

Can Language Models Take A Hint? Prompting for Controllable Contextualized Commonsense Inference

Anonymous ACL submission

Abstract

Generating commonsense assertions, given a certain story context, is a tough challenge even for modern language models. One of the reasons for this may be that the model has to "guess" what topic or entity in a story to generate an assertion about. Prior work has tackled part of the problem, by providing techniques to align commonsense inferences with stories and training language generation models on these. However, none of the prior work provides means to control the parts of a generated assertion. In this work, we present "hinting", a data augmentation technique for improving inference of contextualized commonsense assertions. *Hinting* is a prefix prompting strategy that uses both hard and soft prompts. We demonstrate the effectiveness of *hinting* by showcasing its effect on two contextual commonsense inference datasets: ParaCOMET (Gabriel et al., 2021) and GLUCOSE (Mostafazadeh et al., 2020), for both general and context-specific inference.

1 Introduction

The task of Contextual or Discourse-Aware Commonsense Inference, which consists of generating relevant and coherent commonsense assertions (i.e. facts) for a certain sentence in a story context, while easy for humans, remains challenging for machines (Gabriel et al., 2021). Within this task, we define an *assertion* as a tuple that contains a *subject*¹, a *relation*, and an *object* (e.g., *a dog, is a, animal*), similar to a subject-verb-object triple. An assertion in this task can be seen as contextually specific facts or generally applicable rules, that can be inferred from a sentence in a given story context.

Automated systems (such as pre-trained transformer-based language models (Devlin

¹We note that here we utilize the term "subject" as a part of the relation tuple, and it is not necessarily "subject" in a grammatical sense. In the case of ATOMIC, a subject could be a sentence describing an event that causes another event or a reaction, whereas in ConceptNet it could be a concept.

et al., 2019; Radford et al., 2019)) struggle with generating these contextual assertions, since there is an implicit assumption that clues for making predictions can always be found explicitly in the text. (Da and Kasai, 2019; Davison et al., 2019; Liu and Singh, 2004; Zhang et al., 2021). This becomes problematic because the model is essentially forced to use knowledge that it may not have seen during pre-training. Additionally, models are forced to guess what to predict about (e.g., what *the subject* of an assertion is), which may lead to decreased performance (e.g., the model generates an assertion about cats when it should have talked about dogs).

To clarify the task of contextual commonsense inference even further, below we give an example with a story, a *target sentence*, and some corresponding story specific and general inferences. The story is picked directly from the ROCStories corpus. (Mostafazadeh et al., 2016)

Story: The hockey game was tied up. The **red team** had the puck. They sprinted down the ice. They cracked a shot on goal! *They scored a final goal!*

Story Specific Commonsense Inference: *The red team, is capable of, winning the game*

General Commonsense Inference: *Some people scored a final goal, causes, some people to be happy*

In this example we can see the aforementioned problems that models have to deal with. Although it is commonsense that a final goal will lead to a victory for the red team, it is not explicitly stated in the text. Pre-training of models may include text related to a sudden death goal from sources such as Wikipedia, but the model has to extrapolate that the final goal in this example is a type of sudden death goal and will concede victory to the red team. Similarly, although we may want to talk about the red team, the model has to somehow know that it *needs* to talk about this, and **not** something else.

Previous attempts have tackled the problem of contextual commonsense inference by constructing datasets of stories aligned with assertions (i.e. an assertion is given for a sentence in a story), either through automated or human-annotated ways, and building a model to, given the story and a target sentence, predict part or the whole assertion. One previous attempt to tackle this problem, ParaCOMET (Gabriel et al., 2021), trained a GPT-2 language model (Radford et al., 2019) to infer an object of a commonsense assertion tuple from the ATOMIC (Sap et al., 2019) knowledge base². They formulate the task as follows. Given a story, a sentence identifier token, and a specified relation, the model has to predict the object of a commonsense assertion. An example of an input and expected output from the ParaCOMET formulation can be seen below:

Model Input: The hockey game was tied up. The red team had the puck. They sprinted down the ice. They cracked a shot on goal! *They scored a final goal!* <|sent5|> <|xEffect|>
Model Target/Output: win the game

In this example, since the model is predicting ATOMIC objects, the output is a single phrase (i.e., win the game). Additionally, the symbols <|sent5|> and <|xEffect|> mean that the target sentence is sentence number five³, and that the relation we want to generate a tuple about is the "has the effect on a certain person(s)" respectively.

Another work that tries to approach this is GLUCOSE (Mostafazadeh et al., 2020). Here a dataset is constructed to consist of stories and human annotations for sentences in the stories. The human annotations provide specific and general commonsense assertions. The authors utilize this dataset to train a T5 (Raffel et al., 2020) model to perform contextualized and generalized story assertion inference. The model takes an input sequence in the form of a story, a relation to predict, and a target sentence, and has to predict both the general and specific assertions that may be present in the target sentence with the given story context. An example of the GLUCOSE formulation’s inputs and expected outputs is given below:

Model Input: 1: The hockey game was tied up. The red team had the puck. They sprinted down the ice. They cracked a shot on goal! **They scored a final goal!**
Model Target/Output: The red team scores, Causes/Enables, they win the game ** People_A score, Causes/Enables, People_A win a game

This formulation of contextual commonsense inference is harder than the ParaCOMET one in that it has to generate two sets of a subject, relation, and an object tuples, in which one is the story specific one and the other is the general version of the assertion. These are seen above separated by the ** respectively. In this example additionally, we can see the symbol *I*: which tells the model to predict along a dimension of commonsense described by GLUCOSE (i.e., 1: Event that directly causes or enables X), and the sentence enclosed by asterisks (*) which signifies it is the target sentence. In both works, the models are expected to do their inference from the story, a target sentence, and relation alone.

