

# Guiding Neural Story Generation with Reader Models

Anonymous ACL submission

## Abstract

Automated storytelling has long captured the attention of researchers for the ubiquity of narratives in everyday life. However, it is challenging to maintain coherence and stay on-topic toward a specific ending when generating narratives with neural language models. In this paper, we introduce Story generation with Reader Models (StoRM), a framework in which a *reader model* is used to reason about the story should progress. A reader model infers what a human reader believes about the concepts, entities, and relations about the fictional story world. We show how an explicit reader model represented as a knowledge graph affords story coherence and provides controllability in the form of achieving a given story world state goal. Experiments show that our model produces significantly more coherent and on-topic stories, outperforming baselines in dimensions including plot plausibility and staying on topic. Our system also outperforms outline-guided story generation baselines in composing given concepts without ordering.

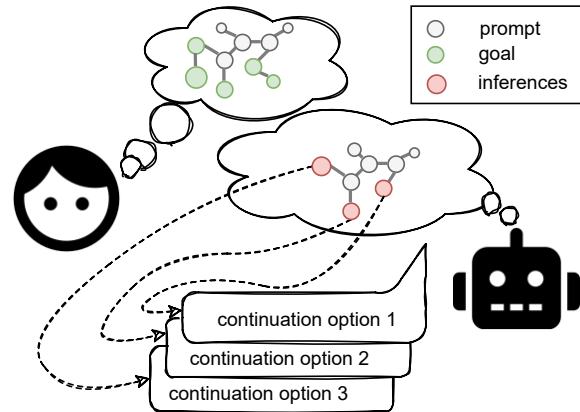


Figure 1: The overview of StoRM system. Our goal is to build a story world  $\text{⦿}$  covering all the given concepts  $\text{⦿}$ . 1. A set of concepts  $\text{⦿}$  and a prompt  $\text{⦿}$  starts the story generation process. 2. The system builds a goal story world  $\text{⦿}$  and prompt story world  $\text{⦿}$  with knowledge graph. 3. The system infers a graph of concepts  $\text{⦿}$ . 4. A language model generates continuation options on inferences. 5. *Topk* continuations minimize its difference with goal story world are added to story.

## 1 Introduction

Automated Story Generation is the challenge of designing an artificial intelligence system that can generate a natural language text that is perceived by readers as a story. Early work on story generation used symbolic planning (Meehan, 1976; Lebowitz, 1987; Cavazza et al., 2003; Porteous and Cavazza, 2009; Riedl and Young, 2010; Ware and Young, 2010; Ware and Siler, 2021). These systems would be provided with a description of the initial world state—usually a list of predicates—and a goal—a description of what predicates should be true to be successful. These approaches had two benefits. First, the plots tended to be coherent because of logical constraints on the actions. Second, the plots were guaranteed to end in a state in which the goal held. However, these systems require substantial knowledge engineering of logical constraints, limit-

ing their generality, and don’t always generate plot or stories in natural language.

Recently, neural language modeling approaches (Roemmele, 2016; Khalifa et al., 2017; Martin et al., 2018; Clark et al., 2018; Yao et al., 2019; Rashkin et al., 2020; Fan et al., 2019; Amanabrolu et al., 2021a) have been applied to story generation because they circumvent the need for manual knowledge engineering and tend to produce relatively fluent, varied, and naturalistic language. Language models are, however, not goal-directed. That is, one cannot natively provide both a context prompt and a goal to be achieved after an arbitrary number of continuations. Further, language models struggle with maintaining story coherence—the logical progression of events—and may also become repetitive. Large, pre-trained language models improve fluency and generalization but do not provide goal-directedness and stories generated

061 can still be perceived as lacking in coherence in the  
062 sense that they meander without direction.

063 In this paper we consider the challenge of co-  
064 herent and controllable text generation for neu-  
065 ral language model based story generation. We  
066 hypothesize that neural language models, while  
067 powerful text-completion systems, are not natively  
068 well-suited for coherent story generation because  
069 a neural network trained with a cross-entropy loss  
070 function is unlikely to model the unfolding context  
071 of a story the same way as a human reader. Stud-  
072 ies of human reader comprehension (Zwaan and  
073 Radvansky, 1998) show that readers comprehend  
074 stories by tracking the relations between entities  
075 and events in ways that can be expressed as a graph.  
076 The perceived coherence of a story is a function  
077 of the connectedness of this graph (Graesser et al.,  
078 1994). Ensuring the causality between sentences  
079 can significantly improve the coherence of stories  
080 (Peng et al., 2021).

081 Inspired by cognitive science, we aim to aug-  
082 ment neural language models with a **reader model**  
083 in which a story generation system infers a graph  
084 of concepts, entities, and relations that a reader is  
085 likely to believe about the story world as they read  
086 an incrementally generated story. The reader model  
087 enables the story generation algorithm to explicitly  
088 reason about the entities and relations and generate  
089 story continuations that use those entities to move  
090 the story forward; a reader can track how entities  
091 and relations change over time and thus perceive  
092 stories as more coherent. We use large language  
093 models to produce the continuation text of the story  
094 generation. However instead of providing the pre-  
095 vious story as context, our algorithm selects one  
096 or more entity from the world model and uses tem-  
097 plate filling to generates candidate continuations.  
098 This is a deviation from the typical way in which  
099 large language models are used to generate text  
100 continuations from a context prompt and possibly  
101 some additional external biases.

102 The reader model doesn’t guarantee controlled  
103 text generation or goal-directedness but provides  
104 a means for directing the generation process. In  
105 addition to a starting context prompt, we require a  
106 goal to be given in the form of what the graphical  
107 reader’s model of the story world should be at the  
108 end of the story. In that way, the goal provides a  
109 rough outline of the entities and relations that need  
110 to be present in the story but without providing  
111 particulars about everything that must be in the

story or the ordering in which they must occur. 112

Our contributions are as twofold: (1) we pro- 113  
pose an automated story generation model with 114  
Reader Models (**StoRM**) which maintain coher- 115  
ence and controllability of generated stories at the 116  
same time; and (2) we conduct a thorough exper- 117  
imental study against strong baselines that shows 118  
that StoRM produces significantly more coherent 119  
and goal-directed story. 120

## 2 Related Work and Background 121

We situate our paper in the literature of neural 122  
networks—recurrent and transformer-based—to 123  
produce stories (Roemmele, 2016; Khalifa et al., 124  
2017; Martin et al., 2018; Clark et al., 2018). There 125  
are a few works that are highly related to our pro- 126  
posed framework, in terms of the following two 127  
dimensions: the generation controllability and the 128  
usage of commonsense knowledge. The controlla- 129  
bility in story generation focuses on how to enable 130  
the generation process to adhere to the user’s in- 131  
puts. Fan et al. (2018) proposes a hierarchical story 132  
generation framework that divides the generation 133  
process into two levels of hierarchy: generating a 134  
writing prompt (premise) and then transforming it 135  
into a passage of text conditioned on the prompt. 136  
Similarly, Plan-And-Write (Yao et al., 2019) con- 137  
ducts generation in two steps: planning a story out- 138  
line based on a title (topic), then generating a story 139  
based on the storyline. Plot Machines (Rashkin 140  
et al., 2020) accepts as an input an un-ordered out- 141  
line of concepts and conditions a language model. 142

