

A Knowledge-Enhanced Pretraining Model for Commonsense Story Generation

Jian Guan¹ Fei Huang¹ Zhihao Zhao² Xiaoyan Zhu¹ Minlie Huang^{1*}

¹Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China

²School of Software, Beihang University, Beijing, China

¹Institute for Artificial Intelligence, State Key Lab of Intelligent Technology and Systems

¹Beijing National Research Center for Information Science and Technology

j-guan19@mails.tsinghua.edu.cn, f-huang18@mails.tsinghua.edu.cn,
extsuioku@gmail.com, zxy-dcs@tsinghua.edu.cn,
aihuang@tsinghua.edu.cn

Abstract

Story generation, namely, generating a reasonable story from a leading context, is an important but challenging task. In spite of the success in modeling fluency and local coherence, existing neural language generation models (e.g., GPT-2) still suffer from repetition, logic conflicts, and lack of long-range coherence in generated stories. We conjecture that this is because of the difficulty of associating relevant commonsense knowledge, understanding the causal relationships, and planning entities and events with proper temporal order. In this paper, we devise a knowledge-enhanced pretraining model for commonsense story generation. We propose to utilize commonsense knowledge from external knowledge bases to generate reasonable stories. To further capture the causal and temporal dependencies between the sentences in a reasonable story, we use multi-task learning, which combines a discriminative objective to distinguish true and fake stories during fine-tuning. Automatic and manual evaluation shows that our model can generate more reasonable stories than state-of-the-art baselines, particularly in terms of logic and global coherence.

1 Introduction

Story generation is a strong indicator of machine understanding of natural language. It is often approached as selecting a sequence of events to form a story with a reasonable logic or plot. Although existing generative models (Roemmele, 2016; Fan et al., 2018; Fan et al., 2019) can generate stories with good local coherence, they

are still struggling to plan a coherent plot and maintain a reasonable event sequence throughout the story, or they are often biased towards generating a limited set of stories with generic plots (See et al., 2019) (e.g., I have a great time), even when using the powerful generative model OpenAI’s GPT-2 (Radford et al., 2019), as shown in Table 1.

Pretrained GPT-2 has been shown to capture useful semantic and syntactic features (Alt et al., 2019), as demonstrated by state-of-the-art performance on some generation tasks such as machine translation and text summarization (Radford et al., 2019). However, compared with such tasks whose source inputs have contained sufficient information to generate desired target texts, story generation is a typical *open-ended generation* task, where only very limited information is given in the input. As shown in this paper, we observe some severe issues when applying GPT-2 to generate reasonable stories, particularly commonsense stories from a limited beginning. These issues include *repetition, logic conflicts, and lack of long-range coherence* (See et al., 2019; Holtzman et al., 2019), as exemplified in Table 1. Specifically, although GPT-2 performs reasonably well at generating some related concepts to bus (e.g., driver, and the probable destinations home or party), it completely ignores the other key entity accident in the leading context, which could be caused by its lower frequency in GPT-2’s initial training corpus (less than 7% of bus). Besides, even though the concepts are relevant, they are usually generic, and used repeatedly and illogically in the generated stories. Therefore, given limited information as input, it is extremely challenging for the subsequent generation without any external guidance—for instance, commonsense knowledge. And the difficulties lie in associating inter

*Corresponding author: Minlie Huang.

Context:

The first time I saw the results of an **accident**, I was in a bus.

Golden Continuation by Human:

We were stopped and a friend pointed out the accident on the road. A man was under a car after having been run over. People tried to move the car off him but only made it worse. After seeing that, I am much more careful when crossing the street.

pretrained GPT-2:

She was trying to get home but the other passengers would not let her go. He thought she was going to die because of her weight. She was trying to get home but the other passengers would not let her go. The bus driver would not let her go.

Fine-tuned GPT-2:

I was on my way to a party. I was on my way to a party and I 'd gotten out of my seat, and started driving. I got a flat tire, so I stopped driving. I drove to the party and had a great time.

Table 1: Story examples generated by human and GPT-2 models. The stories written by the pretrained GPT-2 and fine-tuned GPT-2 (post-trained on ROCStories [Mostafazadeh et al., 2016b]) suffer from repetition (in *italic*), bad inter-sentence coherence to the context (e.g., ignoring key entities such as **accident** in **bold**), as well as conflicting logic (underlined, e.g., *first stopped driving and then drove to the party*), in spite of their good fluency and intra-sentence coherence.

dependent commonsense knowledge for expanding a reasonable story, handling the causal relationships, as well as deciding the temporal orders between entities and events in context.

Explicitly introducing external commonsense knowledge has been shown helpful to improve language understanding and long-range coherence of generated texts (Zhou et al., 2018; Guan et al., 2019; Yang et al., 2019b). For example, for the entities in the given context of Table 1, many potentially related concepts (e.g., *run over*, *cross street*) can be inferred and predicted based on external commonsense knowledge bases such as ConceptNet (Speer and Havasi, 2012) and ATOMIC (Sap et al., 2019). These knowledge bases contain abundant semantic knowledge of concepts and inferential knowledge for commonsense reasoning. We enhance GPT-2 with such knowledge by post-training the model on the knowledge examples constructed from

these knowledge bases, which can provide additional crucial information for story generation. Empirical experiments demonstrate that training with millions of such examples helps improve the coherence and logicity of generated stories. Meanwhile, we adopt multi-task learning to address the problem of handling causal and temporal dependencies. We combine the generation objective with an auxiliary multi-label classification objective, which requires distinguishing true stories from fake stories that are constructed by randomly shuffling the sentences, replacing a sentence with a negatively sampled one, or repeating a sentence in an original story. The additional classification task empowers our model to better capture the logicity in a story implicitly, namely, modeling the causal and temporal dependencies, inter-sentence coherence, and avoiding repetition.

The main contributions of this paper are summarized as follows:

- We propose a knowledge-enhanced pretraining model for commonsense story generation by extending GPT-2 with external commonsense knowledge. The model is post-trained on the knowledge examples constructed from ConceptNet and ATOMIC, thereby improving long-range coherence of generated stories.
- To generate reasonable stories, we adopt a classification task to distinguish true stories from auto-constructed fake stories. The auxiliary task makes the model implicitly capture the causal, temporal dependencies between sentences and inter-sentence coherence, and lead to less repetition.
- We conduct extensive experiments with automatic and manual evaluation. Results show that our model can generate more reasonable stories than strong baselines, particularly in terms of logicity and global coherence.¹

2 Related Work

2.1 Neural Story Generation

Many existing neural story generation models generated stories by conditioning upon various contents such as images (Huang et al., 2016) and

¹Our implementation is available at <https://github.com/thu-coai/CommonsenseStoryGen>, and demo is available at <http://coai.cs.tsinghua.edu.cn/static/CommonsenseStoryGen>.

short text descriptions (Jain et al., 2017). Different from these studies, we consider the setting of open-ended story generation from only a limited leading context in this paper. For this task, prior studies have attempted to build specific sentence representations by modeling story entities and events to simplify the dependencies between sentences (Ji et al., 2017; Clark et al., 2018). Another line is to decompose story generation into separate steps (Martin et al., 2018; Fan et al., 2018; Wang et al., 2016; Xu et al., 2018; Yao et al., 2019; Fan et al., 2019). These models usually focused on first planning story sketches and then generating sentences from the sketches. However, improving pretrained models to generate commonsense stories is yet to be well investigated.

