

Improving Neural Story Generation by Targeted Common Sense Grounding

Huanru Henry Mao, Bodhisattwa Prasad Majumder,
Julian McAuley, Garrison W. Cottrell

Department of Computer Science and Engineering
UC San Diego

{hhmao, bmajumde, jmcauley, gary}@eng.ucsd.edu

Abstract

Stories generated with neural language models have shown promise in grammatical and stylistic consistency. However, the generated stories are still lacking in *common sense reasoning*, e.g., they often contain sentences deprived of world knowledge. We propose a simple multi-task learning scheme to achieve quantitatively better common sense reasoning in language models by leveraging auxiliary training signals from datasets designed to provide common sense grounding. When combined with our two-stage fine-tuning pipeline, our method achieves improved common sense reasoning and state-of-the-art perplexity on the *Writing Prompts* (Fan et al., 2018) story generation dataset.

1 Introduction

Story generation is the task of automatically producing compelling creative writing. Recent advances in language modeling have yielded thematic and stylistic coherence in story generation through large scale pretraining of Transformer models (Vaswani et al., 2017). The recent introduction of the General Pre-trained Transformer v2 (GPT2) (Radford et al., 2019)—a high-capacity Transformer trained on a large, diverse corpus of text crawled from the web (called WebText)—is capable of generating stylistically coherent text but commonly produces text with logical inconsistencies. For example, in one sample the model writes: “*It was a sunny, warm summer night*”. Obviously, this writing is nonsense as it cannot be sunny at night. The lack of common sense reasoning in such a strong language model suggests that minimizing next-token perplexity alone may be insufficient in producing models that can compose sensible stories.

In this paper, we consider the challenge of *common sense reasoning* (CSR) in language model-

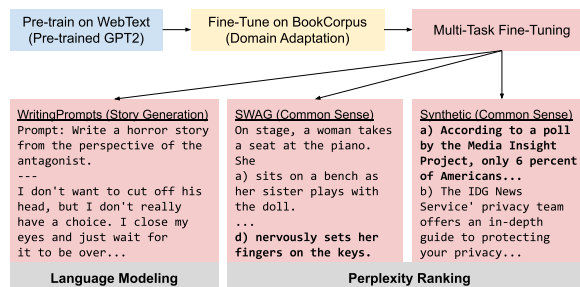


Figure 1: Our two-stage fine-tuning pipeline with auxiliary multi-task learning. For perplexity ranking examples, bolded text indicates the correct answer. For **synthetic**, the correct choice is written by a human, and the wrong choice is generated by a neural network.

ing for story generation. Unlike other work in the CSR literature (Storks et al., 2019), which evaluates CSR in isolation, we are specifically interested in a generative model’s likelihood of producing text that exhibits common sense. Evaluating common sense qualitatively in a model’s samples is difficult, as it is subject to human bias and dependent on sampling procedure. We propose evaluating the common sense of a model automatically by ranking the model’s perplexity on spurious (plausible, but nonsense) text completions from SWAG (Zellers et al., 2018) and Story Cloze (Mostafazadeh et al., 2016) datasets, which are designed for common sense grounding.

Our contributions are as follows: We propose a simple way to define *better* CSR in generative models, which leads to an auxiliary multi-task objective to directly bias our model to generate text with better common sense. When fine-tuning is combined with multi-task learning in a two-stage pipeline, we improve the model’s CSR and outperform state-of-the-art perplexity on the Writing Prompts (Fan et al., 2018) dataset.¹

¹Our source code is at <https://github.com/calclavia/story-generation>.

2 Tasks

2.1 Primary Task: Language Modeling

Our primary task is to perform language modeling (Elman, 1990; Bengio et al., 2003; Dai and Le, 2015) on the WritingPrompts dataset. A language model learns to assign the probability of a text sequence $X = x_1, \dots, x_T$ using the conditional probability factorization:

$$P(X) = \prod_{t=1}^T P(x_t | x_{1:t-1}). \quad (1)$$

We train our model using a standard cross-entropy loss between next-step true tokens and predicted probabilities given current tokens.

