたいものがある (Want to buy something), 雑誌を買う (Buy a magazine), 飲み物を買おう (Going to buy a drink), 飲み物や食べ物を買いたい (Want to buy a drink or food), なんでもある (There is everything X wants)
	xReact	何か買いたいものがある (Want to buy something), 何か買う (Buy something), 何か買おう (Going to buy something), 何か買いたくなる (Come to buy something), ついでに何か買ってしまう (Buy something incidentally)

Table 5: Examples of inferences generated from an LLM. For each relation, five examples are displayed.

	xNeed	xEffect	xIntent	xReact
w/o Fltr	<b>81.62</b>	82.42	83.84	89.29
w/ Fltr	80.81	<b>85.45</b>	<b>86.06</b>	<b>90.30</b>

Table 6: The ratios of appropriate inferences with respect to filtering. Note that the w/ Fltr row is the same as the Val % column in Table 4.

tions except xNeed, filtering improves the quality of triples.

## 4.2 Comparison between humans and a Language Model

In Section 3.1, on one hand, we asked crowdworkers to describe events and inferences. In Section 3.3, on the other hand, we had an LLM generate them. Here, we compare a small-scale knowledge graph by crowdsourcing and a large-scale one from a language model, i.e., inferences generated by humans and a computer. Because the relationships between events can be largely divided into contingent and temporal relationships (Bethard et al., 2008), we adopt contingency and time interval as metrics for comparison.

Of the four relations, we focus on xEffect as a representative, which is a typical causal relation. For each head of the triples acquired by crowdsourcing in Sections 3.1 and 3.2, we generate three tails using the language model in Section 3.3 and compare them with the original tails. From the 554 heads for xEffect in the small-scale graph, we obtained 586 unique inferences.

**Contingency** One measure is how likely a given event is to be followed by a subsequent event. Crowdworkers are given a pair of events in an xEffect relation and asked to judge how likely the following event is to happen on a three-point scale: “must happen,” “likely to happen,” and “does not

happen.” We ask three crowdworkers per inference and calculate the median of them.

**Time Interval** The other measure is the time interval between the occurrence of an event and that of a subsequent event. As in the evaluation of contingency, crowdworkers are given a triple on xEffect. We ask them to judge the time interval between the two events in five levels: almost simultaneous, seconds to minutes, hours, days to months, and longer. Finally, the median is calculated from the results of three crowdworkers.

The comparison between humans and a language model for each measure is shown in Figure 4. Figure 4a shows that the subsequent events by crowdsourcing, or humans, are slightly more probable. In Figure 4b, the inferences generated by an LLM have a longer time interval. This result indicates a difference in the results of prompting humans and a language model; for xEffect, humans infer events that happen relatively soon, while a language model infers events that happen a bit later.

## 5 Japanese Neural Commonsense Models

We train Japanese neural commonsense models using the knowledge graph constructed in Section 4. Japanese versions of GPT-2 (Radford et al., 2019) and T5 (Raffel et al., 2020) are finetuned to generate inferences on unseen events. We conduct automatic and manual evaluations and compare their performances.

### 5.1 Training

**Base models and data** Using the constructed knowledge graph, we finetune pretrained models to construct Japanese neural commonsense models. To evaluate inferences on unseen events, triples in