2.2 Pretraining

Recently, large-scale pretraining models have been widely developed in various NLP tasks. Some work leveraged pretraining to provide better language representations at the word level (Mikolov et al., 2013; Pennington et al., 2014; Peters et al., 2018) or sentence level (Le and Mikolov, 2014; Kiros et al., 2015) for various downstream task-specific architectures. However, Radford et al. (2018) and Devlin et al. (2018) suggest that these complex task-specific architectures are no longer necessary, and it is sufficient to merely fine-tune pretrained task-independent transformer language models for downstream tasks. Mehri et al. (2019) explored different pretraining methods based on language models for dialogue context representation learning. Furthermore, Radford et al. (2019) demonstrate pretrained language models (i.e., GPT-2) can perform downstream tasks better than state-of-the-art models even in a zero-shot setting (i.e., without any fine-tuning on task-specific data). Wolf et al. (2019) fine-tuned GPT-2 for personalized conversation generation, which obtains very competitive results in the challenge. However, as previous studies (See et al., 2019; Holtzman et al., 2019) observed, transferring GPT-2 directly to open-ended text generation still suffers from several issues such as repetition or lack of knowledge and inter-sentence coherence with different decoding algorithms. Besides, although Song et al. (2019) and Dong et al. (2019) extended the language model to support an encoder-decoder framework (Sutskever

et al., 2014), we build our model based on GPT-2 because of its simplicity and broad applicability.

2.3 Commonsense Knowledge

Incorporating commonsense knowledge is necessary and beneficial for language inference (LoBue and Yates, 2011; Bowman et al., 2015; Rashkin et al., 2018b), reading comprehension (Mihaylov and Frank, 2018; Rashkin et al., 2018a), and particularly for open-ended language generation, which usually requires external knowledge to enrich the limited source information. Commonsense knowledge has been demonstrated to significantly improve dialogue generation (Zhou et al., 2018), story ending generation (Guan et al., 2019), and essay generation from given topics (Yang et al., 2019b). And recently, some work also attempted to integrate external commonsense knowledge into pretrained models such as BERT (Devlin et al., 2018) to enhance language representation for reading comprehension (Yang et al., 2019a) and other knowledge-driven NLP tasks like entity typing and relation classification (Zhang et al., 2019). Besides, Sun et al. (2019) improved BERT on Chinese NLP tasks by multi-stage knowledge masking strategy to integrate phrase and entity level knowledge into the language representation. Moreover, Bosselut et al. (2019) transferred the implicit knowledge from GPT-2 by fine-tuning the model to generate an object given the subject and a relation as input in commonsense knowledge graphs, that is, automatic knowledge base construction. However, the low novelty of the generated objects showed that it could still be difficult for GPT-2 to generate commonsense texts solely based on its implicit knowledge. Therefore, we target integrating external knowledge into GPT-2 for generating more reasonable commonsense stories.

2.4 Multi-Task Learning

Incorporating other auxiliary task objectives to complement the primary goal has been shown to improve the performance in many NLP tasks such as sentiment classification (Yu and Jiang, 2016) and conversation generation (Zhao et al., 2017). Recently, multi-task learning was also used to pretrain language models to capture dependencies in context (Devlin et al., 2018; Mehri et al., 2019) and further improve pretrained models' representation power during fine-tuning (Wolf et al., 2019).

3 Methodology

The task in this work can be defined as follows: Given a one-sentence story beginning X as the leading context, the model should continue to complete a K -sentence story Y with a reasonable plot. The sentences in a generated story should have reasonable logical connections, causal relationships, and temporal dependencies with each other and with the given beginning. To this end, we devise a novel framework to leverage knowledge and handle the causal and temporal dependencies, as Figure 1 shows.

3.1 Pretrained Transformer Language Model

The transformer architecture is a general model used in language modeling (Vaswani et al., 2017), which consists of multiple transformer blocks of multi-head self-attention followed by layer-normalization and fully connected layers. Radford et al. (2019) used a 12-layer decoder-only transformer (GPT-2) (i.e., a left-to-right language model) with masked self-attention heads which are constrained in that every token can only attend to its left context. Formally, the objective in this stage is to minimize the following negative likelihood:

$$\mathcal{L}_{GPT} = - \sum_{t=1}^{|u|} \log P(u_t | u_{<t}), \quad (1)$$

$$P(u_t | u_{<t}) = \text{softmax}(\mathbf{H}_t^L \mathbf{W} + \mathbf{b}), \quad (2)$$

$$\mathbf{H}_t^l = \text{block}(\mathbf{H}_{<t}^{l-1}), l \in [1, L], \quad (3)$$

$$\mathbf{H}_t^0 = E_t + P_t, \quad (4)$$

where u is an utterance with $|u|$ tokens in total from the training corpus, u_t is the t -th tokens in u , \mathbf{H}_t^l is the l -th layer’s output at the t -th position computed through the transformer block with the masked self attention mechanism, and \mathbf{H}_t^0 is a summation of token embedding E_t and positional embedding P_t for the t -th token.

GPT-2 network is pretrained on a large-scale corpus but still suffers from many issues such as lack of necessary knowledge for commonsense story generation as aforementioned. Therefore, in this work we improve GPT-2 for generating more reasonable stories with external commonsense knowledge.

3.2 Training with Commonsense Knowledge

Commonsense knowledge can facilitate language comprehension and generation, as reported in a

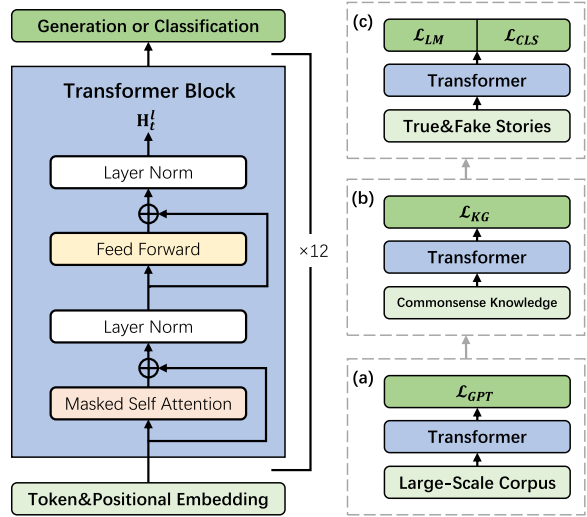


Figure 1: Transformer block architecture (left) and training framework (right). We divide the whole training framework into the following three stages. Train the language model (a) with a large-scale corpus, in which stage we directly inherit the pretrained model parameters from Radford et al. (2019), (b) with commonsense knowledge from external knowledge bases, and (c) with true and auto-constructed fake stories by multi-task learning for story generation and classification. \mathcal{L}_{GPT} , \mathcal{L}_{KG} , \mathcal{L}_{LM} , and \mathcal{L}_{CLS} are the corresponding loss functions in different stages, respectively.

notable work for dialog generation (Zhou et al., 2018). To leverage commonsense knowledge in pretrained language models, we resort to existing large-scale knowledge bases ConceptNet (Li et al., 2016b) and ATOMIC (Sap et al., 2019).