WritingPrompts (Fan et al., 2018) is a dataset of prompts and short stories crawled from Reddit. The aim of the dataset is to produce a story given a free-text prompt. We reduce this conditional text generation task into a generic language modeling task by simply concatenating the prompt before the story and treating a prompt-story pair as one input to the Transformer decoder model. This human-readable format (example in Figure 1) is chosen because GPT2 may have been trained on similarly formatted text from the web. When sampling, we can either seed the model with a prompt or allow it to generate its own prompt.

2.2 Auxiliary Task: Perplexity Ranking

Our auxiliary task aims to bias the model to produce text with better common sense (which we refer to as *sensible* text). Given a set of text sequences $S = \{S_1, \dots, S_N\}$, where S_1 is a sensible text sequence and the rest S_2, \dots, S_N are spurious text sequences, we operationally define *better* as the model assigning higher probability (Eq. 1) to S_1 versus the average probability over spurious text sequences S_2, \dots, S_N . Using this definition, it is possible to directly optimize the model to assign $P(S_1)$ to be higher than any $P(S_i)$. Formally, we define the probability of the model choosing the correct sequence S_1 over spurious sequences as the softmax over the length-normalized log probabilities of all plausible sequences:

$$\frac{\exp(\frac{1}{T_1} \log P(S_1))}{\sum_{i=1}^N \exp(\frac{1}{T_i} \log P(S_i))} \quad (2)$$

where T_i refers to the length of the text sequence S_i . Thus, we can practically minimize the negative log-likelihood of Eq. 2 by reusing the same

Dataset	Size	Role
WritingPrompts	272K Stories	Story Generation
BookCorpus	10K Books	Domain Adaptation
SWAG	73K Questions	Common Sense
Synthetic	250K Pairs	Common Sense

Table 1: Datasets, training set sizes and their roles

softmax layer used for the primary language modeling task. We refer to this objective as *perplexity ranking* as it constrains the model to rank sensible text to have lower perplexity than spurious ones.

SWAG: In order to train on this auxiliary objective, we need training examples in the format of multiple choice questions, where the correct choice corresponds to the text with the best common sense. We choose the SWAG dataset (Zellers et al., 2018), a dataset that provides 4-way multiple choice common sense questions that are adversarially filtered to seem plausible to language models. Unlike other CSR datasets (Talmor et al., 2018), SWAG forms its question and answer as two full sentences, which we can concatenate into a single string to find its probability. This makes it suitable for perplexity ranking.

Synthetic: To facilitate perplexity ranking on SWAG, we additionally use a synthetic dataset that consists of 250K pairs of human written text from WebText and samples generated by the original 1.5B parameter version of the GPT2 model.² These samples are many paragraphs long and truncated to a maximum of 1024 tokens. We frame these pairs as a 2-way classification problem and train the model by perplexity ranking to assign higher likelihood to human written text over synthetic examples. The assumption we make is that human written text is more sensible than text written by neural language models. Hence, on average, this promotes the model to assign higher probabilities to sensible text.

3 Training Pipeline

We introduce a two-stage training pipeline (Figure 1) to improve model performance both in terms of perplexity and CSR on story generation. Our pipeline uses four different datasets (Table 1), each of which plays a role in improving model performance. Our model architecture is the 117M parameter version of GPT2, using the pre-trained

²We obtained the samples from <https://github.com/openai/gpt-2-output-dataset>

weights provided by Radford et al. (2019). We refer readers to Vaswani et al. (2017) for details of the Transformer architecture.

We first perform intermediate fine-tuning (Phang et al., 2018; Howard and Ruder, 2018) of the pre-trained GPT2 on BookCorpus (Zhu et al., 2015) as a method of domain adaptation from WebText to the domain of stories. BookCorpus is a dataset that contains over 10,000 free books crawled from the web.³ We train on this corpus using our language modeling objective. Next, we fine-tune on the target WritingPrompts dataset with a multi-task learning objective. We alternate training between the language modeling objective on WritingPrompts and perplexity ranking on SWAG and our synthetic dataset. Training details and hyperparameters are in the appendix.