Recently, there has been work on exploring prompting (Liu et al., 2021), which is essentially finding ways of altering the input to a language model such that it matches templates that it has seen during pre-training. Prompting a model correctly gives stronger performance in tasks, can help with controllability in the case of text generation, and is more parameter-efficient and data-efficient than fine-tuning, in some cases (Li and Liang, 2021). One novel type of prompting is *prefix prompting* (Li and Liang, 2021; Lester et al., 2021). Prefix prompting consists of modifying a language model’s input (i.e. prefix) by adding additional words. These words can be explicit hard prompts (i.e., actual words such as "give a happy review") or they can be soft prompts, embeddings that are input into a model and can be trained to converge on some virtual template or virtual prompt that can help the model. Prompting holds great potential for improving contextualized commonsense inference.

We introduce the idea of a *hint*, a hybrid of hard and soft prompts. We define a *hint* as an additional input in the form of a part(s) of an assertion that a model has to predict, along with special identifiers for these parts, wrapped within parenthesis characters. Syntactically, a *hint* would take the form of: "([subject symbol, subject], [relation symbol, relation], [object symbol, object])" where the actual content of the *hint*, between the parenthesis, would be a permutation of all but one of the elements of

²ATOMIC is composed of causal assertions, where a certain subject event, causes a certain object event through a given relation.

³We note that in the original ParaCOMET work, the sentences were 0-start indexed. We utilize 1-start indexing for clearer understanding.

the target tuple during training. In the case of supplying hints to GLUCOSE, we include a specific or general symbol which determines whether the part of the hint belongs to a story specific assertion or a general rule. For example, a *hint* for the hockey example for the GLUCOSE formulation would be "`(<|specific|><|subj|>The red team scores, <|general|><|obj|>People_A win the game)`". Altogether, the model's input would be:

Model Input: 1: The hockey game was tied up. The red team had the puck. They sprinted down the ice. They cracked a shot on goal! *They scored a final goal!*(`<|specific|><|subj|>The red team scores, <|general|><|obj|>they win the game)`)

Model Target/Output: The red team scores, Causes/Enables, they win the game ** People_A score, Causes/Enables, People_A win a game

Hints are provided during training by sampling a binomial distribution (with $p = 0.5$) for each element in a minibatch, which determines whether to give a *hint* or not. The actual content of the *hint* would then be generated by randomly sampling without replacement up to all but one of the elements in a target tuple. We give a more detailed description of *hinting* in section 3.2 and 3.3.

We hypothesize that this scheme of *hinting* strikes a balance between the model recalling information from its pre-training, with information that it may not have seen that may only be present in the target tuple. Additionally, by providing and fine-tuning a model on the combination of hard and soft prompts, a generative language model can be guided to "talk" about a certain **subject**, **object**, or **relation**, thus enabling finer control of models in downstream applications. We note that the approach was designed to be simple to implement, and to give control when generating text. In the following sections, we give some background and follow this with a set of experiments to show the effects of *hinting* for the ParaCOMET and GLUCOSE datasets, and finally, analyze the results, and present future directions for this work. Concretely, our contributions are:

- A hybrid prefix prompting technique called *hinting* that provides a partial assertion to augment data for contextual commonsense inference, and
- Demonstrating that *hinting* does in fact improve the performance for contextual commonsense inference as measured by auto-

ated metrics and is comparable in human-based metrics.

2 Related Work

2.1 Prompting

Recently, there has been a shift in paradigm in Natural Language Processing from pre-training and fine-tuning a model, to pre-training, prompting, and predicting (Liu et al., 2021). One primary reason for this shift is the creation of ever-larger language models, which have become computationally expensive to fine-tune. Prompting can be described as converting a pre-trained language model input sequence into another sequence that resembles what the language model has seen during pre-training. Overall, most prompting research is focused on formulating the task as a *cloze* (fill-in-the-blanks) task. However, we consider the task of language generation, an open-ended formulation.

Recall that prefix prompting modifies the input to a language model, by adding either a hard prompt (additional words to the input sequence)(Shin et al., 2020) or a soft prompt (i.e., adding trainable vectors that represent, but are not equivalent to, additional words) (Li and Liang, 2021; Lester et al., 2021; Liu et al., 2021).

Unlike classic prefix prompting, *hinting* uses both hard and soft prompts. The soft prompts are in the form of symbols that represent the different parts of the assertion (i.e., **subject**, **relation type**, and **object**), and the hard prompts are in the form of the actual parts of the assertion that are selected to be appended as part of the *hint*. Our work is similar to KnowPrompt (Chen et al., 2021a), except that they use a masked language model and soft prompts for relationship extraction. AutoPrompt (Shin et al., 2020) is also similar, but finds a set of "trigger" words that give the best performance on a *cloze*-related task, whereas we provide specific structured input for the model to guide text generation. We additionally note that although there are prompt-based relation extraction models (Chen et al., 2021b), we are performing a different task which is contextual commonsense inference.