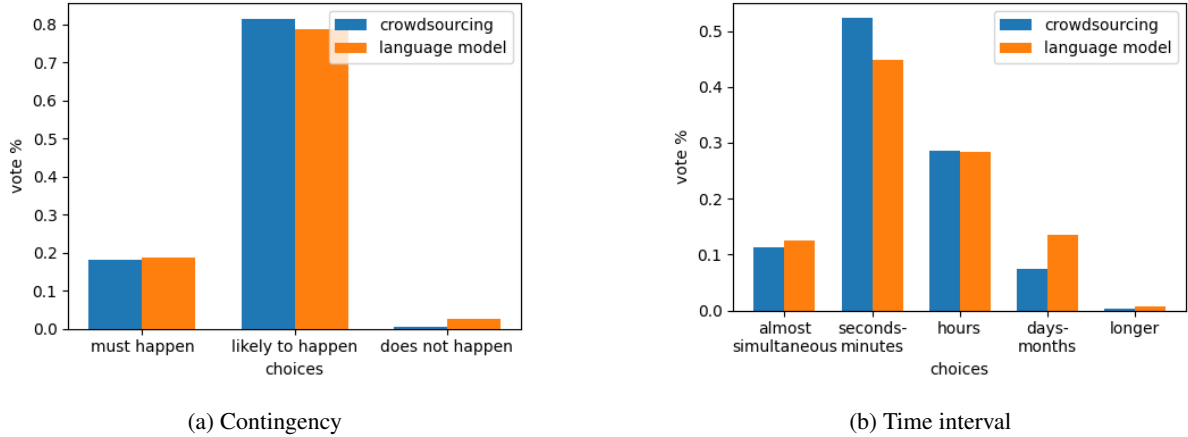


Figure 4: A comparison between crowdsourcing and language model generation.

the knowledge graph are randomly partitioned into training and test sets at a ratio of 9:1. For pretrained models, we adopt Japanese T5<sup>7</sup> and GPT-2<sup>8</sup> of the Hugging Face implementation (Wolf et al., 2020).

**Input format to models** The input for each model differs. See Appendix C for the full input formats for each model. Since T5 is a seq2seq model, the head and the relation are given in the form of “ $r : h$ ” as an input, and the tail is given as the correct output. The relation for T5 is changed to a natural language sentence. For example, “xNeed” is rewritten to “What event occurs before this statement?” The inputs for all relations are shown in Appendix C. For GPT-2, since it predicts the next word, the head and the relation are given as an input, and the model is trained to output the tail. Since the relations are not included in the vocabulary of the pretrained models, they are added as special tokens. In the constructed knowledge graph, the subject of an event is generalized as “X,” but it would be better to change it into a natural expression as the input to the pretrained models. We randomly replace the subject with a personal pronoun during training and inference. To confirm this effect, in section 5.2, we also train GPT-2 with the subject represented as “X.” We denote this as GPT-2<sub>X</sub>.

## 5.2 Evaluation

We generate inferences for the head events in the test set using the trained Japanese neural common-sense models and evaluate the inferences automatically and manually. We also show correlation

<sup>7</sup><https://huggingface.co/megagonlabs/t5-base-japanese-web>

<sup>8</sup><https://huggingface.co/nlp-waseda/gpt2-small-japanese>

Model	AR	MP	BS	BLEU
T5	87.5	1.64	90.26	18.57
GPT-2	<b>91.0</b>	<b>1.73</b>	<b>92.31</b>	18.26
GPT-2 <sub>X</sub>	<b>91.0</b>	1.68	92.03	<b>18.99</b>

Table 7: Total evaluation scores. AR, MP, and BS indicate the accept rate, the mean point, and BERTScore, respectively.

Rel	AR	MP	BS	BLEU
xNeed	88.9	1.58	92.73	22.22
xEffect	92.4	1.72	93.98	22.24
xIntent	88.9	1.66	90.12	9.91
xReact	93.8	1.98	93.00	11.83

Table 8: Evaluation scores of GPT-2 for each relation.

between the automatic and manual evaluations. Examples of the inference results are shown in Appendix C. The average output length and the number of unique words are also reported in Appendix C. In summary, the number of unique words in GPT-2 is larger than that in T5 (392 unique words), with a difference of 35 to 59 words.

**Automatic evaluation** We calculate BLEU (Papineni et al., 2002) and BERTScore (Zhang\* et al., 2020) as automatic metrics. Table 7 shows these results. GPT-2<sub>X</sub> and GPT-2 performed the best in BLEU and BERTScore, respectively.