4 Evaluation

We perform three types of evaluation on the model to assess its readability, reliance on the prompt (prompt ranking) and CSR.

Readability is measured in terms of model perplexity on the test set of WritingPrompts. Because GPT2 uses subword tokenization (Sennrich et al., 2016), it is not directly comparable to the word-level perplexity obtained in Fan et al. (2018). We estimate the corresponding word-level perplexity by taking the product of each subword’s probabilities to obtain probabilities for each word. Both sub-word perplexity and word-level perplexities are reported in our experiments.

Prompt ranking (Fan et al., 2018) assesses how well a model matches a story to its given prompt. This is measured by computing the likelihood of stories conditioned under ten different prompts, nine of which are randomly sampled and one is the true prompt. Following Fan et al. (2018), we count a random story sample as correct when it ranks the true prompt with the lowest perplexity. We compute the accuracy from 1000 random samples.

CSR is evaluated on two multiple choice datasets – SWAG and Story Cloze (Mostafazadeh et al., 2016). We rank the perplexity computed by the model for each example and count it as correct if the lowest perplexity matches the answer. The SWAG validation set provides a proxy of how well the model generalizes in CSR to unseen examples. To ensure generalization beyond SWAG, we also

³We crawled BookCorpus from <https://www.smashwords.com/>

Premise: *John and Billy became very skilled at beer pong. They entered a contest in college. They won the contest and advanced to the next level. The next level sent them to Vegas.*

GPT2 → BC → WP output:

1. They would fall.
2. Later, they figured out what it was all about.

GPT2 → BC → WP + SWAG + SYNTH output:

1. They have been ranked number one in their respective leagues and are considered the best in their respective countries.
 2. They then moved to the bars.
-

Table 2: Top two highest likelihood story completions from 10 random completion samples generated by our models when primed with a premise from the Story Cloze validation set.

perform *zero-shot* evaluation on the Winter 2018 Story Cloze validation set. Story Cloze consists of 5-sentence stories with correct and spurious endings. It is similar to SWAG but serves as an in-domain dataset to specifically test the model’s performance at CSR in story telling.

5 Results and Discussion

We analyze our pipeline and report the results in Table 3. We also generate stories by sampling from our model using nucleus sampling with $p = 0.9$ (Holtzman et al., 2019). We present example story completions in Table 2 and full sampled stories in our appendix and Table 4.

Pre-training: We compare our models with the attention-based Fusion Model (Fan et al., 2018), which has been designed for and trained on WritingPrompts. We observe that a pre-trained GPT2 performing zero-shot inference on WritingPrompts (GPT2 in Table 3) is a strong baseline. By fine-tuning GPT2 on WritingPrompts (GPT2 → WP), we outperform the Fusion Model in perplexity. All our models outperform the Fusion Model in prompt ranking, which suggests that task-specific models are unnecessary given pre-training.

Intermediate Fine-Tuning: The first stage in our pipeline performs intermediate fine-tuning of GPT2 on BookCorpus (GPT2 → BC in Table 3). To confirm that intermediate fine-tuning helps downstream performance, we evaluate the zero-shot performance of the model on WritingPrompts. This yields perplexity and prompt ranking improvements compared to GPT2, demonstrating successful domain adaptation. Perform-

Models	SW PPL	Word PPL	Prompt Ranking	SWAG	Story Cloze
Fusion Model (Fan et al., 2018)	-	36.6	16.3%	-	-
GPT2	35.57	51.29*	49.8%	48.1%	58.8%
GPT2 → BC	29.10	42.01*	62.7%	50.5%	59.6%
GPT2 → WP	21.68	30.65*	80.0%	49.8%	58.3%
GPT2 → BC → WP	20.79	29.56*	80.6%	51.4%	59.1%
GPT2 → BC → WP + SWAG	20.78	29.52*	78.9%	75.3%	63.2%
GPT2 → BC → WP + SWAG + SYNTH	20.78	29.63*	80.1%	76.3%	64.1%

Table 3: Performance of models on the test set of WritingPrompts and validation set of SWAG and Story Cloze. **SW PPL** and **Word PPL** refer to sub-word and word-level perplexity on WritingPrompts respectively. WP refers to WritingPrompts, BC refers to BookCorpus and SYNTH refers to training with our 250K synthetic examples. The asterisk * refers to an estimated score derived from BPE PPL.

ing two-stage fine tuning (GPT2 → BC → WP) further improves perplexity and CSR. We hypothesize the improvement in CSR is due to BookCorpus being a higher quality dataset written by authors when compared against WebText.