2.2 Controllable Generation

Controllable generation can be described as ways to control a language model's text generation given some kind of guidance. One work that tries to implement controllable generation is CTRL (Keskar et al., 2019). The authors supply control signals

during pre-training of a general language model. A body of work in controllable generation has focused on how it can be used for summarization. Representative work that uses techniques similar to ours is GSum (Dou et al., 2021). In contrast to GSum, our method is model independent, allows for the source document to interact with the guidance signal, and contains soft prompts in the form of trainable embeddings that represent the parts of a tuple. The GSum system gives interesting insight into the fact that highlighted sentences, and the provision of triples, does in fact help with the factual correctness of abstractive summarization. We make the distinction that *hinting* falls more under prompting for the reason that we utilize additionally the trainable soft embeddings rather than purely additional hard tokens and that our task of contextual commonsense generation is not explored in the controllable generation works, whose main focus is on controlling unstructured text generation. Some works that are in this area are also (Peng et al., 2018) who utilize what they call "control factors" as keywords or phrases that are supplied by a human-in-the-loop to guide a conversation. More similar to our work, but tailored for the task of interactive story generation and without trainable soft-embeddings, is the work by (Brahman et al., 2020) which uses automatically extracted keywords to generate a story. Future work we could possibly utilize the automatic keyword extraction to supply parts of a hint, rather than our approach of complete parts of an assertion, and expand this to utilize synonyms of keywords. Lastly, there is the work by (See et al., 2019) which looks at controllable text generation for the purpose of conversation and utilizes an embedding give quantitative control signals as part of conditional training.

2.3 Discourse-aware/Contextual commonsense inference

Commonsense inference is the task of generating a commonsense assertion. Discourse-aware/contextual commonsense inference is the task of, given a certain narrative or discourse, inferring commonsense assertions that are coherent within the narrative (Gabriel et al., 2021). This task is particularly hard because commonsense knowledge may not be explicitly stated in text (Liu and Singh, 2004) and the model needs to keep track of entities and their states either explicitly or implicitly. Research into the knowledge that pre-trained

language models learn has yielded good results in that they do contain various types of factual knowledge, as well as some commonsense knowledge (Da and Kasai, 2019; Petroni et al., 2019; Davison et al., 2019). The amount of commonsense knowledge in these models can be improved by supplementing sparsely covered subject areas with structured knowledge sources such as ConceptNet (Speer et al., 2017; Davison et al., 2019).

Knowing that these pre-trained language models may contain some commonsense information has led to the development of knowledge models such as COMET (Bosselut et al., 2019). This line of research has been extended from the sentence-by-sentence level in COMET, to the paragraph-level in ParaCOMET (Gabriel et al., 2021). Contemporaneously, GLUCOSE Mostafazadeh et al. (2020) builds a dataset of commonsense assertions that are contextualized to a set of stories, and generalized.

3 Modeling

3.1 Task

We now detail the task of Contextual Commonsense Inference. We are given a story S composed of n sentences, $S = \{S_1, S_2, \dots, S_n\}$, a target sentence from that story, S_t , where $S_t \in S$, and a dimension/relation type R . Given all this, we want to generate a tuple in the form of (*subject*, R , *object*) that represents an assertion, present or implied, in S_t given the context S , and the relation type R .

We run tests with two variations of this task, one is the ParaCOMET variation and the other the GLUCOSE variation. In the ParaCOMET experiments, we represent S_t with a unique token. Additionally, we only generate the *object* of the tuple.

In our GLUCOSE experiments, we represent S_t by marking it with * on the left and right of the sentence. Additionally, we generate two *subject*, R , *object* tuples, one that is the context-specific tuple, and the other is the general tuple, separated by two asterisks (**). An example of both is in Section 1.

3.2 Hinting

The mechanism we present in this work, called *hinting*, is a kind of mixed/hybrid prompting for generative language models. Prompting is essentially supplying additional text (i.e. prompts) to a language model to aid/guide it in a specific task. In our case, we opt to give a "hint", as to what the assertion that we want to predict contains, at the end of our input text. We chose placing the

Model Input	Model Output
7: The hockey game was tied up. The red team had the puck. They sprinted down the ice. They cracked a shot on goal!. * They scored a final goal!. *	They scored a final goal >Causes>They feel(s) happy ** Some People_A scored a final goal >Causes>Some People_A feel(s) happy
7: The hockey game was tied up. The red team had the puck. They sprinted down the ice. They cracked a shot on goal!. * They scored a final goal!. * hint: (< specific >< subj > the red team scores the final goal)	the red team scores the final goal >Causes>the red team feel(s) happy ** Some People_A (who are a team) score the final goal >Causes>Some People_A feel(s) happy
7: The hockey game was tied up. The red team had the puck. They sprinted down the ice. They cracked a shot on goal!. * They scored a final goal!. * hint: (< specific >< subj > the blue team does not score the final goal)	the blue team does not score the final goal >Causes>the blue team feel(s) disappointment ** Some People_A do not score the final goal >Causes>Some People_A feel(s) disappointment
7: The hockey game was tied up. The red team had the puck. They sprinted down the ice. They cracked a shot on goal!. * They scored a final goal!. * hint: (< specific >< obj > a child)	they scored a final goal >Causes> a child feel(s) happy ** Some People_A scored a final goal >Causes>Someone_A feel(s) happy
7: The hockey game was tied up. The red team had the puck. They sprinted down the ice. They cracked a shot on goal!. * They scored a final goal!. * hint: (< general >< subj > Something_A (that is a point))	They scored a final goal >Causes>They feel(s) happy ** Something_A (that is a point) is scored >Causes>Some People_A feel(s) happy