**Manual evaluation** Using crowdsourcing, we evaluate how likely the generated inferences are. Following the previous study (West et al., 2022), we show crowdworkers two events (a head and a tail) and a relation. Then, we ask them to evaluate the appropriateness of the inference by choos-

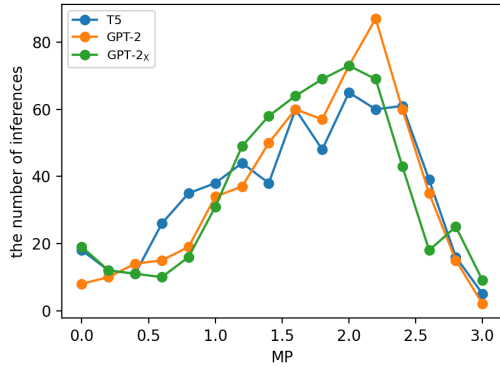


Figure 5: The number of inferences for each MP.

ing from the following options: “always,” “often,” “sometimes,” and “never.” The choices are displayed with an appropriate verb for each relation (e.g., “always happens” for xEffect). Five crowdworkers are asked to judge per inference. For each inference, the numbers of crowdworkers who choose “never” and other than “never” (i.e., at least “sometimes”) are used to determine the majority vote. The acceptance rate (AR), the proportion of inferences in which more crowdworkers choose other than “never.” By assigning 0 to 3 points each to “never,” “sometimes,” “often,” and “always,” we also calculate the mean point (MP) as the average score of all the inferences. Table 7 shows these results. AR is higher than 85% for all models, indicating that the inferences for unseen events are almost correct. GPT-2 obtained the highest scores for both AR and MP. Furthermore, as shown in Table 8, ARs of xNeed and xIntent are lower than xEffect and xReact, respectively, for all models. This can be attributed to the fact that we used an autoregressive model, which makes it difficult to infer in reverse order of time.

Although the replacement of subjects did not make a difference in AR, there is a difference in the distributions of MP as shown in Figure 5. The number of crowdworkers who chose “never” for the inference of GPT-2 is less than half of that for GPT-2<sub>X</sub>. This result indicates that it is better for the model to replace subjects “X” with personal pronouns.

**Correlation between the evaluation metrics** Table 9 shows the correlation coefficients between the manual and automatic evaluation metrics. The correlation coefficients between the manual metrics (AR and MP) and BERTScore are positive, while those between the manual metrics and BLEU

	AR	MP	BS	BLEU
AR	1.00	0.75	0.59	-0.11
MP	-	1.00	0.43	-0.46
BS	-	-	1.00	0.30
BLEU	-	-	-	1.00

Table 9: Correlation coefficients between automatic and manual evaluation metrics.

are negative or no correlation. It seems that BERTScore, which uses vector representations, can evaluate equivalent sentences with different expressions, but BLEU, which is based on n-gram agreement, cannot correctly judge the equivalence. One of the reasons for the negative correlation in BLEU is that many inferences of the mental state consist of a single word in Japanese, such as “tired” and “bored,” for both the gold answer and the generated result. In this case, BLEU tends to be low because the words are rarely matched, but the shorter the sentences are, the easier it is for the model to generate appropriate results.

## 6 Conclusion

We proposed a method for building a knowledge graph from scratch with both crowdsourcing and a language model. Based on our method, we built a knowledge graph on events and mental states in Japanese using Yahoo! Crowdsourcing and HyperCLOVA JP. Since designing tasks for having humans describe commonsense and engineering prompts for having a language model generate knowledge are similar to each other, we compared the characteristics of them. We evaluated the graph generated by HyperCLOVA JP and found that it was similar in quality to the graph written by humans.

Furthermore, we trained a neural commonsense model for event inference based on the built knowledge graph. We attempted inference generation for unseen events by finetuning GPT-2 and T5 in Japanese on the built graph. The experimental results showed that these models are able to generate acceptable inferences for events and mental states.