Multi-tasking Fine-Tuning: Performing multi-task learning on WritingPrompts and SWAG (GPT2 → BC → WP + SWAG) unsurprisingly yields significant improvements on the SWAG validation set. More importantly, the zero-shot performance on Story Cloze also improved, indicating that it was able to generalize its common sense knowledge. We also see qualitatively improved results when generating story completions (Table 2). The addition of the synthetic dataset we introduced (GPT2 → BC → WP + SWAG + SYNTH) further boosts performance on CSR. Other metrics are negligibly affected by the auxiliary tasks.

6 Related Work

Story Generation: Recent work in neural story generation (Kiros et al., 2015; Roemmele, 2016) has shown success in using hierarchical methods (Yao et al., 2018; Fan et al., 2018) to generate stories. In these schemes, a neural architecture is engineered to first generate an outline or a prompt, then to expand the prompt into a full-length story. Our work performs hierarchical generation, but our main focus is on achieving better common sense in the generated text rather than engineering task-specific architectures.

Common Sense Reasoning: Common sense reasoning (CSR) has been studied through many benchmarks such as SWAG (Zellers et al., 2018), Story Cloze (Mostafazadeh et al., 2016), the Winograd Schema Challenge (Levesque et al., 2012), and CommonsenseQA (Talmor et al., 2018). Recent methods (Peters et al., 2018; Rad-

ford et al., 2018) on these benchmarks focus on large-scale pre-training of language models. They show that transfer learning is an effective means to improve CSR and our fine-tuning pipeline builds upon these techniques. Our results on SWAG and Story Cloze are far from state-of-the-art (Devlin et al., 2018). However, our aim is not to directly tackle SWAG or Story Cloze, but instead to use it as a constraint on our model and a proxy to estimate the likelihood of generating sensible text.

Multi-task Learning: Multi-task learning (MTL) introduces inductive bias in a model, helps reduce overfitting and increases robustness (Caruana, 1993; Ruder, 2017; McCann et al., 2018). Our work builds upon MTL principles as we introduce auxiliary tasks to specifically tackle CSR (Søgaard and Bingel, 2017). Contrary to conventional auxiliary multi-task learning (Howard and Ruder, 2018), which adds an additional classifier to the final representations of the model for learning auxiliary tasks, we use perplexity ranking, which does not introduce additional parameters. We argue that perplexity ranking, by definition, guarantees that the model will generate sensible text with a higher probability.

7 Conclusion

We demonstrated a simple way to improve the likelihood that a particular language model generates sensible text. We introduced a training pipeline that leverages transfer and multi-task learning to improve perplexity and CSR. One limitation of our approach is its dependence on SWAG, which only provides short-term common sense grounding. While our synthetic dataset may mitigate this issue, story-level coherence is unlikely to be achieved using this method alone. We leave this investigation for future work.

Prompt: Aliens start abducting humans.

Steve smashed his Xbox controller at the table. He lit a smoke and took another hit. Two exes were in the room, just outside the door.

“What are you doing?” he said. “Get back in there.”

“I-I’ve never actually seen any of the aliens,” he replied, dispelling the smoke from his mouth.

The two aliens were smirking. “Is this your first time?” the ex asked.

“The last one tried to kill us.”

The aliens looked at each other. “Are you sure you are ok? It looks like you are severely hurt.”

“Yeah, I’m fine. I just saw the door open. I’m just fine,” he replied.

“We’re fine? You do that every day, Steve, but we still haven’t gotten past this door,” the ex said, resting a hand on the UFO that stood in the doorway.

“I’m fine.”

“Good! I’ve never seen anything like it, have I Steve?”

Steve stood up. “Yes I have.” He picked up his controller.