Table 1: Example of inputs and outputs for the GLUCOSE trained model with hints. The *hint* is **bolded** and the parts of the *hint* are colored (**subject**, **relation**, **object**). Without a hint we can see that the model tries to infer directly on the content of the sentence, however with *hints*, the model tries to include an inference based on the target sentence with the contents of the *hint*.

hint at the end of the input for simplicity in dataset processing, but it can be placed anywhere and we leave it as future work to explore the effects of placing hints possibly next to the target sentence or at the beginning of the input. Hinting can be seen as a hybrid of prompting the generative model with hard prompts composed of parts of what should be predicted along with soft prompts of symbols that represent those parts. These symbols are for the **subject**, **relation**, and **object** respectively. These soft prompts utilize untrained embeddings for the task. We structure hinting this way such that, after training, whenever a *hint* is given, the model can be guided to generate knowledge about the *hint*'s content based on the target sentence and context.

To balance the model's reliance on the context, its knowledge, and the *hint*, we determine whether to supply the *hint* by sampling a binomial distribution ($p = 0.5$). Thus, we can control the frequency of when to supply a *hint*. Additionally, the content of the *hint* is determined by random sampling of permutations of components, up to a maximum of all but one component. Since our task is to predict the tuple, we do not want to make the model overly reliant on *hints* for the answer.

3.3 An example of Hinting

A simple example of *hinting* is the following:

Story: *The hockey game was tied up. The red team had the puck. They sprinted down the ice. They cracked a shot on goal! They scored a final goal!*

Target sentence: *They scored a final goal!*

Target assertion: (*subject: the red team, relation: are capable of, object: winning the game.*)

A *hint* can be any permutation of the target assertion, except the complete assertion, along with some symbol that indicates which part it is:

Possible Hints: (<|subj|> *the red team*), (<|subj|> *the red team*, <|rell|> *capable of*), (<|subj|> *the red team*, <|obj|> *winning the game*), (<|rell|> *capable of*, <|obj|> *winning the game*), (<|obj|> *winning the game*), (<|rell|> *capable of*)

A *hint* for the given story, target sentence and target assertion, yields the following:

Hint: (<|subj|> *the red team*, <|rell|> *capable of*)

Putting everything altogether, the input for the model would be:

Story with Hint: *The hockey game was tied up. The red team had the puck. They sprinted down the ice. They cracked a shot on goal! They scored a final goal! (<|subj|> the red team, <|rell|> capable of).*

We note that this is a general version of how the hinting mechanism works. The dataset specific hints that we utilize are described in Section 4.1.

3.4 Models

For our first set of experiments, we utilize the ParaCOMET (Gabriel et al., 2021) dataset and framework with the same GPT-2 model as ParaCOMET, along with a BART (Lewis et al., 2020) model and a T5 (Raffel et al., 2020) model to observe the effects of *hinting* in a sequence-to-sequence formulation of the dataset. We use the off the shelf (Huggingface (Wolf et al., 2019)) pretrained "base" version of these models for efficiency. For our second set of experiments with the GLUCOSE dataset, we also utilize the T5 model as was done in GLUCOSE.

4 Experimental Setup

4.1 Experiment Description

We run two sets of experiments to show the effectiveness of *hinting*. The first is utilizing the original ParaCOMET dataset and setup and adding *hints*. The ParaCOMET setup consists of given a story S composed of n sentences, $S = \{S_1, S_2, \dots, S_n\}$, a relation type R , and a target sentence token (i.e. $\langle |sent0| \rangle, \langle |sent1| \rangle, \dots, \langle |sent(n-1)| \rangle$). In the ParaCOMET dataset, we must predict the **object** of a triple, utilizing implicitly the sentence as a **subject** and explicitly the supplied sentence symbol and **relation** R symbol.

Within this framework, after the relation R , we add our *hint* between parenthesis (i.e. “([*hint*])”). In this framing, our *hint* can be composed of: a subject symbol ($\langle |subj| \rangle$) along with the target sentence to serve as a **subject**, a relation symbol along with the **relation** R , or an object symbol along with the **object** of the triple. Using the hockey example a possible *hint* in this set of experiments would be: “($\langle |rell| \rangle \langle |xEffect| \rangle, \langle |obj| \rangle$ they win the game)”.

In our GPT-2 experiments, we utilize the same cross-entropy loss as in (Gabriel et al., 2021). We note that we utilize a sequence-to-sequence (Sutskever et al., 2014) formulation for the T5 and the BART models. This in contrast to the GPT-2-based system requires encoding a source sequence (i.e., story, target sentence, and relation symbol), and decoding it into a target sequence (i.e., the **object** of an assertion). For the T5 model, we add the prefix “source:” before the story S , and the prefix “hint:” for placing our *hints*. For simplicity, we construct the same “heuristic” dataset as ParaCOMET which utilizes a heuristic matching technique to align ATOMIC (Sap et al., 2019) triples to story sentences.