We hope that our method for building a knowledge graph from scratch and the acquired knowledge graph lead to further studies on commonsense inference, especially in low-resource languages.

## Ethical Considerations

For acquiring a small-scale event knowledge graph and analyzing the built graph, we crowdsource com-



520 nonsense knowledge, using Yahoo! Crowdsourc-  
 521 ing. Specifically, we collect the descriptions of  
 522 commonsense, filter them, and explore the charac-  
 523 teristics of the graph by crowdsourcing. Fees and  
 524 the numbers of crowdworkers per process are in  
 525 the text. In total, we employed 1,814 crowdwoek-  
 526 ers paying 40,043 JPY (approximately 288 USD).  
 527 We obtained a consent from crowdworkers on the  
 528 platform of Yahoo! Crowdsourcing.

529 The event knowledge graph and the neural com-  
 530 monsense models built in this paper help computers  
 531 understand commonsense. A commonsense-aware  
 532 computer, for example, can answer open-domain  
 533 questions by humans, interpret human statements  
 534 in detail, and converse with humans naturally. How-  
 535 ever, such graphs and models may contain incorrect  
 536 knowledge even with filtering, which leads the ap-  
 537 plications to harmful behavior.

## 538 References

539 Badr AlKhamissi, Millicent Li, Asli Celikyilmaz, Mona  
 540 Diab, and Marjan Ghazvininejad. 2022. [A review on  
 541 language models as knowledge bases.](#)

542 Steven Bethard, William Corvey, Sara Klingsenstein,  
 543 and James H. Martin. 2008. [Building a corpus of  
 544 temporal-causal structure.](#) In *Proceedings of the  
 545 Sixth International Conference on Language Re-  
 546 sources and Evaluation (LREC’08)*, Marrakech, Mo-  
 547 rocco. European Language Resources Association  
 548 (ELRA).

549 Sumithra Bhakthavatsalam, Chloe Anastasiades, and  
 550 Peter Clark. 2020. [Generickb: A knowledge base of  
 551 generic statements.](#)

552 Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chai-  
 553 tanya Malaviya, Asli Celikyilmaz, and Yejin Choi.  
 554 2019. [COMET: Commonsense transformers for auto-  
 555 matic knowledge graph construction.](#) In *Proceedings  
 556 of the 57th Annual Meeting of the Association for  
 557 Computational Linguistics*, pages 4762–4779, Flo-  
 558 rence, Italy. Association for Computational Linguis-  
 559 tics.

560 Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
 561 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
 562 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
 563 Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
 564 Gretchen Krueger, Tom Henighan, Rewon Child,  
 565 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens  
 566 Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-  
 567 teusz Litwin, Scott Gray, Benjamin Chess, Jack  
 568 Clark, Christopher Berner, Sam McCandlish, Alec  
 569 Radford, Ilya Sutskever, and Dario Amodei. 2020.  
 570 [Language models are few-shot learners.](#) In *Ad-  
 571 vances in Neural Information Processing Systems*,  
 572 volume 33, pages 1877–1901. Curran Associates,  
 573 Inc.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
 Kristina Toutanova. 2019. [BERT: Pre-training of  
 deep bidirectional transformers for language under-  
 standing.](#) In *Proceedings of the 2019 Conference of  
 the North American Chapter of the Association for  
 Computational Linguistics: Human Language Tech-  
 nologies, Volume 1 (Long and Short Papers)*, pages  
 4171–4186, Minneapolis, Minnesota. Association for  
 Computational Linguistics.

Saadia Gabriel, Chandra Bhagavatula, Vered Shwartz,  
 Ronan Le Bras, Maxwell Forbes, and Yejin Choi.  
 2021. [Paragraph-level commonsense transformers  
 with recurrent memory.](#) *Proceedings of the AAAI  
 Conference on Artificial Intelligence*, 35(14):12857–  
 12865.