The ConceptNet dataset² consists of triples obtained from the Open Mind Common Sense entries in ConceptNet 5 (Speer and Havasi, 2012). It contains 34 relations in total and represents each knowledge triple by $R = (h, r, t)$, meaning that head concept h has the relation r with tail concept t for example, (cross street, Causes, accident). And the ATOMIC dataset³ is an atlas of everyday commonsense reasoning containing a mass of textual description of inferential knowledge organized as typed *if-then* triples. For example, a typical *if-then* triple is (PersonX pays PersonY a compliment, xIntent, to be nice), where xIntent is the relation between the head and tail events standing for

²<http://www.conceptnet.io/>.

³<https://homes.cs.washington.edu/~msap/atomic/>.

Knowledge Bases	Original Triples	Examples of Transformed Sentences
ConceptNet	(eiffel tower, AtLocation , paris) (telephone, UsedFor , communication)	eiffel tower is at paris. telephone is used for communication.
ATOMIC	(PersonX dates for years, oEffect , continue dating) (PersonX cooks spaghetti, xIntent , to eat)	PersonX dates for years. PersonY will continue dating. PersonX cooks spaghetti. PersonX wants to eat.

Table 2: Examples of template-based transformation of triples in knowledge bases. Phrases in **bold** represent the original and transformed relations.

If-Event-Then-Mental-State. We implicitly introduce the knowledge to the pretrained language model by post-training on knowledge-augmented data. Some work has attempted to explicitly incorporate commonsense knowledge into language generation (Zhou et al., 2018; Guan et al., 2019; Yang et al., 2019b). However, all these works assume that there is an alignment between the training data and the knowledge bases. Therefore, they suffer from the following issues: (1) It is difficult to match the events extracted from the training data with those stored in KB. (2) Learning and utilizing multi-hop triples in knowledge graphs is costly in time because of the large-scale size. (3) Most of KB triples do not appear in the task-specific training data, so that those absent triples are not fully utilized in existing models. Fortunately, our model is trained on the knowledge bases directly, which can effectively ease these limitations.

We transform the commonsense triples in ConceptNet and ATOMIC into readable natural language sentences using a template-based method (Levy et al., 2017), as illustrated in Table 2. We do not use roughly concatenated triples in order to avoid introducing additional special tokens (e.g., *UsedFor* in ConceptNet and *oEffect* in ATOMIC), or break the syntactic features contained in the pretrained language model (Alt et al., 2019), which are essential for following story generation. And then the language model is post-trained on the transformed sentences to learn commonsense knowledge between entities and events by minimizing the negative likelihood of predicting the next token:

$$\mathcal{L}_{KG} = - \sum_{t=1}^{|r|} \log P(r_t | r_{<t}), \quad (5)$$

where r is a transformed sentence with $|r|$ tokens in total, and r_t is the t -th token in r . In this way, we can incorporate commonsense knowledge into GPT-2 implicitly.

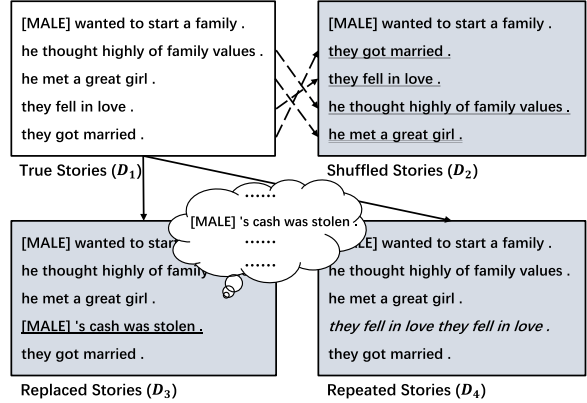


Figure 2: An example of fake story construction. The shuffled sentences are indicated by dashed lines, the replaced sentence is underlined, and the repeated one is in *italic*.

3.3 Multi-Task Learning

In order to encourage our model to generate reasonable stories in logic, we add an auxiliary classification task to the generation task during fine-tuning on the ROCStories corpus. The task requires distinguishing true stories from fake stories. We first construct three additional sets of fake stories by shuffling the sentences, replacing a sentence with a negatively sampled one, and randomly repeating a sentence in an original story. Notably, these operations are performed only on the following K sentences of a story (i.e., not including the leading context [the beginning]). For simplicity, we denote the true story set and three manually constructed fake story sets with D_1 , D_2 , D_3 , and D_4 respectively, as illustrated in Figure 2.

Our main finding is that training a language model to distinguish the reasonable stories from those with disordered logic, unrelated topics, or repeated plots is helpful to generate more reasonable stories in terms of logic and coherence. We add an additional classification layer at the last layer of the transformer language model in a multi-task setting. The classifier takes as input the hidden states of the last transformer block and

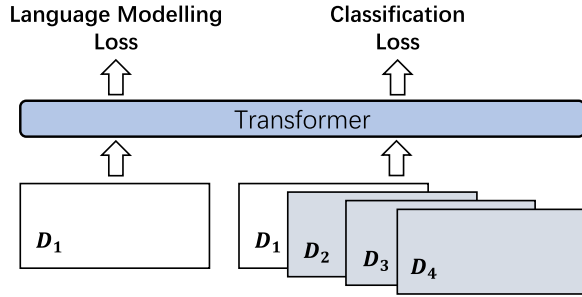


Figure 3: Multi-task learning diagram. D_1 is the true story dataset, while D_2, D_3 , and D_4 are the auto-constructed fake stories transformed from D_1 . Note that the language modeling loss is optimized only on the true stories, but the classification loss on both true and fake ones.

computes a score through a softmax layer over D_1, D_2, D_3 , and D_4 , formally as follows:

$$P(l_s|s) = \mathbf{softmax}\left(\frac{1}{|s|} \sum_{t=1}^{|s|} \mathbf{H}_t^L \mathbf{W}_L + \mathbf{b}_L\right), \quad (6)$$

where s is a true or fake story and contains $|s|$ tokens, \mathbf{H}_t^L is the hidden state of the L -th block layer (i.e., the last layer) of the transformer language model when encoding the story, l_s is predicted to indicate which dataset (D_i) the story (s) belongs to, and \mathbf{W}_L and \mathbf{b}_L are the trainable parameters of the additional classifier.

As illustrated in Figure 3, the loss function \mathcal{L}_{ST} of the full model is computed as follows:

$$\mathcal{L}_{ST} = \mathcal{L}_{LM} + \lambda \mathcal{L}_{CLS}, \quad (7)$$

$$\mathcal{L}_{LM} = - \sum_{t=1}^{|s|} \log P(s_t | s_{<t}), s \in D_1, \quad (8)$$

$$\mathcal{L}_{CLS} = - \log P(l_s = \tilde{l}_s | s), s \in D_1, D_2, D_3, D_4, \quad (9)$$

where s is a story containing $|s|$ tokens, s_t is the t -th token of s , \mathcal{L}_{LM} is the language modeling loss, \mathcal{L}_{CLS} is the classification loss, and \tilde{l}_s indicates the correct D_i which the story s is sampled from. λ is an adjustable scale factor.