Commonsense knowledge plays an important 143  
role in story generation. The most popular way 144  
of utilizing it is to train neural language models 145  
(e.g. GPT-2 (Radford et al., 2019)) on common- 146  
sense knowledge bases such as ConceptNet (Speer 147  
and Havasi, 2013) and ATOMIC (Sap et al., 2019; 148  
Hwang et al., 2021) which contains detailed in- 149  
formation regarding well-known facts or causal 150  
relationships. Thus the resulting language model, 151  
named COMET (Bosselut et al., 2019; Hwang et al., 152  
2021), becomes capable of inferring new common- 153  
sense knowledge on novel phrases. Ammanabrolu 154  
et al. (2021b) proposes Causal, Commonsense Plot 155  
Ordering (C2PO) framework which takes advan- 156  
tage of COMET to infer predecessor and successor 157  
events and then bi-directionally search from pre- 158  
specified start event to end event, however, C2PO 159  
generates plots made up of highly constrained, tem- 160  
plated text; Peng et al. (2021) leverages COMET to 161

infer the character intentions and effects of actions so as to guide the generation process, but they did not consider controllability. There are also other approaches that directly incorporate commonsense knowledge graphs into the encoding process (Mihaylov and Frank, 2018; Guan et al., 2019). Their works paid more attention on improving coherence with the help of common-sense knowledge but did not take controllability into consideration.

### 3 Story Generation with Reader Models

In this section, we introduce a framework—*Story generation with Reader Models* (StoRM)—for generating stories with models of what the reader will believe about the fictional story world. We hypothesize that the incorporation of a *reader model* into the story generation process will increase story coherence. We define *story coherence* as the extent to which readers can identify connections between different events and entities in a story. In this work, the reader model is represented as a *knowledge graph*, a set of triples of the form  $\langle \text{subject}, \text{relation}, \text{object} \rangle$ . By making the beliefs about what the reader likely knows explicit, we provide mechanisms for selecting which entities to include in the continuation of the story.

Because the StoRM framework maintains a knowledge graph that approximates the reader’s beliefs about the story world, we are able to compare the reader model to a desired world state, also described as a knowledge graph. The StoRM framework is thus *controllable*—a user can provide a *story goal* in the form of a knowledge graph (*goal story world*) that describes the story world at the conclusion of the story. This allows the system to make informed decisions about which possible story continuations are likely to achieve the desired goal state by inferring how each possible story continuation changes the reader model to be closer to the desired goal story world

Our framework starts with a prompt and a description of the story outcome (See Figure 2). The prompt is transformed into a knowledge graph by extracting entities (§3.1). The entities are expanded using two commonsense techniques (§3.2). (1) ConceptNet (Speer and Havasi, 2013), a crowd-sourced knowledge base of concepts and relations, and (2) COMET<sub>20</sub><sup>20</sup> (Hwang et al., 2021), a neural network that generates commonsense inferences. The generation technique selects different entities and uses templates to generate possible story con-

tinuations (§3.3). By targeting different entities and using template infilling, we reduce neural network hallucination of new entities and create a diverse set of story continuations. Each potential continuation is scored based on how it changes the knowledge graph relative to the goal, which is also transformed into a knowledge graph. The selected continuation starts the next iteration of the generation process.

#### 3.1 Knowledge Graph Acquisition

With the automatic generation of the story, some important information could be forgotten. The knowledge graph is an explicit and *persistent* memory of entities mentioned or inferred from the story text generated so far. Knowledge Graphs represent information in the form of triples, consisting of a subject entity, relation and object entity. For example, “*Jenny lived in Florida*” is represented as  $\langle \text{jenny}, \text{live}, \text{florida} \rangle$ . The entities represent the nodes of the graph and their relations act as edges.

To acquire the knowledge graph, we firstly trained a Semantic Role Labeling (SRL) model (Gildea and Jurafsky, 2002) on VerbAtlas (Di Fabio et al., 2019)—a hand-crafted lexical-semantic resource whose goal is to bring together all verbal synsets from WordNet (Fellbaum, 1998) into semantically-coherent frames. This SRL model provides the automatic identification and labeling of argument structures of stories. Further detail can be found in Appendix A.1.

StoRM then converts the output of VerbAtlas SRL model into knowledge graph triples. Entities represent the theme and attribute and VerbAtlas frames act as edges. An example is shown in left side of Fig. 2. Multiple character names, object names and pronouns make the knowledge graph representation hard to interpret. Hence, we adopt an end-to-end Coreference Resolution model (Lee et al., 2017) to find all expressions that refer to the same entity in a story to minimize the entities.

StoRM starts with two knowledge graphs. The first,  $\mathbf{G}_1$ , is the converted prompt (first sentence). The second  $\mathbf{G}_{\text{goal}}$  is the converted goal description.  $\mathbf{G}_{\text{goal}}$  can also optionally be generated from a target story, in which case each sentence is converted to nodes and relations and incrementally added to the graph. With the generation of the continuation candidates, we will update the knowledge graph  $\mathbf{G}_t$  with new continuations to get new knowledge graph  $\mathbf{G}_{t+1}$ , where  $t$  is the index of the sentence in the story.

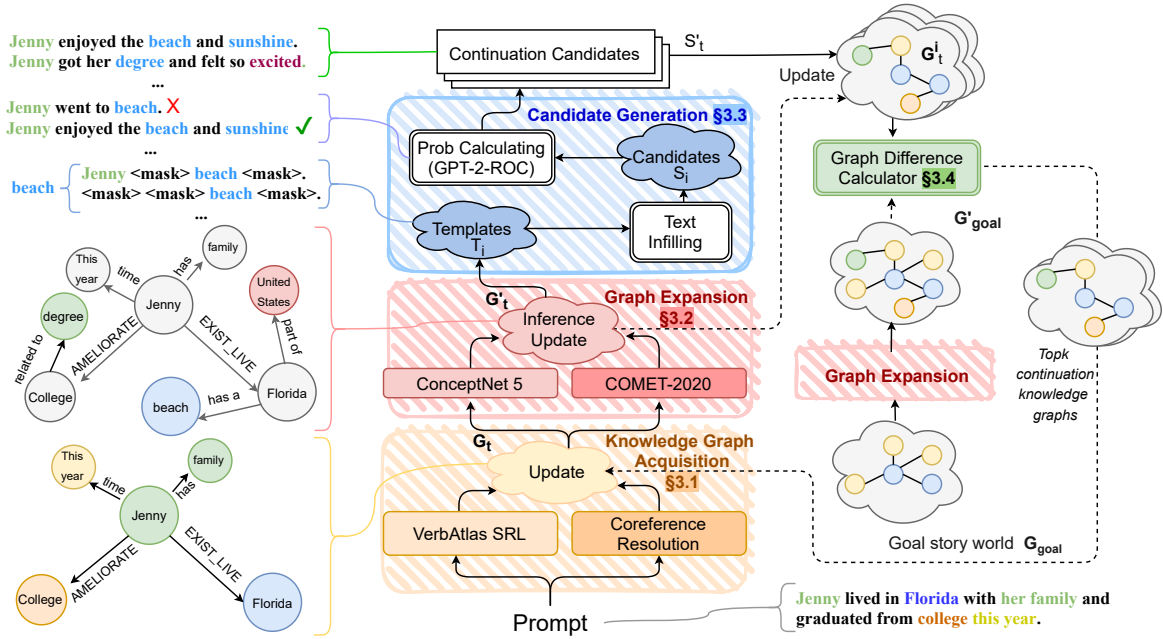


Figure 2: The overall procedure of StoRM.