Deepanway Ghosal, Pengfei Hong, Siqi Shen, Navonil  
 Majumder, Rada Mihalcea, and Soujanya Poria. 2021.  
[CIDER: Commonsense inference for dialogue expla-  
 nation and reasoning.](#) In *Proceedings of the 22nd  
 Annual Meeting of the Special Interest Group on Dis-  
 course and Dialogue*, pages 301–313, Singapore and  
 Online. Association for Computational Linguistics.

Deepanway Ghosal, Siqi Shen, Navonil Majumder,  
 Rada Mihalcea, and Soujanya Poria. 2022. [CICERO:  
 A dataset for contextualized commonsense inference  
 in dialogues.](#) In *Proceedings of the 60th Annual Meet-  
 ing of the Association for Computational Linguistics  
 (Volume 1: Long Papers)*, pages 5010–5028, Dublin,  
 Ireland. Association for Computational Linguistics.

Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras,  
 Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and  
 Yejin Choi. 2021. [\(comet-\) atomic 2020: On sym-  
 bolic and neural commonsense knowledge graphs.](#)  
*Proceedings of the AAAI Conference on Artificial  
 Intelligence*, 35(7):6384–6392.

Boseop Kim, HyungSeok Kim, Sang-Woo Lee,  
 Gichang Lee, Donghyun Kwak, Jeon Dong Hyeon,  
 Sunghyun Park, Sungju Kim, Seonhoon Kim, Dong-  
 pil Seo, Heungsub Lee, Minyoung Jeong, Sungjae  
 Lee, Minsub Kim, Suk Hyun Ko, Seokhun Kim,  
 Taeyong Park, Jinuk Kim, Soyoung Kang, Na-Hyeon  
 Ryu, Kang Min Yoo, Minsuk Chang, Soobin Suh,  
 Sookyo In, Jinseong Park, Kyungduk Kim, Hiun  
 Kim, Jisu Jeong, Yong Goo Yeo, Donghoon Ham,  
 Dongju Park, Min Young Lee, Jaewook Kang, Inho  
 Kang, Jung-Woo Ha, Woomyoung Park, and Nako  
 Sung. 2021. [What changes can large-scale language  
 models bring? intensive study on HyperCLOVA:  
 Billions-scale Korean generative pretrained trans-  
 formers.](#) In *Proceedings of the 2021 Conference  
 on Empirical Methods in Natural Language Process-  
 ing*, pages 3405–3424, Online and Punta Cana, Do-  
 minican Republic. Association for Computational  
 Linguistics.

Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang  
 Cao, and Shuzi Niu. 2017. [DailyDialog: A manually  
 labelled multi-turn dialogue dataset.](#) In *Proceedings  
 of the Eighth International Joint Conference on Nat-  
 ural Language Processing (Volume 1: Long Papers)*,