For our second set of experiments, we utilize the formulation utilized in GLUCOSE (Mostafazadeh et al., 2020). The formulation utilizes the T5 model in a sequence-to-sequence formulation once more. In this formulation, the source text is composed of a prefix of a dimension to predict $D \in 1, 2, \dots, 10^4$, followed by the story S with the marked target sentence. The target sentence, S_t , is marked with * before and after the sentence. An example input is: “1: The first sentence. *The target sentence. *The third sentence.”. This task is slightly different from the ParaCOMET one, in that in addition to

⁴The definition for each dimension number is given in the GLUCOSE work

predicting a context specific triple, the model has to predict a generalized triple. In this task we have to infer a general and context specific **subject**, **object** and a **relation**. For our *hints* we provide up to five out of these six things, along with a symbol that represents whether it is the **subject**, **object** or a **relation**, and another symbol that represents whether it is part of the general or specific assertion. We add our *hint* after the story S , utilizing the prefix “hint:” and supplying the *hint* between parenthesis. In this set of experiments, an example of our *hint* can be, given the example in section 3.3: “($\langle |general| \rangle \langle |obj| \rangle$ People_A win a Something_A)”.

4.2 Experiment Configuration

We run the ParaCOMET experiments for 10 epochs on the dataset’s training data and evaluation data. We utilize a max source sequence length for the BART and T5 models of 256, and a max target length of 128. For the GPT-2 models we utilize a max sequence length of 384. Additionally, we use the ADAM (Kingma and Ba, 2015) optimizer with a learning rate of $2e-5$, and a linear warm-up of 0.2 percent of the total iterations. For the T5 models we utilize a learning rate of $1e-4$ because early experiments showed that the model would not converge with lesser learning rates. We utilize the scripts from (Gabriel et al., 2021) for data generation. We also utilize a batch size of 4 for training and we accumulate gradients for 4 steps for an effective batch size of 16. The results that we present are the average of the 10 runs over 4 seeds for hinted and non-hinted conditions.

We run GLUCOSE experiments for 5 epochs and 4 seeds on the original GLUCOSE data. Additionally, we utilize a linear warm up of 3000 steps. We utilize the ADAM optimizer with a learning rate of $3e-4$, a train batch size 4, with gradient accumulation of 4 steps for an effective batch size of 16, and a max source length of 256 and max target length of 128. In our results we present the average of the 4 seeds across the 5 epochs. In both experiments we report the scores given by SacreBLEU (Post, 2018), ROUGE (Lin, 2004), and METEOR (Banerjee and Lavie, 2005) using the datasets library (Lhoest et al., 2021) metrics system. We run our experiments in a machine with an AMD ThreadRipper 3970 Pro and 4 NVIDIA A6000s. Every epoch per model is approximately an hour.

Additionally, we run a small Mechanical Turk study similar to the one presented in the original

Model	BLEU		METEOR		ROUGE1		ROUGE2		ROUGE L		ROUGE L SUM	
	Hint	No Hint	Hint	No Hint	Hint	No Hint	Hint	No Hint	Hint	No Hint	Hint	No Hint
ParaCOMET	42.705*	41.960	59.411*	59.045	63.339*	61.454	52.483*	50.513	63.292*	61.395	63.294*	61.399
Bart	41.765*	41.639	58.766*	58.639	61.054	61.013	49.970	49.889	61.004	60.964	61.010	60.969
T5	41.070	41.102	58.004	58.000	59.535	59.631	48.695	48.823	59.488	59.588	59.494	59.597

Table 2: Averages of 4 different seeds over 10 epochs for *hinted* (Hint) and non-*hinted* (No Hint) runs of the ParaCOMET dataset from (Gabriel et al., 2021). The largest scores are **bolded** and significantly different scores have an asterisk (*) next to them. We can see from the results that *hinted* systems tend to achieve higher performance even if slightly and in some cases significantly, and do not decrease performance significantly.

Bleu		Meteor		Rouge 1		Rouge 2		ROUGE L		Rouge LSUM	
No Hint	Hint	No Hint	Hint	No Hint	Hint	No Hint	Hint	No Hint	Hint	No Hint	Hint
58.542	59.099*	66.829	66.917	66.387	66.681*	47.850	48.141	62.542	62.874*	62.528	62.868*

Table 3: Averages of 4 different seeds over 5 epochs for *hinted* (Hint) and non-*hinted* (No Hint) runs of the GLUCOSE contextual inference task dataset. This is the same dataset as the work in (Mostafazadeh et al., 2020) The largest scores are **bolded** and significantly different scores have an asterisk (*) next to them. Once more, we see that hinting provides a small, increase in performance, all the while permitting controllability.

Model	Non-Hinted	Hinted
ParaCOMET	3.71	3.76
Bart	3.72	3.48
T5	3.73	3.68
T5-GLUCOSE	4.10	4.06

Model	Non-Hinted	Hinted
ParaCOMET	81%	84%
Bart	83%	74%
T5	81%	81%
T5-GLUCOSE	92%	90%

Table 4: Results of human evaluation of ParaCOMET and GLUCOSE datasets. The largest scores are **bolded** and significantly different scores have an asterisk (*) next to them. We sampled 100 test points for each model from their test datasets and had the *hinted* and non-*hinted* models infer assertions. Humans judged these assertions on a 5 point Likert scale where above 3 was plausible similar to (Gabriel et al., 2021). On the left we can see the average values of the human judgments and on the right we can see the percentage of plausible inferences (rated ≥ 3). We can see that hinting provides comparable performance.