4 Experiments

4.1 Dataset

We evaluated our model on the ROCStories corpus (Mostafazadeh et al., 2016a). The corpus contains 98,162 five-sentence stories for evaluating story understanding. The original task is designed to select a correct story ending from two candidates,

Dataset	Training	Validation	Test
ROCStories	88,344	4,908	4,909
ConceptNet	600,000	2,400	2,400
ATOMIC	574,267	70,683	64,456

Table 3: Statistics of datasets and knowledge bases.

whereas our task is to generate a reasonable story given the first sentence of a story (i.e., K , namely, the number of generated sentences, is four in our setting). Following Radford et al. (2019), the stories are tokenized using byte pair encoding (BPE) with a vocabulary of 50,257 items. The average number of tokens in X/Y (i.e., the beginning/the following K sentences in a story) is 13.39/50.00 with BPE, while the model uses pretrained positional embeddings with a maximal sequence length of 1024 tokens.

As for the knowledge bases, we used the 605k version of ConceptNet. The second KB we used contains 709k records from the 877k tuples of ATOMIC after transformation and deduplication. We randomly selected stories and knowledge sentences for training/validation/test respectively, as shown in Table 3. Because the ROCStories dataset is rather small for generation, we made delexicalization by replacing all the names in stories with special placeholders ‘‘[MALE]’’, ‘‘[FEMALE]’’, and ‘‘[NEUTRAL]’’ for male, female, and unknown names, respectively. Additionally ‘‘PersonX’’ and ‘‘PersonY’’ in ATOMIC are replaced by ‘‘[MALE]’’ and ‘‘[FEMALE]’’ as well.

4.2 Baselines

We compared our models with the following state-of-the-art baselines:

Convolutional Seq2Seq (ConvS2S): It directly generates a story conditioned upon the beginning based on a convolutional seq2seq model (Gehring et al., 2017) with decoder self-attention.

Fusion Convolutional Seq2Seq Model (Fusion): It generates a story by first pretraining a convolutional seq2seq model, and then fixing the model and providing it to the second clone model with fusion mechanism (Fan et al., 2018).

Plan&Write: It first generates a sequence of keywords as planning, conditioned upon the input; and then generates a story based on the planned keywords (Yao et al., 2019). During training, one

keyword is extracted from each sentence with RAKE algorithm (Rose et al., 2010).

Skeleton-based Model with Reinforcement Learning (SKRL): The model first generates a compressed story including the most critical phrases, called skeleton, and then generates a story conditioned upon the skeleton. The skeleton is automatically learned by reinforcement learning (Xu et al., 2018).

Decomposed Model with Semantic Role Labeling (DSRL): It first generates a predicate-argument structure conditioned upon the beginning and then generates a story by surface realization on top of the structure. The structures are identified by semantic role labelling (Fan et al., 2019).

We also made comparisons with GPT-2 in different settings as follows:

GPT-2 (Scratch): The network architecture is the same as GPT-2, but the model is only trained on ROCStories without any pretrained parameters.

GPT-2 (Pretrain): This model directly used the public checkpoint of pretrained parameters⁴ for story generation. Following Radford et al. (2019), stories are generated in a zero-shot setting. To induce story generation behavior, we conditioned the language model on a context of example stories, and then sample sentences from the model after a final prompt of story beginning. We used the first K generated sentences as the generated story.

GPT-2 (Fine-tuning): This model is fine-tuned on the ROCStories corpus from the public checkpoint of pretrained parameters.

Furthermore, we also conducted ablation tests by removing the proposed components respectively to investigate the influence of each component with the same network structure.

4.3 Experiment Settings

We set the parameters by following the small version of Radford et al. (2019)’s design: The language model is equipped with 12 layers, 768-dimensional hidden states, and 12 attention heads. The batch size is 10 during training on the ROCStories corpus using Adam optimizer with an initial learning rate of $1e-4$. The scale factor λ

⁴The pretrained model is available at <https://github.com/openai/gpt-2>.

is set to 0.05. And we generated stories using a top- k sampling scheme (Fan et al., 2018) with $k = 40$ and a softmax temperature of 0.7 (Goodfellow et al., 2016) to balance the trade-off between diversity and fluency. We applied these settings to all the baselines.

4.4 Automatic Evaluation

Evaluation Metrics We adopted the following automatic metrics to evaluate the generation performance in the entire test set. (1) **Perplexity (PPL)**. Smaller perplexity scores indicate better fluency in general. (2) **BLEU**. BLEU (Papineni et al., 2002) evaluates n -gram overlap between a generated story and a human-written story. However, BLEU is usually inappropriate for open-ended text generation (Fan et al., 2018) because there are multiple plausible stories for the same input but only one story is given in the dataset. And BLEU scores will become extremely low for large n . We thus experimented with $n = 1, 2$. (3) **Coverage**. To access the effect of incorporating commonsense knowledge, we calculated the coverage score as the average number of commonsense triples matched in each generated story, which requires both head and tail entities/events appears in the same story. (4) **Repetition**. We measured the redundancy of stories by computing repetition-4, the percentage of generated stories that repeat at least one 4-gram (Shao et al., 2019). (5) **Distinct**. To measure the generation diversity, we adopted distinct-4 (Li et al., 2016a), the ratio of distinct 4-grams to all the generated 4-grams.

Results The results of automatic evaluation are shown in Table 4. Note that the perplexity scores of some baselines are not comparable with ours because they tokenize stories by words rather than by byte pair encodings as used in GPT-2. Thus, we did not provide these scores. Our model outperforms the variants of GPT-2 in terms of perplexity, and has higher BLEU scores than all the baselines, indicating better fluency and more overlaps with the reference stories. Our model also has higher knowledge coverage and distinct-4 scores, showing that our model can generate more diverse stories with more abundant knowledge. However, we observed that pretraining might lead to more severe repetition by comparing three variants of GPT-2. Our model effectively improves the situation but still performs worse than the baselines with task-specific architectures,

Models	PPL	BLEU-1	BLEU-2	Coverage	Repetition-4(%)	Distinct-4(%)
ConvS2S	N/A	0.312	0.132	13.64	22.87	72.78
Fusion	N/A	0.322	0.137	12.02	24.23	72.82
Plan&Write	N/A	0.308	0.126	13.38	17.06	67.20
SKRL	N/A	0.267	0.088	10.82	18.34	69.42
DSRL	N/A	0.293	0.117	10.38	15.36	73.08
GPT-2 (Scratch)	11.82	0.311	0.134	10.76	22.87	73.33
GPT-2 (Pretrain)	33.50	0.257	0.085	8.04	39.22	64.99
GPT-2 (Fine-tune)	7.96	0.322	0.141	12.40	29.41	73.85
Ours	7.85	0.326	0.143	18.48	21.93	78.96
w/o Pretrain	11.04	0.316	0.134	16.33	21.52	77.17
w/o Knowledge	7.70	0.314	0.136	13.95	25.08	73.24
w/o Multi-task	8.04	0.324	0.140	17.19	24.40	79.43
<i>Golden Story</i>	<i>N/A</i>	<i>N/A</i>	<i>N/A</i>	<i>19.28</i>	<i>7.64</i>	<i>89.51</i>

Table 4: Automatic evaluation results. The best performance is highlighted in **bold**. The results of golden story are in *italic*. The perplexity scores marked with N/A are not comparable with ours because the corresponding models tokenize stories by words rather than by byte pair encodings used in GPT-2.