633	pages 986–995, Taipei, Taiwan. Asian Federation of	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	690
634	Natural Language Processing.	Lee, Sharan Narang, Michael Matena, Yanqi	691
635	Alisa Liu, Swabha Swayamdipta, Noah A. Smith, and	Zhou, Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the</a>	692
636	Yejin Choi. 2022. <a href="#">Wanli: Worker and ai collaboration</a>	<a href="#">limits of transfer learning with a unified text-to-text</a>	693
637	<a href="#">for natural language inference dataset creation.</a>	<a href="#">transformer.</a> <i>Journal of Machine Learning Research</i> ,	694
638	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	21(140):1–67.	695
639	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	Hannah Rashkin, Maarten Sap, Emily Allaway, Noah A.	696
640	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	Smith, and Yejin Choi. 2018. <a href="#">Event2Mind: Com-</a>	697
641	<a href="#">Roberta: A robustly optimized bert pretraining ap-</a>	<a href="#">monsense inference on events, intents, and reactions.</a>	698
642	<a href="#">proach.</a>	In <i>Proceedings of the 56th Annual Meeting of the</i>	699
643	Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong	<i>Association for Computational Linguistics (Volume 1:</i>	700
644	He, Devi Parikh, Dhruv Batra, Lucy Vanderwende,	<i>Long Papers)</i> , pages 463–473, Melbourne, Australia.	701
645	Pushmeet Kohli, and James Allen. 2016. <a href="#">A corpus</a>	Association for Computational Linguistics.	702
646	<a href="#">and cloze evaluation for deeper understanding of</a>	Melissa Roemmele, Cosmin Adrian Bejan, and Andrew	703
647	<a href="#">commonsense stories.</a> In <i>Proceedings of the 2016</i>	S. Gordon. 2011. Choice of plausible alternatives:	704
648	<i>Conference of the North American Chapter of the</i>	An evaluation of commonsense causal reasoning. In	705
649	<i>Association for Computational Linguistics: Human</i>	<i>AAAI Spring Symposium on Logical Formalizations</i>	706
650	<i>Language Technologies</i> , pages 839–849, San Diego,	<i>of Commonsense Reasoning</i> , Stanford University.	707
651	California. Association for Computational Linguistics.	Maarten Sap, Ronan Le Bras, Emily Allaway, Chan-	708
652		dra Bhagavatula, Nicholas Lourie, Hannah Rashkin,	709
653	Nasrin Mostafazadeh, Aditya Kalyanpur, Lori Moon,	Brendan Roof, Noah A. Smith, and Yejin Choi. 2019.	710
654	David Buchanan, Lauren Berkowitz, Or Biran, and	<a href="#">Atomic: An atlas of machine commonsense for if-</a>	711
655	Jennifer Chu-Carroll. 2020. <a href="#">GLUCOSE: Gener-</a>	<a href="#">then reasoning.</a> <i>Proceedings of the AAAI Conference</i>	712
656	<a href="#">aLized and COntextualized story explanations.</a> In	<i>on Artificial Intelligence</i> , 33(01):3027–3035.	713
657	<i>Proceedings of the 2020 Conference on Empirical</i>	Robyn Speer, Joshua Chin, and Catherine Havasi. 2017.	714
658	<i>Methods in Natural Language Processing (EMNLP)</i> ,	<a href="#">Conceptnet 5.5: An open multilingual graph of gen-</a>	715
659	pages 4569–4586, Online. Association for Computa-	<a href="#">eral knowledge.</a> <i>Proceedings of the AAAI Conference</i>	716
660	tional Linguistics.	<i>on Artificial Intelligence</i> , 31(1).	717
661	Kazumasa Omura, Daisuke Kawahara, and Sadao Kuro-	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and	718
662	hashi. 2020. <a href="#">A method for building a commonsense</a>	Jonathan Berant. 2019. <a href="#">CommonsenseQA: A ques-</a>	719
663	<a href="#">inference dataset based on basic events.</a> In <i>Proceed-</i>	<a href="#">tions answering challenge targeting commonsense</a>	720
664	<i>ings of the 2020 Conference on Empirical Methods</i>	<a href="#">knowledge.</a> In <i>Proceedings of the 2019 Conference</i>	721
665	<i>in Natural Language Processing (EMNLP)</i> , pages	<i>of the North American Chapter of the Association for</i>	722
666	2450–2460, Online. Association for Computational	<i>Computational Linguistics: Human Language Tech-</i>	723
667	Linguistics.	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	724
668	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	4149–4158, Minneapolis, Minnesota. Association for	725
669	Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evalu-</a>	Computational Linguistics.	726
670	<a href="#">ation of machine translation.</a> In <i>Proceedings of the</i>	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	727
671	<i>40th Annual Meeting of the Association for Computa-</i>	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	728
672	<i>tional Linguistics</i> , pages 311–318, Philadelphia,	Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all</a>	729
673	Pennsylvania, USA. Association for Computational	<a href="#">you need.</a> In <i>Advances in Neural Information Pro-</i>	730
674	Linguistics.	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.	731
675	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel,	Peter West, Chandra Bhagavatula, Jack Hessel, Jena	732
676	Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and	Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu,	733
677	Alexander Miller. 2019. <a href="#">Language models as knowl-</a>	Sean Welleck, and Yejin Choi. 2022. <a href="#">Symbolic</a>	734
678	<a href="#">edge bases?</a> In <i>Proceedings of the 2019 Confer-</i>	<a href="#">knowledge distillation: from general language mod-</a>	735
679	<i>ence on Empirical Methods in Natural Language Pro-</i>	<a href="#">els to commonsense models.</a> In <i>Proceedings of the</i>	736
680	<i>cessing and the 9th International Joint Conference</i>	<i>2022 Conference of the North American Chapter of</i>	737
681	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	<i>the Association for Computational Linguistics: Hu-</i>	738
682	pages 2463–2473, Hong Kong, China. Association	<i>man Language Technologies</i> , pages 4602–4625, Seat-	739
683	for Computational Linguistics.	tle, United States. Association for Computational	740
684	Alec Radford and Karthik Narasimhan. 2018. Im-	Linguistics.	741
685	proving language understanding by generative pre-	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	742
686	training.	Chaumond, Clement Delangue, Anthony Moi, Pier-	743
687	Alec Radford, Jeff Wu, Rewon Child, David Luan,	ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-	744
688	Dario Amodei, and Ilya Sutskever. 2019. Language	icz, Joe Davison, Sam Shleifer, Patrick von Platen,	745
689	models are unsupervised multitask learners.	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	746

Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. **Transformers: State-of-the-art natural language processing**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. **Xlnet: Generalized autoregressive pretraining for language understanding**. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.

Ann Yuan, Daphne Ippolito, Vitaly Nikolaev, Chris Callison-Burch, Andy Coenen, and Sebastian Gehrmann. 2021. **Synthbio: A case study in faster curation of text datasets**. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1.

Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. **SWAG: A large-scale adversarial dataset for grounded commonsense inference**. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 93–104, Brussels, Belgium. Association for Computational Linguistics.

Hongming Zhang, Daniel Khashabi, Yangqiu Song, and Dan Roth. 2020. **Transomcs: From linguistic graphs to commonsense knowledge**. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pages 4004–4010. International Joint Conferences on Artificial Intelligence Organization. Main track.

Hongming Zhang, Xin Liu, Haojie Pan, Yangqiu Song, and Cane Wing-Ki Leung. 2019. **Aser: A large-scale eventuality knowledge graph**.

Tianyi Zhang\*, Varsha Kishore\*, Felix Wu\*, Kilian Q. Weinberger, and Yoav Artzi. 2020. **Bertscore: Evaluating text generation with bert**. In *International Conference on Learning Representations*.

## A An Example of Crowdsourced Evaluation

We evaluate and filter the inferences obtained in Sections 3.1 and 3.3 by crowdsourcing. An example of the interface for evaluating an xEffect inference is shown in Figure 6.

## B Hyperparameter Details

We generate a large-scale knowledge graph using HyperCLOVA JP in Section 3.3. The hyperparameters for the generation is shown in Table 10.

With the built knowledge graph, we finetune Japanese T5 and GPT-2 on the task of commonsense inference in Section 5. The hyperparameters for T5 and GPT-2 are shown in Table 11.

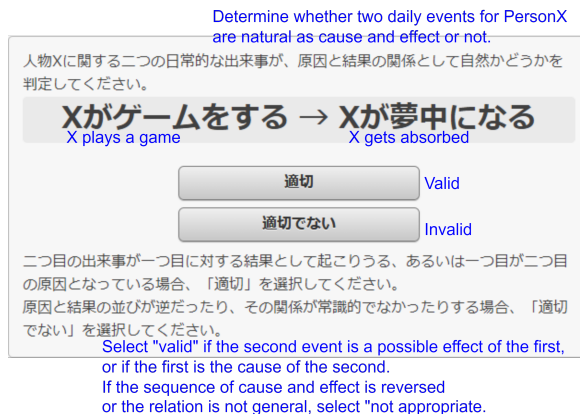


Figure 6: An example of evaluation regarding xEffect relations. We ask three crowdworkers whether a given inference is acceptable or not.