547 ParaCOMET (Gabriel et al., 2021) in which a hu- 566
548 man judges a generated assertion and judges the 567
549 plausibility of it on a 5-point Likert scale: obvi- 568
550 ously true (5), generally true (4), plausible (3), neu- 569
551 tral or unclear (2), and doesn't make sense (1). We 570
552 present the results in the same manner where Ta- 571
553 ble 4 displays the percent of inferences judged as 572
554 plausible or true (3-5), and the average rating per 573
555 inference. Participants were given \$0.1 to complete 574
556 the task. We give an image of the HIT in Appendix 575
557 A. We sample from each of the ParaCOMET and 576
558 GLUCOSE test sets, 100 entries. Then based on 577
559 the models for each dataset, we pick the epoch that 578
560 had the highest automated scores and we proceed 579
561 to randomly sample one of the trained hint and non- 580
562 hinted models. We then select one sentence of the 581
563 randomly sampled test entries and ask both models 582
564 to generate an inference along a randomly sampled 583
565 relation or dimension for that sentence. 584

5 Results and Analysis 566

5.1 Experiment 1: ParaCOMET with hints 567

568 The aggregated results for this set of experiments 569
569 can be found in Table 2. We can see here that on 570
570 average, *hinting* does tend to improve the score 571
571 even if slightly. It seems that providing a *hint* is 572
572 beneficial and not detrimental for contextual com- 573
573 monsense inference. Given the way that this task 574
574 is framed, a possibility that could explain the rela- 575
575 tive similarity of the performances, is that *hinting* 576
576 *only* adds the **object** of the triple as additional pos- 577
577 sible data that the model may see during training; 578
578 the **subject** and the **relation** can be repeated with 579
579 *hinting*. We note that the performance of the T5 580
580 model was less than that of the other models, and 581
581 we believe that it may be lack of hyperparameter 582
582 tuning, as it was seen that the model was sensitive 583
583 to the learning rate and had to use a higher than 584
584 usual learning rate.

5.2 Experiment 2: GLUCOSE with hints 585

585 The aggregated results for this set of experiments 586
586 can be found in Table 3. Onc more we notice that 587

hinting does tend to improve the performance of the contextual commonsense inference task. This suggests that *hinting* is indeed beneficial for the task of contextualized commonsense inference, especially when faced with the harder task of generating both a general and context dependent assertion. We believe that this improvement is because *hinting* gives the model the clues it may need to decide on what to focus or attend to, to generate useful inferences.

5.3 Experiment 3: Human Judgements

The results for a small Mechanical Turk study for human evaluation of model inferences can be seen in Table 4. Overall we can see here that hinted systems are judged as less plausible. Interestingly, after inspecting the results where there was a large difference (more than two points between), we see that there are some cases in which the same or very similar responses got completely different scores. We also see upon looking some of the inferences that the hinted model tends to be more general and provide shorter responses than the non-hinted model (e.g., hinted inference: "satisfied" vs. non-hinted inference: "happy and satisfied").

5.4 Discussion: Why *hint*?

From the results of our experiments, we can see that *hinting* tends to increase the performance of contextualized commonsense inference at least with regards to automated metrics and does not significantly degrade or improve human judgements. This brings the question of: Why *hint* at all? The primary reason is for controllability in the generation. By supplying these *hints*, we are teaching the model pay attention and generate inferences about a certain **subject**, **relation**, or **object**. This in turn, after training, can be leveraged by a user or downstream application to guide the model to generate assertions from parts that are manually supplied. Although this is not very clear within the ParaCOMET formulation, it becomes clearer in the GLUCOSE formulation of the problem. We give an illustrative example of the usefulness of *hinting* in Table 1. We can see that by giving a model the *hint*, the model could be capable of inferring about information that may not be present in the story. We note that this behavior is useful in downstream tasks such as story understanding and contextual knowledge graph generation in which we may need a model to have a specific **subject** or **object**. Lastly, *hinting* was designed to be simple to implement, and is model independent.

5.5 Discussion: Is *hinting* optimal?

This work was a proof of concept for this technique. We acknowledge there is a large body of research on the area of prompting. The way the *hinting* mechanism was designed however, leaves much space to explore alternate mechanisms such as AutoPrompt(Shin et al., 2020), including additional soft prompts such as those in Li and Liang (2021), or even replacing the contents of the hint with synonyms or related words. Because of the naivety of the approach, we do not think it is an optimal approach, and there is a large body of research that points to manual templating of prompts being less effective than learned prompts (Liu et al., 2021). However, from our tests, our approach does not degrade performance, and only improves it.

6 Future Work

When designing the *hinting* system certain aspects were formulated to leave space for improvements. One such area is finding a smarter way of selecting when to *hint*, and finding a smarter way of selecting what to *hint*. Additionally, more soft prompts could be added to the *hint* such that they would learn a better virtual template.