### 3.2 Graph Expansion

Human readers use commonsense knowledge to infer the presence of entities and concepts not explicitly mentioned in story text. For example, Florida has beaches and eating dinner implies dishes. In accordance, we use common-sense inferences of entities to expand the knowledge graph to provide entities for characters to interact with and thus drive the story forward. Because the presence of entities and concepts are inferred from prior events, the reader should be able to track the connections between entities and events, thus supporting perceived story coherence.

We use two inference techniques. We first consider ConceptNet5 (Speer and Havasi, 2013), a multilingual knowledge base, representing words and phrases that people use and the common-sense relations between them. We expand each entity except characters nodes with the inference extracted from ConceptNet5 to get inference entity set  $E_{1,t}$ . Second, we repeat the process with COMET<sub>20</sub><sup>20</sup> (Hwang et al., 2021), which is a transformer-based generative model trained on the ATOMIC<sub>20</sub><sup>20</sup> commonsense dataset (Hwang et al., 2021) to infer relations about sentences. We expand the knowledge graph with these inference on the current story to get inference set  $E_{2,t}$ . Finally, we obtain the updated knowledge graph  $G'_t = \{e_{1,G'_t}, \dots, e_{l,G'_t}\}$  and inference entities set  $E_t = E_{1,t} \cup E_{2,t} = \{\hat{e}_{1,G'_t}, \dots, \hat{e}_{l,G'_t}\} \subset G'_t$ , where  $\hat{e}$  is an inferred entity and  $l$  is the number

of the inference entities.

### 3.3 Continuation Candidate Generation

Given each entity in the knowledge graph including inference entities (§3.2) and current story history, we first generate a set of continuation candidates. We consider the conditional sentence generation as a infilling task (Taylor, 1953).

**Templates.** A set of templates  $T_i$  are generated on each entity  $e_{i,G'_t}$  (See Appendix A.2). For example, one of the templates generated on beach are [subject] <mask> <mask> beach <mask>. The [subject] of the sentence is (1) the same subject with the previous sentence, (2) no fixed (<mask>), or (3) any characters in previous story history. The number of <mask> before and after inference entity ranges from 1 to 10.

**Text Infilling.** We fine-tune RoBERTa (Liu et al., 2019) on ROCStories (Mostafazadeh et al., 2016)—a set of five-sentence stories involving commonsense scenarios—to infill the mask tokens. Details of training are shown in Appendix A.3. All the templates  $T_i$  are filled by fine-tuned RoBERTa and we obtain a number of continuation candidates  $S_i = \{s_{i,1}, \dots, s_{i,m}\}$  for each entity  $e_{i,G'_t}$  where  $m$  is the number of templates of each entity.

**Filtering.** We fine-tune GPT-2 (Radford et al., 2019) on ROCStories (henceforth called GPT-2-ROC) and filter the continuation candidates by cal-

320 culating their conditional probability  $\mathbb{P}_s$  with it:

$$321 \quad \mathbb{P}_s = \prod_{j=1}^n \mathbb{P}(X_j | X_1, \dots, X_{j-1}) \quad (1)$$

322 where  $n$  is the length of the sentence  $s$  and  $X_j$  is the  
323  $j$ th token in sentence  $s$ . We only keep one sentence  
324  $s_i \in S_i$  with the highest probability for each entity  
325  $e_{i, \mathbf{G}'_t}$  and obtain the continuation candidates  $S'_t =$   
326  $\{s_1, \dots, s_l\}$ .

### 327 3.4 Graph Difference

328 We achieve controllability of the continuation candi-  
329 dates  $S'$  by calculating the *graph difference* be-  
330 tween the candidate knowledge graph  $\mathbf{G}_t^{s_i}$  and the  
331 goal knowledge graph  $\mathbf{G}_{\text{goal}}$ . The candidate knowl-  
332 edge graph  $\mathbf{G}_t^{s_i}$  is obtained by updating knowledge  
333 graph  $\mathbf{G}_t$  with continuation candidate  $s_i$ . We cal-  
334 culate the knowledge graph difference score:

$$335 \quad R(s_i) = (1 - \alpha) \times r_1(\mathbf{G}_t^{s_i}, \mathbf{G}_{\text{goal}}) \\ 336 \quad + \alpha \times r_2(E_t^i, \mathbf{G}'_{\text{goal}}) \quad (2)$$

337 where  $r_1$  is story entity overlapping score and  
338  $r_2$  is inference overlapping score.  $\alpha$  is a hyper-  
339 parameter to control the inference’s contribution  
340 on calculating overlapping rate.

341 *Story entity overlapping score* ( $r_1$ ) calculates the  
342 overlapping rate between the candidate knowledge  
343 graph  $\mathbf{G}_t^{s_i}$  and the full knowledge graph  $\mathbf{G}_{\text{goal}}$   
344 without considering inference nodes. We define  
345 a match as same entities (nodes) and their corre-  
346 sponding edges (relations) between two knowledge  
347 graph. Then calculate the story entity overlapping  
348 rate by

$$349 \quad r_1(\mathbf{G}_t^{s_i}, \mathbf{G}_{\text{goal}}) = \frac{\sum_j \sum_k \mathbb{I}(e_{j, \mathbf{G}_t^{s_i}} = e_{k, \mathbf{G}_{\text{goal}}})}{\text{size of } \mathbf{G}_{\text{goal}}} \quad (3)$$

350 where  $\mathbb{I}(e_{j, \mathbf{G}_t^{s_i}} = e_{k, \mathbf{G}_{\text{goal}}}) = 1$  when there is  
351 a match between entity  $e_{j, \mathbf{G}_t^{s_i}} \in \mathbf{G}_t^{s_i}$  and entity  
352  $e_{k, \mathbf{G}_{\text{goal}}} \in \mathbf{G}_{\text{goal}}$ , otherwise 0.

353 After updating knowledge graph  $\mathbf{G}_t$  with  $s_i$ , we  
354 expand updated knowledge graph  $\mathbf{G}_t^{s_i}$  with infer-  
355 ence nodes  $E_t^{s_i}$ , which we repeat the process in  
356 Section 3.2 on  $\mathbf{G}_t^{s_i}$ . We calculate the overlapping  
357 rate between  $E_t^{s_i}$  and goal knowledge graph with  
358 inferences nodes,  $\mathbf{G}'_{\text{goal}}$ , as *inference overlapping*  
359 *score*  $r_2$ ,

$$360 \quad r_2(E_t^{s_i}, \mathbf{G}'_{\text{goal}}) = \frac{\sum_j \sum_k \mathbb{I}(\hat{e}_{j, E_t^{s_i}} = e_{k, \mathbf{G}'_{\text{goal}}})}{\text{size of } \mathbf{G}'_{\text{goal}}} \quad (4)$$

---

**Prompt:** Doug has always wanted to be married.

---

>>> **Goal story world (without order):**  
<Doug, meet, women>, <Doug, ask, girl>  
<Doug, marry, women>, <girl, is, nice>

---

**StoRM:**

He decides to go on a **first date** with his girlfriend.  
Doug happily **spends time** with his girlfriend.  
Doug **loves** every minute of her.  
He is happy that he **married a beautiful girl**.

---

>>> **Outline:**

# too shy # ask out a nice girl # first date # marry

---

**GPT-outline:**

Doug decided to **date** Amy more than anything else.  
He was *shy and shy*.  
*She* decided he must be more serious.  
*She* went on a date with him and *two of them* hit it off.  
Later that night harry **proposed** to Amy and she said yes.