Models	Grammaticality				Logicity			
	Win (%)	Lose (%)	Tie (%)	κ	Win (%)	Lose (%)	Tie (%)	κ
Ours vs. Fusion	50.0**	27.0	23.0	0.421	57.0**	28.0	15.0	0.455
Ours vs. DSRL	58.0**	24.0	18.0	0.441	58.0**	29.0	12.0	0.475
Ours vs. GPT-2 (Scratch)	54.0**	24.5	21.5	0.385	54.0**	26.0	20.0	0.304
Ours vs. GPT-2 (Pretrain)	52.0**	31.5	16.5	0.483	56.5**	32.5	11.0	0.493
Ours vs. GPT-2 (Fine-tune)	42.0**	28.0	30.0	0.344	51.0**	27.5	21.5	0.371
Ours vs. Ours w/o Pretrain	51.0**	31.0	18.0	0.378	56.0**	28.0	16.0	0.375
Ours vs. Ours w/o Knowledge	46.0**	23.0	21.0	0.289	48.0**	29.0	23.0	0.314
Ours vs. Ours w/o Multi-task	37.5	31.0	31.5	0.313	48.5**	25.5	26.0	0.297

Table 5: Manual evaluation results. The scores indicate the percentages of *Win*, *Lose*, or *Tie* when our model is compared with a baseline. κ denotes Fleiss’ kappa (all are *fair agreement* or *moderate agreement*). The scores marked with * mean p-value < 0.05 and ** indicates p-value < 0.01 in sign test.

for instance, the planning-based models (e.g., DSRL). Fortunately, See et al. (2019) showed that increasing k for top- k sampling could alleviate the repetition issue. Additionally compared with training from scratch, fine-tuned GPT-2 performs much better in fluency (lower perplexity scores) but suffers from worse repetition, and only improve slightly in coverage and diversity. Furthermore, pretrained GPT-2 has the lowest coverage and distinct-4, which further verifies our hypothesis that GPT-2 lacks the necessary knowledge to expand a story plot.

As for the ablation test, our model without pretraining has significantly higher perplexity, indicating that pretraining contributes to story fluency. When removing external knowledge, coverage and distinct-4 drop while repetition-4

rises substantially, suggesting that post-training on millions of knowledge sentences can effectively enhance the language model’s ability to generate stories with more commonsense knowledge, although we do not explicitly utilize knowledge during fine-tuning on ROCStories. Removing multi-task learning also leads to slightly better distinct-4 but causes much higher repetition-4, indicating that the classification loss is of great help for reducing redundancy.

We also provide the performance of our model on the auxiliary story classification task and the predicted proportional distribution of the generated stories by different models on the four story types with the auxiliary story classifier, as shown in Table 6. Both metrics are computed on 1,000 samples from the test set. We can

Story types	D ₁	D ₂	D ₃	D ₄
F1 score	0.80	0.81	0.88	0.98
Models	Proportional Distribution (%)			
GPT-2 (Pretrain)	15.83	40.8	39.36	4.01
GPT-2 (Fine-tune)	86.94	9.98	2.93	0.15
Ours	90.12	7.98	1.86	0.04
w/o Knowledge	87.76	9.51	2.67	0.06
w/o Multi-task	88.69	9.07	2.02	0.22

Table 6: Final prediction **F1 score** of our model on the auxiliary story classification task in terms of the four types of story sets respectively, and the **proportional distribution** of the predicted story types of the generated stories by different models.

Models	BR (%)	LR (%)
GPT-2 (Pretrain)	59.3	44.8
GPT-2 (Fine-tune)	73.4	69.6
Ours	76.2	72.7
w/o Knowledge	74.9	71.5
w/o Multi-task	75.7	70.4

Table 7: Accuracy of beginning ranking and logic ranking. Larger scores are better.

observe that it is relatively easier to detect fake stories with repeated plots (D_4) than those with disordered logic (D_2) and unrelated topics (D_3). When using the auxiliary story classifier to classify the generated stories, pretrained GPT-2 is considered to generate more fake stories, with only 15.83% stories of type D_1 , which agrees with the previous automatic evaluation especially in terms of repetition. Besides, our model performs better than baselines, indicating that the external knowledge and the auxiliary task can encourage our model to generate more reasonable stories.

Following Fan et al. (2018) and See et al. (2019), we computed *beginning ranking accuracy* (BR) to measure how strongly the output of a model is coherent with the beginning, and *logic ranking accuracy* (LR) to measure the ability of capturing the causal and temporal dependencies in the context. For BR, we first sampled 9 negative beginnings (first sentence) for a true story, and then calculated the perplexity of the 10 stories. If the true story has the lowest perplexity by our model, it is regarded as a correct prediction. As for LR, since each story in ROCStories consists of five sentences, we produced four shuffled versions

by switching each pair of adjacent sentences. We then used our model to score the five stories with perplexity. A prediction is regarded as correct if the true story has the lowest score. We randomly sampled 1,000 human-written stories from the test set in our evaluation. As shown in Table 7, the external knowledge and multi-task learning effectively promote the coherence and help capture inter-sentence dependencies in the context.

4.5 Manual Evaluation

To evaluate the fluency and logic of generated stories, we conducted pairwise comparisons with two strong baseline models (Fusion and DSRL) that performed best in automatic evaluation, three variants of GPT-2, and three ablated models of ours. For manual evaluation, we randomly sampled 200 stories from the test set and obtained 1,800 stories from the nine models. For each pair of stories (one by our model and the other by a baseline, along with the beginning), three annotators were hired to give a preference (win, lose, or tie) in terms of two metrics, respectively. We resorted to a crowdsourcing service Amazon Mechanical Turk (AMT) for annotation, and we adopted majority voting to make final decisions among the three annotators.

Evaluation Metrics We evaluated the models from the following two perspectives: **grammaticality** to indicate whether a story is natural and fluent, and **logicality** to indicate whether a story is coherent to the given beginning and reasonable in terms of causal and temporal dependencies in the context. Note that the two aspects are independently evaluated. We show a screenshot of the annotation on AMT in Figure 4.

Results The manual evaluation results are shown in Table 5. To measure the inter-annotator agreement, we calculated Fleiss’ kappa (Fleiss, 1971) for each pairwise comparison and all the results show fair agreement ($0.2 \leq \kappa \leq 0.4$) or moderate agreement ($0.4 \leq \kappa \leq 0.6$). We also conducted a sign test to check the significance of the differences. The results indicate that our model performs significantly better than other baselines in both metrics. More specifically, post-training on knowledge bases leads to significant improvements in grammar and logic by offering more knowledge for expanding the story plots. And multi-task learning further enhances the

Evaluating the logicity and grammaticality of stories Instructions (Click to expand)

Task Description

Each story contains about five sentences . For each story, we will put the first sentence into different systems, and the following sentences will be generated by the systems. The requirement for this manual evaluation is to judge which story better complies with the **English grammar norm**, and is **more logically related to the first sentence**.

NOTE that the names in all stories are replaced with "[MALE]" or "[FEMALE]" or "[NEUTRAL]", and all the initials for each sentence are in lowercase. **They are not grammar errors.**

Evaluation Criterion

You need compare the stories from two metrics: **grammaticality** and **logicity**. **And the two metrics are independent of each other.** One of the judgements should not have any influence on the other one. Specific criteria for evaluating are as follows:

1. Grammaticality

In the process of evaluating grammaticality, it should be considered whether the statement itself complies with the English standard usage. Then annotate which story is better at grammaticality. You may not care about what the generated sentences are saying but only if **there are any grammatical problems in the sentence itself**.