Another area to explore is providing deeper ablation studies to determine what parts of the *hint* are more effective and when. This work is more a proof-of-concept that *hinting*, or more broadly prompting, is useful towards the task of contextual commonsense inference. Furthermore, given that models trained with *hinting* for contextual commonsense inference can be guided by the information supplied in *hints*, such models can be utilized in a variety of downstream applications such as story understanding and contextual knowledge graph generation.

7 Conclusion

In this work we presented *hinting*, a simple hybrid prompting mechanism that consists of appending parts of a target tuple into an input sequence for the task of contextual commonsense inference. We showed that *hinting* tends to improve performance in automated metrics and provides comparable performance with human-based judgements. With this, we open the doors for exploring prompting within the realm of contextual commonsense inference.

References

685
686
687
688
689
690
691
692

Satanjeev Banerjee and Alon Lavie. 2005. [METEOR: An automatic metric for MT evaluation with improved correlation with human judgments](#). In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

693
694
695
696
697
698

Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Çelikilimaz, and Yejin Choi. 2019. Comet: Commonsense transformers for automatic knowledge graph construction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)*.

699
700
701
702
703
704
705
706

Faeze Brahman, Alexandru Petrusca, and Snigdha Chaturvedi. 2020. [Cue me in: Content-inducing approaches to interactive story generation](#). In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 588–597, Suzhou, China. Association for Computational Linguistics.

707
708
709
710
711

Xiang Chen, Xin Xie, Ningyu Zhang, Jiahuan Yan, Shumin Deng, Chuanqi Tan, Fei Huang, Luo Si, and Huajun Chen. 2021a. [Adaprompt: Adaptive prompt-based finetuning for relation extraction](#). *CoRR*, abs/2104.07650.

712
713
714
715
716

Xiang Chen, Xin Xie, Ningyu Zhang, Jiahuan Yan, Shumin Deng, Chuanqi Tan, Fei Huang, Luo Si, and Huajun Chen. 2021b. [Adaprompt: Adaptive prompt-based finetuning for relation extraction](#). *arXiv preprint arXiv:2104.07650*.

717
718
719
720
721
722
723

Jeff Da and Jungo Kasai. 2019. [Cracking the contextual commonsense code: Understanding commonsense reasoning aptitude of deep contextual representations](#). In *Proceedings of the First Workshop on Commonsense Inference in Natural Language Processing*, pages 1–12, Hong Kong, China. Association for Computational Linguistics.

724
725
726
727
728
729
730

Joe Davison, Joshua Feldman, and Alexander M Rush. 2019. Commonsense knowledge mining from pre-trained models. 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 1173–1178.

731
732
733
734
735
736
737
738
739

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

Zi-Yi Dou, Pengfei Liu, Hiroaki Hayashi, Zhengbao Jiang, and Graham Neubig. 2021. [GSum: A general framework for guided neural abstractive summarization](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4830–4842, Online. 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.

Nitish Shirish Keskar, Bryan McCann, Lav R Varshney, Caiming Xiong, and Richard Socher. 2019. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*.

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980.

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Quentin Lhoest, Albert Villanova del Moral, Patrick von Platen, Thomas Wolf, Yacine Jernite, Abhishek Thakur, Lewis Tunstall, Suraj Patil, Mariama Drame, Julien Chaumond, Julien Plu, Joe Davison, Simon Brandeis, Victor Sanh, Teven Le Scao, Kevin Canwen Xu, Nicolas Patry, Steven Liu, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Nathan Raw, Sylvain Lesage, Anton Lozhkov, Matthew Carrigan, Théo Matussière, Leandro von Werra, Lysandre Debut, Stas Bekman, and Clément Delangue. 2021. [huggingface/datasets: 1.13.2](#).

Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). *arXiv preprint arXiv:2101.00190*.

Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Hugo Liu and Push Singh. 2004. Conceptnet—a practical commonsense reasoning tool-kit. *BT technology journal*, 22(4):211–226.