“Ok, shoot!” the aliens yelled. “This is no time to hesitate.”

Steve put the controller to his ear. “Who are you? What are you doing here?”

The aliens pulled the gun towards him. “We-we’re here to kidnap you.”

Steve put his hands on the UFO, and pulled the trigger. The gun went off. The aliens, scared, immediately realized where they were and ran towards the other two aliens. They kept running. Steve continued on his way towards his childhood toy house. When he was just a teen, he had built his own helicopter. He always felt they had lost. The abandoned building on the edge of the highway he and the ex had entered had been the perfect landing spot. He didn’t know where they were. He didn’t know how they were going to get out of there. There were bodies everywhere.

Table 4: Sample generated by GPT2 → BC → WP + SWAG + SYNTH primed with the same prompt as Fan et al. (2018).

Acknowledgements. This work is partly supported by NSF #1750063, the Guangzhou Science and Technology Planning Project #201704030051, NSF awards CNS-1730158, ACI-1540112, ACI-1541349, OAC-1826967, the University of California Office of the President, and the University of California San Diego’s California Institute for Telecommunications and Information Technology/Qualcomm Institute. We thank CENIC for the 100Gpbs networks and all the reviewers for their constructive suggestions.

References

- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. 2003. [A neural probabilistic language model](#). *Journal of Machine Learning Research*, 3:1137–1155.
- Rich Caruana. 1993. Multitask learning: A knowledge-based source of inductive bias. In *ICML*.
- Andrew M. Dai and Quoc V. Le. 2015. [Semi-supervised sequence learning](#). In *NIPS*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [BERT: pre-training of deep bidirectional transformers for language understanding](#). *CoRR*, abs/1810.04805.
- Jeffrey L. Elman. 1990. Finding structure in time. *Cognitive Science*, 14(2):179–212.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. [Hierarchical neural story generation](#). In *ACL*.
- Ari Holtzman, Jan Buys, Maxwell Forbes, and Yejin Choi. 2019. [The curious case of neural text degeneration](#). *CoRR*, abs/1904.09751.
- Jeremy Howard and Sebastian Ruder. 2018. [Fine-tuned language models for text classification](#). *CoRR*, abs/1801.06146.
- Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. [Skip-thought vectors](#). In *NIPS*.
- Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In *Thirteenth International Conference on the Principles of Knowledge Representation and Reasoning*.
- Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2018. The natural language de-

- cathlon: Multitask learning as question answering. *CoRR arXiv:1806.08730*.
- Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James F. Allen. 2016. A corpus and evaluation framework for deeper understanding of commonsense stories. *NAACL-HLT*.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *NAACL-HLT*.
- Jason Phang, Thibault Févry, and Samuel R. Bowman. 2018. Sentence encoders on stilts: Supplementary training on intermediate labeled-data tasks. *CoRR*, abs/1811.01088.
- A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. 2019. Language models are unsupervised multitask learners. Technical report, OpenAI.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. Technical report, OpenAI.
- Melissa Roemmele. 2016. Writing stories with help from recurrent neural networks. In *AAAI*.
- Sebastian Ruder. 2017. An overview of multi-task learning in deep neural networks. *CoRR*, abs/1706.05098.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *ACL*.
- Anders Søgaard and Joachim Bingel. 2017. Identifying beneficial task relations for multi-task learning in deep neural networks. In *EACL*.
- Shane Storks, Qiaozi Gao, and Joyce Y. Chai. 2019. Commonsense reasoning for natural language understanding: A survey of benchmarks, resources, and approaches. *CoRR*, abs/1904.01172.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2018. Commonsenseqa: A question answering challenge targeting commonsense knowledge. *CoRR*, abs/1811.00937.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*.
- Lili Yao, Nanyun Peng, Ralph M. Weischedel, Kevin Knight, Dongyan Zhao, and Rui Yan. 2018. Plan-and-write: Towards better automatic storytelling. *CoRR*, abs/1811.05701.
- Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. SWAG: A large-scale adversarial dataset for grounded commonsense inference. *CoRR*, abs/1808.05326.
- Yukun Zhu, Ryan Kiros, Richard S. Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *ICCV*.