Max tokens	32
Temperature	0.5
Top-P	0.8
Top-K	0
Repeat penalty	5.0

Table 10: Hyperparameters for event and inference generation with HyperCLOVA JP.

## C Details of Neural Commonsense Models

Table 12 shows the average output length and the number of unique words for each model. The average output length of T5 is longer than those of GPT-2s, but GPT-2s have the greater numbers of unique words than T5.

Table 13 shows the input formats to the models. The prompts to T5 may not be the best; prompt-engineering could improve the results.

Examples of outputs are shown in Table 14. We can see that the obtained outputs are acceptable to humans. The outputs vary for each model.

	<b>T5</b>	<b>GPT-2</b>
Batch size	64	64
Learning rate	5e-5	5e-5
Weight decay	0.0	0.0
Adam betas	(0.9, 0.999)	(0.9, 0.999)
Adam epsilon	1e-8	1e-8
Max grad norm	1.0	1.0
Num epochs	30	3
LR scheduler type	Linear	Linear
Warmup steps	0	0

Table 11: Hyperparameters for finetuning T5 and GPT-2 on the knowledge graph.

<b>Model</b>	<b>Avg Out Len</b>	<b>Uniq Word #</b>
T5	5.29	392
GPT-2	5.03	451
GPT-2 <sub>X</sub>	5.03	436

Table 12: Average output length and the number of unique words.



Model	Rel	Encoder Input	Decoder Input
T5	xNeed	この文の前に起こる イベントは何ですか?: $h$ (What event occurs before this statement?: $h$ )	$t$
	xEffect	このイベントの次に発生する事象は何ですか?: $h$ (What is the next event to occur after this event?: $h$ )	$t$
	xIntent	次の文の発生した理由は何ですか?: $h$ (What is the reason for the occurrence of the following statement?: $h$ )	$t$
	xReact	次の文の後に感じることは何ですか?: $h$ (What will be felt after the following statement?: $h$ )	$t$
GPT-2	xNeed	-	$h$ xNeed $t$
	xEffect	-	$h$ xEffect $t$
	xIntent	-	$h$ xIntent $t$
	xReact	-	$h$ xReact $t$

Table 13: The input formats for training. Note that  $h$  and  $t$  denote a head and a tail.

Model	Input	Output
T5	この文の前に起こる イベントは何ですか?:あなたが友人たちと旅行に出かける (What event occurs before this statement?: You go on a trip with your friends)	あなたが車を運転する (You drive a car)
	このイベントの次に発生する事象は何ですか?:あなたが友人たちと旅行に出かける (What is the next event to occur after this event?: You go on a trip with your friends)	あなたが楽しい時間を過ごす (You have a good time)
	次の文の発生した理由は何ですか?:あなたが友人たちと旅行に出かける (What is the reason for the occurrence of the following statement?: You go on a trip with your friends)	楽しい (Have fun)
	次の文の後に感じることは何ですか?:あなたが友人たちと旅行に出かける (What will be felt after the following statement?: You go on a trip with your friends)	楽しい (Have fun)
GPT-2	僕が友人たちと旅行に出かけるxNeed (I go on a trip with your friends xNeed)	僕がパスポートを取得する (I get my passport)
	僕が友人たちと旅行に出かけるxEffect (I go on a trip with your friends xEffect)	僕が楽しい時間を過ごす (I have a good time)
	僕が友人たちと旅行に出かけるxIntent (I go on a trip with your friends xIntent)	楽しいことがしたい (Want to have fun)
	僕が友人たちと旅行に出かけるxReact (I go on a trip with your friends xReact)	楽しい (Feel fun)

Table 14: Examples of the inferences generated by T5 and GPT-2.