2. Logicity

In the process of evaluating logicity, you need to **carefully read the whole story** including the first sentence and the generated sentences, and compare stories in logicity. Then annotate which story is better at logicity in terms of the coherence to the given beginnings and the inter-sentence causal and temporal dependencies. In this process, you may encounter sentences that are not completely grammatical. **Please make a logical evaluation based on the main part of the sentence (such as some key words, etc.) and what you can intuitively feel.** Under the circumstances, the story can be judged totally illogical only if the grammar too poor to understand the meaning or the logic is unreasonable.

Notes

- Again, the grammaticality and logicity of the story are **two independent metrics**. Some very logically inappropriate generated stories are good in the grammaticality part, and there are some stories with obvious grammatical errors but they don't affect the respective judgement.
- Sometimes, there may be more than one kind of reasonable story for a beginning. Please do not limit your imagination. **As long as the story is logically reasonable, direct and able to make sense, it can be judged good in logicity.**
- Some stories may not be accurately judged. In the process of determining the comparison of this type of two stories, according to your own understanding of the examples and the subjective feelings of the stories, choose a better story you think the most appropriate. Please ensure that your evaluation criterion for different stories is the same.
- Most importantly, in your process of evaluating, **please NOT add story details between the first sentence and the generated stories based on your imagination!**

Story

ID:
4

Story A:
i needed a good plastic drink dispenser for halloween . i went to the store . i bought a bottle of water to buy some . i looked everywhere for it for a few minutes . i also did n't like that one in any room .

Story B:
i needed a good plastic drink dispenser for halloween . i went to the store to find one . i found a cheap plastic drink . [FEMALE] friends said they would try it . i picked the cheap one and they said i would .

Which story has better grammar?

The same good or same bad

Story A

Story B

Which story has better logic?

The same good or same bad

Story A

Story B

Figure 4: A screenshot of the annotation on AMT for manual evaluation.

performance in logic and does not affect fluency of generated stories.

4.6 Relation Understanding

It is still necessary to further investigate whether our model really understands the relations between head and tail entities/events. For example, when our model learns car accident *causes* injury from ConceptNet, it will agree with car accident *leads to* injury and denies car accident *is driven by* injury if our model can identify the specific relation between the head (car accident) and

Models	Acc	Comparison Pairs		
		C vs. W	T vs. C	T vs. W
GPT-2 (Pretrain)	39.28	53.83	44.74	49.87
GPT-2(Fine-tune)	47.48	60.31	39.57	56.01
Ours	67.07	71.91	55.76	79.89
w/o Knowledge	48.07	62.07	42.43	55.64

Table 8: Accuracy (Acc, %) of relation ranking and winning rates (%) of pairwise comparisons which require selecting a more reasonable sentence from two candidates, each from Correct (C), Wrong (W), or Training (T) templates.

tail (injury). By contrast, the model will not distinguish the three statements if it only learns simple relevance (or, co-occurrence) between car accident and injury instead of the specific causal relation.

Therefore, we constructed two sets of sentences including **correct** and **wrong** knowledge respectively based on the test set of ConceptNet. Specifically, the correct sentences are produced with a synonymous template whose relation tokens are replaced by synonyms (e.g., *causes* can also be translated to *leads to*), while the wrong sentences with a random template whose relation tokens are randomly replaced by another one. Besides, we use **training** template referring to the templates that are used during post-training on knowledge bases. Then, we regard the sentence with lower perplexity as more reasonable. We calculate the accuracy of relation ranking as the percentage of cases where the sentence with the wrong template has the highest perplexity compared with the sentences with correct and training templates. Furthermore, we also conducted an automatic pairwise comparison to distinguish the reasonable sentences from unreasonable ones based on the perplexity scores of different models.

As shown in Table 8, the external knowledge can help our language model distinguish false sentences from true ones with higher accuracy than GPT-2 (Random chance scores 33.3%). Furthermore, our model prefers the correct template compared with the wrong one (winning rate of 71.91%), and has a close preference between the training and correct templates (winning rate of 55.76%). By contrast, GPT-2 without any external knowledge relies more on frequency to score relations, and thus can hardly tell the differences among different types

Context:	[MALE] was driving around in the snow .
Golden Story:	He needed to get home from work. He was driving slowly to avoid accidents . Unfortunately the roads were too slick and [MALE] lost control . His tires lost traction and he hit a tree .
ConvS2S:	He did n't feel like driving on his way out of the street. He put on his gloves and drove away. [MALE] made a lot of <i>paperwork</i> and <i>took a taxi to the mechanic</i> . He gave [MALE] a few thousand dollars to <i>help him get up</i> .
Fusion:	He was driving for a bit but was getting icy. [MALE] ended up driving over a <i>cliff</i> . [MALE] <i>was in the middle of the road</i> . He was out of the snow by the time he got home.
Plan&write:	He noticed a man sitting on the side of the road. He went to the nearest gas station. The man was n't sure if he could <i>have the money</i> . The man was <i>left a big deal</i> .
SKRL:	A drop of battery stacked in the quarter. [MALE] went to check the phone and turned back. He finally got a car hit off from his way to <i>play again</i> . [MALE] had forgotten his <i>hair</i> .
DSRL:	A man came in a cop. He called for an ambulance. He said he might go to a <i>park</i> . He finally showed up to <i>post what was wrong</i> .
GPT-2 (Scratch):	He noticed a large snow pile. He <i>got out</i> and <i>began driving</i> . Then he sped off and went home. [MALE] was happy that he had <i>saved his money</i> .
GPT-2 (Pretrain):	The doctors <i>took him to a hospital</i> . They told him he was going to be fine. He got all the doctors and nurses. The [MALE] parents said <i>they have no hope</i> .
GPT-2 (Fine-tune):	He noticed a snowglobe blocking his road. [NEUTRAL] he <i>backed up to the car</i> and got a snow shovel. [MALE] used to <i>shovel the snow</i> and <i>shovel the snow</i> . Now the snow is gone <i>but he has snow boots and snow gloves</i> .
Ours:	Suddenly his car broke down on the side of the road. [MALE] had to call a tow truck . The tow truck came and took [MALE] home . [MALE] was happy he was able to get home .
w/o Pretrain:	He tried to stop to see how he could get out of his car. His car started to sputter down on him. [MALE] <i>pulled over and kept driving</i> . He was able to make up and the car sped away.
w/o Knowledge:	He noticed a young lady running off. He stopped and picked her up. When he checked his <i>id</i> it was lost. [MALE] realized he had forgotten his <i>id</i> .
w/o Multi-task:	He noticed a car in the road. He decided to stop. He <i>got out of his car</i> . He <i>drove for half an hour</i> .

Table 9: Generated stories from different models. **Bold** words denote the **key** entities/events in the story. And *italic* words denote the *improper* entities/events in terms of logic and coherence in the context while the underlined words are the proper ones.

of sentences. The results indicate that our model can capture not only simple relevance between head and tail entities/events, but also the specific causal relations.

4.7 Case Study

We presented some generated examples in Table 9. Our model can generate more natural and reasonable stories than baselines.