---

**Yao et al. (2019):**

Doug has always wanted to be married.  
Doug was natural grown in the kid.  
He **met his girl** someday.  
Doug asked to watch her.  
They started **dating**.  
They both will honestly meet new friends together.

---

Table 1: Story examples generated by StoRM, GPT-outline, and Yao et al. (2019). StoRM successfully finds the correct order to follow a single topic (**bolded**)—“meet a girl and marry”, and shows a good plot coherence. The story generated by GPT-outline also follows the outline(**bolded**) but suffers in repetition (in *italic*) and lacks coherence (underlined). Yao et al. (2019) suffers in plot coherence (underlined).

361 where  $\hat{e}_{j, E_t^{s_i}} \in E_t^{s_i}$  and  $e_{k, \mathbf{G}'_{\text{goal}}} \in \mathbf{G}'_{\text{goal}}$ .

362 We obtain the *topk* continuation knowledge  
363 graphs with the *topk* highest graph difference  
364 scores. They will be used to produce continuations  
365 further. We always keep a total of  $k$  knowledge  
366 graphs (states) when generating stories for each  
367 index of the sentence in the story. Thus the full  
368 generation process is implemented as a form of  
369 beam search through reader model space.

## 370 4 Experiments

371 We evaluated our system with four experiments.  
372 The first experiment accesses whether the knowl-  
373 edge graph acquisition technique captures the in-  
374 formation that natural language story conveys. The  
375 second experiment is an ablation study that assesses  
376 how each component contributes to the story gener-  
377 ation process in Figure 2. The third and fourth ex-  
378 periments compare StoRM to two neural language  
379 model story generators on the the dimensions of  
380 coherence and controllability. Story examples can

381 be found in Table 1.

382 **Datasets.** We conduct the experiments on the  
383 ROCStories corpus (Mostafazadeh et al., 2016).  
384 It contains 98, 159 five-sentence stories involving  
385 common-sense scenarios. We additionally extract  
386 outlines for ROCStories using the RAKE algo-  
387 rithm (Rose et al., 2010) for use controlling the  
388 baseline models.

389 **Baselines.** For fair comparison to baselines, we  
390 require story generation systems that operate on  
391 un-ordered outlines that abstractly indicate events  
392 that should appear in the story, and/or goal states as  
393 inputs.<sup>1</sup> We selected two strong baselines. The first  
394 is GPT-2-small (Radford et al., 2019) fine-tuned  
395 on ROCStories (Mostafazadeh et al., 2016) with  
396 outlines. GPT-2 is fed into outlines with the for-  
397 mat of {topic\_1 # topic\_2 #...#}, and then  
398 minimizes the cross entropy loss between network  
399 output logits and golden truth story from which  
400 the topics were extracted. Training details can be  
401 found in Appendix A.4.

402 The second baseline is the system by Yao  
403 et al. (2019), which trained RNN based conditional  
404 generation models on ROCStories to generate story  
405 on outline. Their Plan-and-Write system is capable  
406 of generating its own outline of topics. However,  
407 for fair comparison, we provide the system with the  
408 outline of topics extracted from our dataset. This  
409 controls for goal-seeking behavior because the out-  
410 lines are extracted from the same stories that are  
411 used to generate goal knowledge graphs for StoRM.  
412 Training details can be found in Appendix A.5.

#### 413 4.1 Knowledge Graph Acquisition Evaluation

414 We assess whether knowledge graph can acquire  
415 the story world state accurately and compre-  
416 hensively. We randomly select 125 sentences  
417 from ROCStories and convert them into knowl-  
418 edge graph triples. Human participants were  
419 asked to validate each graph triples given the sen-  
420 tence and then write down the missing informa-  
421 tion. For example, they need to check whether  
422  $\langle jenny, LIKE, beach \rangle$  given “Jenny likes beach  
423 and sunshine” is correct and write down the miss-  
424 ing concept, “sunshine”. The detail of this study  
425 is shown in Appendix B.1 and B.2.

<sup>1</sup>Two potential baselines were considered but not pursued. The system by Tambwekar et al. (2019) is goal-driven but does not produce natural language without manual intervention. The system by Rashkin et al. (2020) accepts unordered outline terms but the results of the original paper could not be

Precision %	Recall %	# of triplets
81.96 <sup>‡</sup>	72.89 <sup>‡</sup>	255

Table 2: Results of evaluating knowledge graph triples. ‡ indicates  $\kappa > 0.4$  or moderate agreement.

426 Table 2 shows the accuracy (precision) and sen-  
427 sitivity (recall) of the extracted knowledge graph  
428 triples. We treat the majority vote from human  
429 participants as the ground-truth. *Precision* is the  
430 fraction of extracted triples that are correct rated  
431 by human participants. *Recall* is the fraction of  
432 the triples that are successfully extracted from sto-  
433 ries. Precision, 81.96%, shows that the knowledge  
434 graph can represent the information in sentences  
435 accurately. Recall, 72.89%, proved that the knowl-  
436 edge graph can represent most of the information  
437 in sentences. Both of these two metrics have mod-  
438 erate agreement. This indicates that the knowledge  
439 graph extracted from sentences matches reader ex-  
440 pectations and can be used as story world state  
441 upon which to base further story generation.

#### 442 4.2 Ablation Study

443 We perform ablation studies to validate the contri-  
444 butions of different components in Figure 2. We  
445 build goal story world states (§3.1) on stories from  
446 ROCStories (Mostafazadeh et al., 2016) to guide  
447 the story generation process. StoRM keeps gener-  
448 ating story continuations until knowledge graph  
449 difference score  $R(s)$  reaches 0.8. We measure the  
450 following two metrics:

- 451 • *Average story length* (Avg. len): Calculate the  
452 average story length which is required to reach  
453  $R(s) = 0.8$  (§3.4). Smaller average story length  
454 stands for faster, and thus more direct, goal  
455 achievement. We stop generation when story  
456 length reaches 10.
- 457 • *Sentence transformer cosine similarity* (S-  
458 F) (Reimers and Gurevych, 2019): Evaluate  
459 the semantic similarity between generation and  
460 golden truth story by calculating embedding co-  
461 sine similarity. Higher S-F score indicates higher  
462 similarity between generation and gold story.

463 Table 3 shows the result of the ablation study. Re-  
464 moving ConceptNet and COMET<sub>20</sub><sup>20</sup>, which infer  
465 nodes in the reader model, significantly increases  
466 average story length and reduces similarity score. It  
467 indicates with the help of these two inference tech-  
reproduced at the time of writing.