794	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang,	Abigail See, Stephen Roller, Douwe Kiela, and Ja-	849
795	Hiroaki Hayashi, and Graham Neubig. 2021. Pre-	son Weston. 2019. What makes a good conver-	850
796	train, prompt, and predict: A systematic survey of	sation? how controllable attributes affect human	851
797	prompting methods in natural language processing.	judgments. In <i>Proceedings of the 2019 Conference</i>	852
798	<i>arXiv preprint arXiv:2107.13586</i> .	of the North American Chapter of the Association	853
		for Computational Linguistics: <i>Human Language</i>	854
799	Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong	<i>Technologies, Volume 1 (Long and Short Papers)</i> ,	855
800	He, Devi Parikh, Dhruv Batra, Lucy Vanderwende,	pages 1702–1723, Minneapolis, Minnesota. Associ-	856
801	Pushmeet Kohli, and James Allen. 2016. A cor-	ation for Computational Linguistics.	857
802	pus and cloze evaluation for deeper understanding of		
803	commonsense stories. In <i>Proceedings of the 2016</i>		
804	<i>Conference of the North American Chapter of the</i>		
805	<i>Association for Computational Linguistics: Human</i>		
806	<i>Language Technologies</i> , pages 839–849.		
807	Nasrin Mostafazadeh, Aditya Kalyanpur, Lori Moon,	Taylor Shin, Yasaman Razeghi, Robert L Logan IV,	858
808	David Buchanan, Lauren Berkowitz, Or Biran, and	Eric Wallace, and Sameer Singh. 2020. Autoprompt:	859
809	Jennifer Chu-Carroll. 2020. GLUCOSE: GeneraL-	Eliciting knowledge from language models with	860
810	alized and COntextualized story explanations .	automatically generated prompts. <i>arXiv preprint</i>	861
811	In <i>Proceedings of the 2020 Conference on Empirical</i>	<i>arXiv:2010.15980</i> .	862
812	<i>Methods in Natural Language Processing (EMNLP)</i> ,		
813	pages 4569–4586, Online. Association for Computa-		
814	tional Linguistics.		
815	Nanyun Peng, Marjan Ghazvininejad, Jonathan May,	Robyn Speer, Joshua Chin, and Catherine Havasi. 2017.	863
816	and Kevin Knight. 2018. Towards controllable story	Conceptnet 5.5: An open multilingual graph of gen-	864
817	generation. In <i>Proceedings of the First Workshop on</i>	eral knowledge. In <i>Thirty-first AAAI conference on</i>	865
818	<i>Storytelling</i> , pages 43–49.	<i>artificial intelligence</i> .	866
819	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel,		
820	Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and		
821	Alexander Miller. 2019. Language models as knowl-		
822	edge bases? In <i>Proceedings of the 2019 Confer-</i>		
823	<i>ence on Empirical Methods in Natural Language</i>		
824	<i>Processing and the 9th International Joint Confer-</i>	Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014.	867
825	<i>ence on Natural Language Processing (EMNLP-</i>	Sequence to sequence learning with neural networks.	868
826	<i>IJCNLP)</i> , pages 2463–2473, Hong Kong, China. As-	In <i>Advances in neural information processing sys-</i>	869
827	sociation for Computational Linguistics.	<i>tems</i> , pages 3104–3112.	870
828	Matt Post. 2018. A call for clarity in reporting BLEU		
829	scores . In <i>Proceedings of the Third Conference on</i>		
830	<i>Machine Translation: Research Papers</i> , pages 186–		
831	191, Brussels, Belgium. Association for Computa-	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	871
832	tional Linguistics.	Chaumond, Clement Delangue, Anthony Moi, Pier-	872
833	Alec Radford, Jeff Wu, Rewon Child, David Luan,	eric Cistac, Tim Rault, Rémi Louf, Morgan Fun-	873
834	Dario Amodei, and Ilya Sutskever. 2019. Language	towicz, et al. 2019. Huggingface’s transformers:	874
835	models are unsupervised multitask learners.	State-of-the-art natural language processing. <i>arXiv</i>	875
		<i>preprint arXiv:1910.03771</i> .	876
836	Colin Raffel, Noam Shazeer, Adam Roberts, Kather-		
837	ine Lee, Sharan Narang, Michael Matena, Yanqi		
838	Zhou, Wei Li, and Peter J. Liu. 2020. Exploring		
839	the limits of transfer learning with a unified text-to-		
840	text transformer . <i>Journal of Machine Learning Re-</i>	Zhihan Zhang, Xiubo Geng, Tao Qin, Yunfang Wu, and	877
841	<i>search</i> , 21(140):1–67.	Daxin Jiang. 2021. Knowledge-aware procedural	878
842	Maarten Sap, Ronan Le Bras, Emily Allaway, Chan-	text understanding with multi-stage training . In <i>Pro-</i>	879
843	dra Bhagavatula, Nicholas Lourie, Hannah Rashkin,	<i>ceedings of the Web Conference 2021, WWW ’21</i> ,	880
844	Brendan Roof, Noah A Smith, and Yejin Choi. 2019.	page 3512–3523, New York, NY, USA. Association	881
845	Atomic: An atlas of machine commonsense for if-	for Computing Machinery.	882
846	then reasoning. In <i>Proceedings of the AAAI Con-</i>		
847	<i>ference on Artificial Intelligence</i> , volume 33, pages		
848	3027–3035.		

A Mechanical Turk Survey

Commonsense knowledge verification

You are helping to determine whether some statement is true for most people or not given a context.

Please make sure to read the full instructions before starting.

This HIT is part of a scientific research project. Your decision to complete this HIT is voluntary. There is no way for us to identify you. The only information we will have, in addition to your responses, is the time at which you completed the survey. The results of the research may be presented at scientific meetings or published in scientific journals. Clicking on the 'SUBMIT' button on the bottom of this page indicates that you are at least 18 years of age and agree to complete this HIT voluntarily.

Instructions

1. Read the story context
2. Read dimension of commonsense
3. Read the generated contextual inference
4. Answer the question to the best of your understanding

Story Context:

Read the following story context and focus on the sentence that is **bolded**:

Dan recently entered his dog into a ugly dog contest. **As Dan arrived, he was shocked to see how other dogs looked.** In addition, Dan was shocked to see how others were looking at him. Dan realized that his dog was not as ugly as he thought. Dan laughed because his dog looked real good compared to other dogs.

Figure 1: Screenshot of the Mechanical Turk Task

Inference Dimension:

Read the following commonsense dimension:

has the effect on a person

Statement:

Read the following statement:

looks at dog

Question 1:

For the given story context and inference dimension, how would you rate the statement?

Obviously True
 Generally True
 Plausible
 Neutral or Unclear
 Doesn't Make Sense

Figure 2: Screenshot of the Mechanical Turk Task pt.2