As illustrated, the baselines (from ConvS2s to DSRL) predict wrong entities and events that are irrelevant to the leading context (e.g., *paperwork*), thereby leading to bad overall coherence in the generated stories. Pretrained GPT-2 without any fine-tuning generates an entirely irrelevant, unreasonable story (e.g., *hospital, doctor*) due to the lack of knowledge. GPT-2 trained from scratch and fine-tuned GPT-2 suffer from conflicting logic (e.g., *first got out and then began driving, and backed up to the car when driving*),

repetition (e.g., *shovel the snow*), and poor coherence with some irrelevant keywords (e.g., *save money*). In comparison, the story by our model is coherent in logic and fluent in grammar. Furthermore, without pretraining, our model can still incorporate external knowledge to generate a story with an understandable main idea but not always reasonable locally (e.g., *pulled over and kept driving*). When removing knowledge out of our full model, some confusing entities (e.g., *id*) will be generated. Additionally, removing multi-task learning also significantly affects the logic of generated stories (e.g., *first got out and then drove*) due to the inability of capturing the causal and temporal dependencies in context.

In order to verify the ability of our model to incorporate external knowledge when generating stories, we showed the utilized commonsense knowledge of this example in Figure 5. We can observe that the external knowledge is useful

<p>[MALE] was <u>driving</u> around in the <u>snow</u>.</p> <p>Suddenly his car broke down on the side of the road.</p> <p>[MALE] had to call a tow truck.</p> <p>The tow truck came and took [MALE] home.</p> <p>[MALE] was happy he was able to get home.</p>	<p>(car, UsedFor, <u>drive</u>)</p> <p>(<u>drive</u>, HasPrerequisite, car)</p> <p>(<u>snow</u>, HasProperty, slippery to <u>drive</u> on)</p> <p>(<u>drive</u>, HasSubevent, something break down)</p> <p>(PersonX calls a tow truck, <i>xNeed</i>, have his car break down)</p> <p>(PersonX asks to come, <i>xNeed</i>, call)</p> <p>(PersonX takes ___ to get home, <i>xWant</i>, go home)</p> <p>.....</p>	<p>Car is used for <u>drive</u>.</p> <p><u>Drive</u> has prerequisite of car.</p> <p><u>Snow</u> has property slippery to <u>drive</u> on.</p> <p><u>Drive</u> has subevent something break down.</p> <p>[MALE] calls a tow truck, [MALE] <i>needs</i> to have his car break down.</p> <p>[MALE] asks to come, [MALE] <i>needs</i> to call.</p> <p>[MALE] takes ___ to get home, [MALE] <i>wants</i> to go home.</p> <p>.....</p>
--	---	--

Figure 5: An example illustrating how commonsense knowledge facilitates generating reasonable stories. The right block demonstrates interrelated knowledge for the generated story, and the corresponding transformed sentences used in the training. The knowledge is retrieved from ConceptNet and ATOMIC according to the keywords denoted in **bold** in the generated story. The underlined words represent the keywords in the leading context, while the *italic* words represent the relations.

for expanding a reasonable story plot such as driving, broke down, call, came and took home, and get home.

5 Error Analysis

Although the proposed model outperforms the state-of-the-art baselines, it needs to be noted that there are still many unreasonable stories losing to other models in manual evaluation. Therefore, we analyzed error types by manually checking all *lost* stories in pairwise comparisons between our model and two strong baselines including Fusion and GPT-2 (Fine-tune) to reveal the factors that affect the performance. The numbers of stories which lost to our model in logic are 114/102 of 200/200 in total for Fusion/GPT-2 (Fine-tune) respectively. And there are 111 stories of 400 generated by our model losing to these two baselines in logic.

We manually annotated four types of error from the lost stories: **repetition** (repeating the same scenes), **unrelated** entities or events (with some wrong keywords but a reasonable main plot), **conflicting** logic (wrong causal relation or temporal order), and **chaotic** scenes (difficult to understand). The distribution of different error types is shown in Table 10. We can observe that unrelated entities/events and conflicting orders make up most of the errors for all the models. Compared with Fusion, GPT-2 (Fine-tune) reduces chaotic scenes effectively but still suffers from severe repetition. Equipped with external knowledge and multi-task learning, our model can further reduce chaotic logic and meanwhile avoid repetition. However, the analysis result illustrates that generating a coherent and reasonable story is challenging.

Error Type	Ours	Fusion	GPT-2 (Fine-tune)
Repetition (%)	1.75	5.50	6.50
Unrelated (%)	11.25	16.00	15.50
Conflicting (%)	13.75	22.00	24.50
Chaotic (%)	1.00	13.50	4.50

Table 10: Distribution of error types for different models.

Error Type	Cases
Repetition	[MALE] made up his mind to join the army . <i>He was determined to get into the army</i> . He had never been away from home. <i>He was determined to get into the army</i> . He was sent out to Afghanistan.
Unrelated	[MALE] felt he was getting sick . He had to go to an emergency room. It was his first major surgery. He had a terrible stomach ache. He was nervous about a <i>test</i> in an hour.
Conflicting	[FEMALE] swept and mopped the floor . She put her clothes in the washing machine. She was ready to go to bed. When she was <i>done</i> , she <i>washed the clothes</i> . She went to bed.
Chaotic	[MALE] was on thin ice with his job . He had a friend over to help him. [MALE] was able to <i>hold his breath</i> the entire time. he was <i>so cold that he froze</i> in his tracks. [MALE] finally <i>felt good</i> about himself.

Table 11: Typical errors by our model. **Bold** sentences are the leading context. *Italic* words denote the improper entities/events in terms of logic and coherence in the context.

We also present some typical cases by our model for each error type in Table 11. These cases show that our model still does not completely prevent logical errors including sentence-level repetition (get into the army), unrelated entities to the context (*test* is obviously unrelated to surgery and stomach ache), conflicting events (first *done* but then *washed the clothes*), and chaotic logic (due to lack

of knowledge about on thin ice). These errors also indicate external knowledge, causal relationships, and temporal dependencies play a central role in commonsense story generation.

6 Conclusions and Future Work

We present a knowledge-enhanced pretraining model with multi-task learning for commonsense story generation. The proposed framework leverages the implicit knowledge from deep pretrained language models as well as the explicit knowledge by post-training on external commonsense knowledge bases, which leads to better performance for commonsense story generation. Besides, in order to further capture the causal and temporal dependencies between the sentences in a story, we employ an auxiliary classification task to distinguish true and auto-constructed fake stories. Extensive experiments show that the proposed method can outperform strong baselines. Further analysis demonstrates that the generated stories are more coherent and reasonable thanks to the use of commonsense knowledge and multi-task learning.

As future work, it would be very interesting to make generative pretraining models have commonsense knowledge without any fine-tuning, namely, integrating the knowledge at the pretraining stage.

Acknowledgments

This work was supported by the National Science Foundation of China (grant no. 61936010/61876096) and the National Key R&D Program of China (grant no. 2018YFC0830200). We would like to thank THUNUS NExT Joint-Lab for the support. We would also like to thank our action editor, Noah Smith, and the anonymous reviewers for their invaluable suggestions and feedback.

References

Christoph Alt, Marc Hübner, and Leonhard Hennig. 2019. Fine-tuning pre-trained transformer language models to distantly supervised relation extraction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1388–1398, Florence, Italy. Association for Computational Linguistics.

Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, 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*, pages 4762–4779, Florence, Italy. Association for Computational Linguistics.

Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *EMNLP*, pages 632–642.

Elizabeth Clark, Yangfeng Ji, and Noah A. Smith. 2018. Neural text generation in stories using entity representations as context. In *NAACL*, pages 1631–1640.