Model		Avg. len ↓	S-F % ↑
StoRM Full	$\alpha = 0.50$	$7.96 \pm 0.73$	$77.54 \pm 4.27$
	$\alpha = 1.00$	$8.96 \pm 0.65^b$	$74.67 \pm 3.25^a$
	$\alpha = 0.25$	$9.08 \pm 0.65^a$	$70.12 \pm 5.28^b$
COMET <sub>20</sub> only		$8.96 \pm 0.74^a$	$70.10 \pm 4.06^b$
ConceptNet only		$8.85 \pm 0.69^b$	$73.71 \pm 4.12^a$
BART Candidates		$8.08 \pm 0.69$	$73.14 \pm 3.90^a$

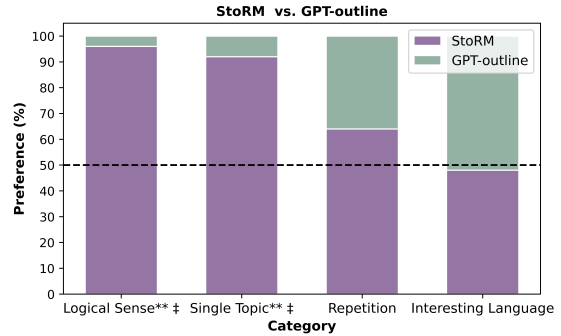
Table 3: Results of the ablation study. ConceptNet only and COMET<sub>20</sub> only indicate StoRM model only using ConceptNet or COMET<sub>20</sub> for inferring nodes. BART shows the result of replacing the whole candidate generation module with BART.  $\alpha$  is tuning the inference contribution when calculating graph difference. <sup>a</sup> and <sup>b</sup> indicate StoRM( $\alpha = 0.5$ ) results are significant better at  $p < 0.05$  and  $p < 0.01$  using the Mann-Whitney  $U$  test, respectively.

nique, StoRM is faster to achieve a better goal. Ablating the *candidate generation* module (Blue box in Figure 2) by replacing our technique with BART likewise diminishes similarity between gold story and generated story significantly. We experiment with three values of  $\alpha$  in our StoRM framework. The best performing model has  $\alpha = 0.5$ , balancing between inference-node-guided and goal-node-guided story generation, where larger  $\alpha$  indicates more inference-node-driven.

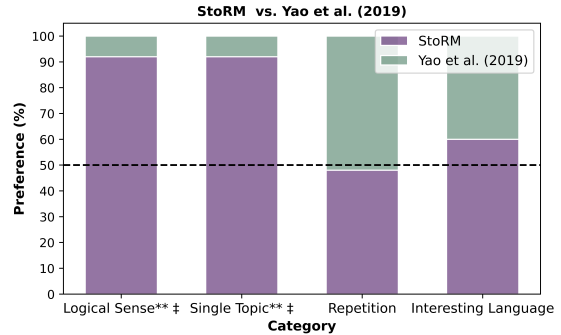
### 4.3 Story Coherence Evaluation

Having established that knowledge graph is able to represent the story, we seem to understand whether StoRM improves the coherence and quality of the generated story. In this paper, we evaluate coherence using human participant evaluation, asking a set of questions that includes dimensions such a logical coherence, loyalty to plot, and enjoyability. Variations of these questions have been used to evaluate other story generation systems (Purdy et al., 2018; Tambwekar et al., 2019; Ammanabrolu et al., 2020, 2021a; Castricato et al., 2021; Peng et al., 2021). We focus on dimensions involving overall perceptions of narrative coherence:

- *Logical Sense*: Get at narrative coherence without using the term “coherence” which can sometimes also be confused with grammaticality.
- *Follows A Single Topic*: Query about coherence since incoherent plots can look like an intertwining of several unrelated topics.
- *Repetition*: Measure which story has less repetitive words.
- *Interesting Language*: Focus on wording.



(a) StoRM vs. GPT-outline



(b) StoRM vs. Yao et al.(2019)

Figure 3: Human evaluation results comparing StoRM with two baselines, \* indicates  $p < 0.05$ , \*\* indicates  $p < 0.01$ , † indicates  $\kappa > 0.2$  or fair agreement. ‡ indicates  $\kappa > 0.4$  or moderate agreement.

Each human participant reads a randomly selected subset of story pairs, comprised of one story from StoRM and one from baselines—GPT-outline or Yao et al.. For the above four questions, participants answered which story best met the criteria. Details can be found in Appendix B.3. We build a goal knowledge graph story world to guide StoRM and use outlines extracted from the same story to seed baselines. The baselines have an advantage over StoRM; they use outlines that are ordered which should correlate with story coherence whereas our system must intuit the order of events that achieves the goal.

The results are shown in Figure 3, which indicate that StoRM performs significantly better than baselines on the the dimensions of “Logical Sense” and “Single Topic”. These are the primary dimensions that ask about story coherence. We conclude that our system improves the perception of narrative coherence of generated narratives and stays more on topic, while retaining comparably interesting language and avoidance of repetition (neither of which are statistically significantly different from the baseline). Since GPT-outline and Yao et al.(2019) are

Model	KG-o % ↑	B-2 % ↓	B-3 % ↓	B-4 % ↓	S-M % ↑	S-F % ↑	R-1 % ↑	R-2 % ↑	R-L % ↑
StoRM	<b>53.21</b>	30.96	12.21	5.41	<b>7.64</b>	<b>74.54</b>	<b>50.06</b>	<b>23.89</b>	<b>48.74</b>
GPT-outline	47.24 <sup>a</sup>	<b>27.48<sup>b</sup></b>	11.06	4.56	7.38	74.04	42.88 <sup>a</sup>	21.13	41.61 <sup>a</sup>
Yao et al.	31.67 <sup>b</sup>	27.68 <sup>a</sup>	<b>9.56<sup>b</sup></b>	<b>3.98<sup>a</sup></b>	4.90 <sup>b</sup>	66.66 <sup>b</sup>	35.49 <sup>b</sup>	7.81 <sup>b</sup>	33.83 <sup>b</sup>

Table 4: Evaluation on goal-guided story generation. Columns show average score of each model’s generated summaries according to various metrics. <sup>a</sup> and <sup>b</sup> indicate StoRM results are significant different at  $p < 0.05$  and  $p < 0.01$  using the Mann-Whitney  $U$  test. A large ROUGE score together with a low self-BLEU score demonstrate a model’s ability to generate realistic looking as well as diverse generations.

both guided by outlines, which are used as prompt to seed the neural language model, StoRM overcomes a disadvantage in that it must figure out how to reasonably order concept nodes in the goal story world and compose a coherent story.

#### 4.4 Controllability Evaluation

We assess whether StoRM is able to achieve the given goal, as measured by the coverage of all the given concepts in goal story state with the following five metrics.

- *Knowledge graph overlapping rate* (KG-o): Generated stories are transformed into knowledge graph (§3.1). Calculate the overlapping rate of knowledge graph nodes between generated story and the golden truth story.
- *Sentence transformer cosine similarity* (S-F) (Reimers and Gurevych, 2019): Evaluate semantic similarity between generation and golden truth by calculating embedding cosine similarity.
- *Sentence mover’s similarity* (Clark et al., 2019): Evaluate stories in a continuous space using word and sentence embeddings.
- *Self-BLEU score* (B-2, B-3, and B-4) (Papineni et al., 2002; Zhu et al., 2018): 2-gram, 3-gram and 4-gram self-BLEU scores are reported. These evaluate the diversity of generated stories.
- *ROUGE* (R-1, R-2, and R-L) (Lin, 2004): Higher ROUGE score indicates more coverage.

Table 4 shows the result of all the systems. We first convert all the generated stories of StoRM and baselines to knowledge graph and then calculate their knowledge graph overlapping rate (KG-o) with goal story world. Higher KG-o indicates better controllability. StoRM performs statistical significantly better than baselines in this dimension and “ROUGE”. Similarity between generated story and golden truth is also considered a way to evaluate controllability. StoRM outperforms Yao et al.(2019) in “sentence transformer similarity” and “sentence mover’s similarity” but comparable

to GPT-outline. Without seeding outlines to the language model like baselines, StoRM is able to cover most of the concepts in the goal story state. Lower self-BLEU score indicates more diversity in generated stories. The constraints imposed for controllability also impose constraints on diversity, though generally diversity is secondary to overall story quality. Yao et al. (2019) shows better diversity at the expense of controllability.