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.

Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. 2019. Unified language model pre-training for natural language understanding and generation. *arXiv preprint arXiv:1905.03197v3*.

Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898.

Angela Fan, Mike Lewis, and Yann Dauphin. 2019. Strategies for structuring story generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2650–2660, Florence, Italy. Association for Computational Linguistics.

Joseph L. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76(5):378.

- Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N. Dauphin. 2017. Convolutional sequence to sequence learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1243–1252. JMLR. org.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press.
- Jian Guan, Yansen Wang, and Minlie Huang. 2019. Story ending generation with incremental encoding and commonsense knowledge. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6473–6480.
- Ari Holtzman, Jan Buys, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751v1*.
- Ting-Hao Kenneth Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, Zitnick C. Lawrence, Parikh Devi, Vanderwende Lucy, Galley, Michel, and Mitchell Margaret. 2016. Visual storytelling. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1233–1239.
- Parag Jain, Priyanka Agrawal, Abhijit Mishra, Mohak Sukhwani, Anirban Laha, and Karthik Sankaranarayanan. 2017. Story generation from sequence of independent short descriptions. *arXiv preprint arXiv:1707.05501v2*.
- Yangfeng Ji, Chenhao Tan, Sebastian Martschat, Yejin Choi, and Noah A. Smith. 2017. Dynamic entity representations in neural language models. In *EMNLP*, pages 1830–1839.
- Ryan Kiros, Yukun Zhu, Ruslan R. Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Skip-thought vectors, In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 3294–3302. Curran Associates, Inc.
- Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In *International Conference on Machine Learning*, pages 1188–1196.
- Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. 2017. Zero-shot relation extraction via reading comprehension. In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, pages 333–342.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016a. A diversity-promoting objective function for neural conversation models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119.
- Xiang Li, Aynaz Taheri, Lifu Tu, and Kevin Gimpel. 2016b. Commonsense knowledge base completion. In *Proceedings of ACL*.
- Peter LoBue and Alexander Yates. 2011. Types of common-sense knowledge needed for recognizing textual entailment. In *ACL, HLT’11*, pages 329–334. Stroudsburg, PA, USA.
- Lara Martin, Prithviraj Ammanabrolu, Xinyu Wang, William Hancock, Shruti Singh, Brent Harrison, and Mark Riedl. 2018. Event representations for automated story generation with deep neural nets. In *AAAI*, pages 868–875.
- Shikib Mehri, Evgeniia Razumovskaia, Tiancheng Zhao, and Maxine Eskenazi. 2019. Pretraining methods for dialog context representation learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3836–3845, Florence, Italy. Association for Computational Linguistics.
- Todor Mihaylov and Anette Frank. 2018. Knowledgeable reader: Enhancing cloze-style reading comprehension with external commonsense knowledge. In *ACL*, pages 821–832.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781v3*.
- Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James

- Allen. 2016a. A corpus and cloze evaluation for deeper understanding of commonsense stories. In *NAACL*, pages 839–849.
- Nasrin Mostafazadeh, Lucy Vanderwende, Wen Tau Yih, Pushmeet Kohli, and James Allen. 2016b. Story cloze evaluator: Vector space representation evaluation by predicting what happens next. In *The Workshop on Evaluating Vector-Space Representations for NLP*, pages 24–29.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual meeting of the Association for Computational Linguistics*, pages 311–318. Association for Computational Linguistics.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of NAACL-HLT*, pages 2227–2237.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8).
- Hannah Rashkin, Antoine Bosselut, Maarten Sap, Kevin Knight, and Yejin Choi. 2018a. Modeling naive psychology of characters in simple commonsense stories. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2289–2299.
- Hannah Rashkin, Maarten Sap, Emily Allaway, Noah A. Smith, and Yejin Choi. 2018b. Event2mind: Commonsense inference on events, intents, and reactions. In *ACL*, pages 463–473.
- Melissa Roemmele. 2016. Writing stories with help from recurrent neural networks. In *AAAI*, pages 4311–4312. Phoenix, AZ. AAAI Press.
- Stuart Rose, Dave Engel, Nick Cramer, and Wendy Cowley. 2010. Automatic keyword extraction from individual documents. *Text Mining: Applications and Theory* pages 1–20.
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A. Smith, and Yejin Choi. 2019. Atomic: an atlas of machine commonsense for if-then reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3027–3035.
- Abigail See, Aneesh Pappu, Rohun Saxena, Akhila Yerukola, and Christopher D. Manning. 2019. Do massively pretrained language models make better storytellers? In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 843–861.
- Zhihong Shao, Minlie Huang, Jiangtao Wen, Wenfei Xu, and Xiaoyan Zhu. 2019. Long and diverse text generation with planning-based hierarchical variational model. *EMNLP 2019*.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2019. Mass: Masked sequence to sequence pre-training for language generation. In *International Conference on Machine Learning*, pages 5926–5936.
- R. Speer and C. Havasi. 2012. Representing general relational knowledge in conceptnet 5. In *LREC*, pages 3679–3686.
- Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, and Hua Wu. 2019. Ernie: Enhanced representation through knowledge integration. *arXiv preprint arXiv:1904.09223v1*.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *NIPS*, pages 3104–3112.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.

- Zhe Wang, Wei He, Hua Wu, Haiyang Wu, Wei Li, Haifeng Wang, and Enhong Chen. 2016. Chinese poetry generation with planning based neural network. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1051–1060.
- Thomas Wolf, Victor Sanh, Julien Chaumond, and Clement Delangue. 2019. Transfertransfo: A transfer learning approach for neural network based conversational agents. *arXiv preprint arXiv:1901.08149v2*.
- Jingjing Xu, Xuancheng Ren, Yi Zhang, Qi Zeng, Xiaoyan Cai, and Xu Sun. 2018. A skeleton-based model for promoting coherence among sentences in narrative story generation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4306–4315.
- An Yang, Quan Wang, Jing Liu, Kai Liu, Yajuan Lyu, Hua Wu, Qiaoqiao She, and Sujian Li. 2019a. Enhancing pre-trained language representations with rich knowledge for machine reading comprehension. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2346–2357. Florence, Italy. Association for Computational Linguistics.
- Pengcheng Yang, Lei Li, Fuli Luo, Tianyu Liu, and Xu Sun. 2019b. Enhancing topic-to-essay generation with external commonsense knowledge. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2002–2012. Florence, Italy. Association for Computational Linguistics.
- Lili Yao, Nanyun Peng, Ralph Weischedel, Kevin Knight, Dongyan Zhao, and Rui Yan. 2019. Plan-and-write: Towards better automatic storytelling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7378–7385.
- Jianfei Yu and Jing Jiang. 2016. Learning sentence embeddings with auxiliary tasks for cross-domain sentiment classification. In *Proceedings of 2016 Conference on Empirical Methods in Natural Language Processing*, pages 236–246.
- Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. ERNIE: Enhanced language representation with informative entities. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1441–1451. Florence, Italy. Association for Computational Linguistics.
- Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi. 2017. Learning discourse-level diversity for neural dialog models using conditional variational autoencoders. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 654–664.
- Hao Zhou, Tom Yang, Minlie Huang, Haizhou Zhao, Jingfang Xu, and Xiaoyan Zhu. 2018. Commonsense knowledge aware conversation generation with graph attention. In *IJCAI*.