## 5 Conclusions

Neural language models are widely used to produce text, including stories. However, they struggle with maintaining *story coherence*—the logical progression of events—and goal-directedness. Our framework—*Story Generation with Reader Models* (StoRM)—augments neural language models with a reader model. This reader model—in this case an explicit knowledge graph—approximates the reader’s beliefs about the story world. StoRM increases the story coherence by expanding the reader model with commonsense technique and producing continuations by selecting entities in this reader model. In order to achieve goal-directedness, the reader model enables the system to make informed decisions about which possible story continuations are likely to achieve the desired goal state by inferring how each possible story continuation changes the reader model to be closer to the desired goal.

A thorough experimental study shows that StoRM produces significantly more coherent and goal-directed stories than two strong baselines. The goal-directness results are significant because the StoRM framework takes a goal as a knowledge graph, which can be thought of as an unordered outline of concepts that should appear in the story; our system does well to find an appropriate sequencing of events. Thus a reader model based approach provides improved story coherence while providing users a powerful means of control.



## 6 Broader Impact

Our system faces the same potential pitfalls as other contemporary language learning systems. It is prone to echoing the biases present in the dataset (Sheng et al., 2019) and generate non-normative text (i.e. in violation of social norms). No existing automated storytelling systems is able to entirely eliminate these biases, though stories can be used to teach language models to reduce non-normative continuations (Peng et al., 2020). Fictional stories that are presented to readers as non-fictional can be used to influence (Green and Brock, 2000) or misinform. Future work may enable real-world facts to be injected into the knowledge graph of a similar system for the purposes of journalism or misinformation. However, because our graph expansion method relies on ConceptNet5 (Speer and Havasi, 2013) and COMET<sub>20</sub> (Hwang et al., 2021) for inference, our system is prone to process and produce simple stories.

The ability to produce coherent and goal-directed stories has downstream applications beyond automated story-telling. Our work is also applicable to figure out how to reasonably order concept nodes and validate whether there exists a multi-hop explanation for how two concepts are related.

## References

Prithviraj Ammanabrolu, Wesley Cheung, William Broniec, and Mark O Riedl. 2021a. Automated storytelling via causal, commonsense plot ordering. In *Proceedings of AACL*, volume 35.

Prithviraj Ammanabrolu, Wesley Cheung, William Broniec, and Mark O Riedl. 2021b. Automated storytelling via causal, commonsense plot ordering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 5859–5867.

Prithviraj Ammanabrolu, Ethan Tien, Wesley Cheung, Zhaochen Luo, William Ma, Lara J Martin, and Mark O Riedl. 2020. Story realization: Expanding plot events into sentences. In *Proceedings of AAAI*, volume 34.

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.

Louis Castricato, Spencer Frazier, Jonathan Balloch, and Mark O. Riedl. 2021. Tell me a story like i’m five: Story generation via question answering. In

*Proceedings of the 3rd Workshop on Narrative Understanding*.

Marc Cavazza, Olivier Martin, Fred Charles, Steven J Mead, and Xavier Marichal. 2003. Interacting with virtual agents in mixed reality interactive storytelling. In *International Workshop on Intelligent Virtual Agents*, pages 231–235. Springer.

Elizabeth Clark, Asli Celikyilmaz, and Noah A Smith. 2019. Sentence mover’s similarity: Automatic evaluation for multi-sentence texts. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2748–2760.

Elizabeth Clark, Yangfeng Ji, and Noah A. Smith. 2018. Neural text generation in stories using entity representations as context. In *Proceedings of NAACL-HTL*, pages 2250–2260.

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.

Andrea Di Fabio, Simone Conia, and Roberto Navigli. 2019. Verbatlas: a novel large-scale verbal semantic resource and its application to semantic role labeling. 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 627–637.

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. *arXiv preprint arXiv:1902.01109*.

Christiane Fellbaum. 1998. Wordnet: An electronic lexical database cambridge. MA: MIT Press.

Daniel Gildea and Daniel Jurafsky. 2002. Automatic labeling of semantic roles. *Computational linguistics*, 28(3):245–288.

Arthur C Graesser, Murray Singer, and Tom Trabasso. 1994. Constructing inferences during narrative text comprehension. *Psychological review*, 101(3):371.

Melanie C. Green and Timothy C. Brock. 2000. The role of transportation in the persuasiveness of public narratives. *Journal of Personality and Social Psychology*, 79(5):701–721.

Jian Guan, Yansen Wang, and Minlie Huang. 2019. Story ending generation with incremental encoding and commonsense knowledge. In *Proceedings of*

709		<i>the AAAI Conference on Artificial Intelligence</i> , volume 33, pages 6473–6480.	
710			
711	Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras,		
712	Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and		
713	Yejin Choi. 2021. Comet-atomic 2020: On sym-		
714	bolic and neural commonsense knowledge graphs. In		
715	AAAI.		
716	Ahmed Khalifa, Gabriella AB Barros, and Julian		
717	Togelius. 2017. Deeptingle. <i>arXiv preprint</i>		
718	<i>arXiv:1705.03557</i> .		
719	Michael Lebowitz. 1987. Planning stories. In <i>Proceed-</i>		
720	<i>ings of the 9th annual conference of the cognitive</i>		
721	<i>science society</i> , pages 234–242.		
722	Kenton Lee, Luheng He, Mike Lewis, and Luke Zettle-		
723	moyer. 2017. End-to-end neural coreference resolu-		
724	tion. <i>arXiv preprint arXiv:1707.07045</i> .		
725	Chin-Yew Lin. 2004. Rouge: A package for automatic		
726	evaluation of summaries. In <i>Text summarization</i>		
727	<i>branches out</i> , pages 74–81.		
728	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-		
729	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,		
730	Luke Zettlemoyer, and Veselin Stoyanov. 2019.		
731	Roberta: A robustly optimized bert pretraining ap-		
732	proach. <i>arXiv preprint arXiv:1907.11692</i> .		
733	Ilya Loshchilov and Frank Hutter. 2017. Decou-		
734	pled weight decay regularization. <i>arXiv preprint</i>		
735	<i>arXiv:1711.05101</i> .		
736	Lara Martin, Prithviraj Ammanabrolu, Xinyu Wang,		
737	William Hancock, Shruti Singh, Brent Harrison, and		
738	Mark Riedl. 2018. Event representations for auto-		
739	mated story generation with deep neural nets. In		
740	<i>Proceedings of AAAI</i> , volume 32.		
741	James Richard Meehan. 1976. <i>The Metanovel: Writing</i>		
742	<i>Stories by Computer</i> . Yale University.		
743	Todor Mihaylov and Anette Frank. 2018. Knowledge-		
744	able reader: Enhancing cloze-style reading compre-		
745	hension with external commonsense knowledge. In		
746	<i>Proceedings of the 56th Annual Meeting of the As-</i>		
747	<i>sociation for Computational Linguistics (Volume 1:</i>		
748	<i>Long Papers)</i> , pages 821–832.		
749	Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong		
750	He, Devi Parikh, Dhruv Batra, Lucy Vanderwende,		
751	Pushmeet Kohli, and James Allen. 2016. A corpus		
752	and evaluation framework for deeper understanding		
753	of commonsense stories. In <i>Proceedings of NAACL-</i>		
754	<i>HLT</i> , pages 839–849.		
755	Martha Palmer, Daniel Gildea, and Paul Kingsbury.		
756	2005. The proposition bank: An annotated corpus of		
757	semantic roles. <i>Computational linguistics</i> , 31(1):71–		
758	106.		
759	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-		
760	Jing Zhu. 2002. Bleu: a method for automatic evalu-		
761	ation of machine translation. In <i>Proceedings of the</i>		
762	<i>40th annual meeting of the Association for Computa-</i>		
763	<i>tional Linguistics</i> , pages 311–318.		
	Xiangyu Peng, S. Li, Spencer Frazier, and Mark O.		
	Riedl. 2020. <a href="#">Reducing non-normative text generation</a>		
	<a href="#">from language models</a> . In <i>International Conference</i>		
	<i>on Natural Language Generation</i> .		
	Xiangyu Peng, Siyan Li, Sarah Wiegrefe, and Mark		
	Riedl. 2021. Inferring the reader: Guiding automated		
	story generation with commonsense reasoning. <i>arXiv</i>		
	<i>preprint arXiv:2105.01311</i> .		
	Julie Porteous and Marc Cavazza. 2009. Controlling		
	narrative generation with planning trajectories: the		
	role of constraints. In <i>Joint International Conference</i>		
	<i>on Interactive Digital Storytelling</i> , pages 234–245.		
	Springer.		
	Christopher Purdy, Xinyu Wang, Larry He, and Mark		
	Riedl. 2018. Predicting generated story quality with		
	quantitative measures. In <i>Proceedings of the AAAI</i>		
	<i>Conference on Artificial Intelligence and Interactive</i>		
	<i>Digital Entertainment</i> , volume 14.		
	Alec Radford, Jeff Wu, Rewon Child, David Luan,		
	Dario Amodei, and Ilya Sutskever. 2019. Language		
	models are unsupervised multitask learners.		
	Hannah Rashkin, Asli Celikyilmaz, Yejin Choi, and		
	Jianfeng Gao. 2020. Plotmachines: Outline-		
	conditioned generation with dynamic plot state track-		
	ing. In <i>Proceedings of the 2020 Conference on Empir-</i>		
	<i>ical Methods in Natural Language Processing</i>		
	(EMNLP), pages 4274–4295.		
	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-bert:</a>		
	<a href="#">Sentence embeddings using siamese bert-networks</a> .		
	In <i>Proceedings of the 2019 Conference on Empirical</i>		
	<i>Methods in Natural Language Processing</i> . Associa-		
	tion for Computational Linguistics.		
	Mark O Riedl and Robert Michael Young. 2010. Narra-		
	tive planning: Balancing plot and character. <i>Journal</i>		
	<i>of Artificial Intelligence Research</i> , 39:217–268.		
	Melissa Roemmele. 2016. Writing stories with help		
	from recurrent neural networks. In <i>Proceedings of</i>		
	AAAI, volume 30.		
	Stuart Rose, Dave Engel, Nick Cramer, and Wendy		
	Cowley. 2010. Automatic keyword extraction from		
	individual documents. <i>Text mining: applications and</i>		
	<i>theory</i> , 1:1–20.		
	Maarten Sap, Ronan Le Bras, Emily Allaway, Chan-		
	dra 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 <i>Proceedings of the AAAI Con-</i>		
	<i>ference on Artificial Intelligence</i> , volume 33, pages		
	3027–3035.		
	Emily Sheng, Kai-Wei Chang, Premkumar Natarajan,		
	and Nanyun Peng. 2019. The woman worked as a		
	babysitter: On biases in language generation. <i>arXiv</i>		
	<i>preprint arXiv:1909.01326</i> .		

- 817 Peng Shi and Jimmy Lin. 2019. Simple bert models for  
818 relation extraction and semantic role labeling. *arXiv*  
819 *preprint arXiv:1904.05255*.
- 820 Robert Speer and Catherine Havasi. 2013. Conceptnet 5:  
821 A large semantic network for relational knowledge.  
822 In *The People’s Web Meets NLP*, pages 161–176.  
823 Springer.
- 824 Pradyumna Tambwekar, Murtaza Dhuliawala, Lara J  
825 Martin, Animesh Mehta, Brent Harrison, and Mark O  
826 Riedl. 2019. Controllable neural story plot genera-  
827 tion via reinforcement learning. In *Proceedings of*  
828 *the 28th IJCAI*.
- 829 Wilson L Taylor. 1953. “cloze procedure”: A new  
830 tool for measuring readability. *Journalism quarterly*,  
831 30(4):415–433.
- 832 Stephen Ware and Cory Siler. 2021. Sabre: A narra-  
833 tive planner supporting intention and deep theory of  
834 mind. In *Proceedings of the 17th AAAI International*  
835 *Conference on Artificial Intelligence and Interactive*  
836 *Digital Entertainment*.
- 837 Stephen G Ware and R Michael Young. 2010. Modeling  
838 narrative conflict to generate interesting stories. In  
839 *Sixth Artificial Intelligence and Interactive Digital*  
840 *Entertainment Conference*.
- 841 Lili Yao, Nanyun Peng, Ralph Weischedel, Kevin  
842 Knight, Dongyan Zhao, and Rui Yan. 2019. Plan-  
843 and-write: Towards better automatic storytelling. In  
844 *Proceedings of the AAAI Conference on Artificial*  
845 *Intelligence*, volume 33, pages 7378–7385.
- 846 Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan  
847 Zhang, Jun Wang, and Yong Yu. 2018. Taxygen: A  
848 benchmarking platform for text generation models.  
849 In *The 41st International ACM SIGIR Conference on*  
850 *Research & Development in Information Retrieval*,  
851 pages 1097–1100.
- 852 Rolf A Zwaan and Gabriel A Radvansky. 1998. Situa-  
853 tion models in language comprehension and memory.  
854 *Psychological bulletin*, 123(2):162.

## 855 A Implementation Details

### 856 A.1 Semantic Role Labeling Using VerbAtlas

857 The SRL model provides the automatic identi-  
858 fication and labeling of argument structures of  
859 stories. For example, it extracts ‘verbatlas’:  
860 ‘EXIST\_LIVE’, ‘args\_words’: {‘Theme’:  
861 ‘Jenny’, ‘Attribute’: ‘Florida’} from  
862 “Jenny lived in Florida”. Verbs in the story will be  
863 represented as the VerbAtlas frame. For example,  
864 “live” is represented as “EXIST\_LIVE”.

865 For the semantic role labeling model (SRL), we  
866 use a fine-tuned transformer model proposed by  
867 (Shi and Lin, 2019) which is the current state-  
868 of-the-art for English SRL. It is a BERT (Devlin  
869 et al., 2019) model with a linear classification layer  
870 trained on the Ontonotes 5.0 dataset to predict Prop-  
871 Bank SRL. We use an open-source implementation  
872 <sup>2</sup>, which is based on the official AllenNLP BERT-  
873 SRL model <sup>3</sup>. Trained with the following hyperpa-  
874 rameters:

- 875 • Batch size: 32
- 876 • Dropout for the input embeddings: 0.1
- 877 • Learning rate:  $5e^{-5}$
- 878 • Optimizer: Adam
- 879 • Total Epochs: 15

880 Then, we use the mappings from Propbank  
881 frames to VerbAtlas (Di Fabio et al., 2019) classes  
882 to return the correct corresponding VerbAtlas  
883 classes instead of Propbank’s (Palmer et al., 2005).  
884 The direct mapping is possible because, for every  
885 VerbAtlas class, there is only one PropBank frame,  
886 which allows us to utilize the rich content provided  
887 by VerbAtlas while using the same model initially  
888 trained to predict PropBank.

### 889 A.2 Continuation Candidate Generation 890 Details

891 For each entity  $e_{i,G'_t}$ , we generate  $5 \times 10 \times (c + 1)$   
892 templates, where  $c$  is the number of subjects in the  
893 story history. Rules of making the templates are as  
894 follows,

- 895 • The first token in the template has two  
896 choices: (1) Previous subject, i.e. “Jenny”;  
897 (2)  $\langle mask \rangle$ .
- 898 • Between the first token and the entity  $e_{i,G'_t}$ ,  
899 we put  $0 \sim 4 \langle mask \rangle$ .
- 900 • After the entity  $e_{i,G'_t}$ , we put  $0 \sim 10 \langle mask \rangle$   
901 tokens

902 Examples are as follows when  $e_{i,G'_t}$  = “beach”,

- 903 • Jenny  $\langle mask \rangle$  beach  $\langle mask \rangle$ .
- 904 • Jenny  $\langle mask \rangle$   $\langle mask \rangle$  beach.
- 905 • Jenny  $\langle mask \rangle$   $\langle mask \rangle$  beach  $\langle mask \rangle$   
906  $\langle mask \rangle$   $\langle mask \rangle$ .

### 907 A.3 RoBERTa Fine-tuning

908 We fine-tune RoBERTa (Liu et al., 2019) on ROC-  
909 Stories (Mostafazadeh et al., 2016) to infill the  
910 mask tokens in the given text template. We pre-  
911 process the ROCStories by masking 15% of all the  
912 tokens randomly, concatenating all texts together,  
913 and splitting them into chunks of the same length  
914 (equal to 128). Each chunk is then used as one  
915 training sample.

916 During fine-tuning, we use the AdamW opti-  
917 mizer (Loshchilov and Hutter, 2017) to train the  
918 RoBERTa for 3 epochs with batch size = 8. Other  
919 optimizer-related hyperparameters are attached as  
920 follows.

- 921 • learning rate:  $\gamma = 2 \times 10^{-5}$
- 922 • betas:  $\beta_1 = 0.9, \beta_2 = 0.999$
- 923 • epsilon:  $\epsilon = 10^{-8}$
- 924 • weight decay:  $\lambda = 0.01$

### 925 A.4 Baselines—GPT-outline

926 We use the small version of GPT-2 (Radford  
927 et al., 2019) with 124M parameters as the  
928 base for all fine-tuned models. We converted  
929 the data as the following format: topic\_1 #  
930 topic\_2 # ...# topic\_n # Stories. For ex-  
931 ample, Florida # beach # Jenny lived in  
932 Florida. She loves beach. When fine-  
933 tuning GPT-2 on ROCStories and the common-  
934 sense knowledge resources (done separately), we  
935 train with a batch-size of 16, a learning rate of  
936 0.00005, and using the Adam optimizer with gradi-  
937 ent clipping at a max norm of 1. GPT-2 is fed into  
938 outlines with the format of {topic\_1 # topic\_2  
939 #...#}, and then minimizes the cross entropy loss  
940 between network output logits and gold truth story  
941 from which the topics were extracted. All models  
942 were trained on single GeForce RTX 2080 GPUs  
943 in Pytorch using the Huggingface Transformers  
944 library.<sup>4</sup>

### 945 A.5 Baselines—Yao et al.(2019)

946 We replicate the Plan and Write model using the  
947 code published on the paper’s public repository <sup>5</sup>

<sup>2</sup><https://github.com/Riccorl/transformer-srl>

<sup>3</sup><https://demo.allennlp.org/semantic-role-labeling>

<sup>4</sup><https://huggingface.co/transformers/>

<sup>5</sup><https://bitbucket.org/VioletPeng/language-model/>

948 to train our baseline (Yao et al., 2019). We first  
949 filtered the training dataset, removing our test story  
950 outlines to prevent data leakage. Then we train the  
951 model with the hyperparameters specified in their  
952 documentation. Below are those hyperparameter  
953 values:

- 954 • Dropout for input embedding: 0.4
- 955 • Dropout for the RNN layers: 0.25
- 956 • Random Seed: 141
- 957 • Total Epochs: 500
- 958 • Word Embedding size: 1000
- 959 • Hidden units size per RNN layer: 1000
- 960 • Batch size: 80
- 961 • Learning rate: 30
- 962 • Optimizer: SGD
- 963 • Sequence Length: 70

964 However, the best checkpoint was around 100  
965 epochs as the validation loss stopped decreasing at  
966 2.85.

## B Human Evaluation Details

### B.1 Task Instruction

We ask participants a set of questions to make sure they understand our task. The details can be found in Figure 4.

You will be asked to read a sentence, and then answer questions about **triplets** based on that.

**Triples** are composed of  $\langle \text{entity}_1, \text{relation}, \text{entity}_2 \rangle$ .

Any sentence can be represented as several *triples*.

For example, 'Jenny loves beach and sunshine' can be represented as  $\langle \text{Jenny}, \text{LOVE}, \text{beach} \rangle$  and  $\langle \text{Jenny}, \text{LOVE}, \text{sunshine} \rangle$ .

PS: **entity\_1 and entity\_2's order can be swapped**. For example,  $\langle \text{Jenny}, \text{LOVE}, \text{beach} \rangle$   $\square$   $\langle \text{beach}, \text{LOVE}, \text{Jenny} \rangle$

>> What is the **goal** of this survey?

Please select the correct **triplet** to represent "Linda graduate from college in the USA".

(Multiple choices, select all that apply)

(USA, IN, college)

(Linda, GRADUATE, college)

(LINDA, LEAVE, USA)

Please write down the **triplets** to represent "Jenny lived in Florida"

No need to worry about format, you can use  $\langle \text{aa}, \text{bbb}, \text{cc} \rangle$

Figure 4: Screenshot of the human study instruction.

### B.2 Knowledge Graph Acquisition Evaluation Set-up

We assess whether knowledge graph can acquire the story world state accurately and comprehensively. We randomly select 125 sentences from ROCStories and convert them into knowledge graph triplets. We recruited 30 participants on a crowdsourcing platform. Each participant read a randomly selected subset of knowledge graph triplets (20 sentences per participant). They were asked to validate each graph triplets given the sentence and then write down the missing information. An example is shown in Figure 5. At least 3 crowd workers validate each triple and we take the majority vote as the result.

For each triplet, please check whether it is correct given the following sentence:

**Glen was told about a first game against a rival school**

	Correct	Wrong
$\langle \text{school}, \text{is}, \text{rival} \rangle$	<input type="radio"/>	<input type="radio"/>
$\langle \text{Glen}, \text{tell}, \text{game} \rangle$	<input type="radio"/>	<input type="radio"/>
$\langle \text{Glen}, \text{tell}, \text{school} \rangle$	<input type="radio"/>	<input type="radio"/>

Please write down the triplets you think are missing.  
If not, please write down **N/A**.

Figure 5: Screenshot of Knowledge Graph Acquisition evaluation.

### B.3 Story Coherence Evaluation Set-up

We evaluate coherence using human participant evaluation, asking a set of questions that includes dimensions such as a logical coherence, loyalty to plot, and enjoyability. Example of human study is shown in Figure 6. We ask the following four questions:

- Which story makes better logical sense?
- Which story follows a single topic?
- Which story avoids repetition?
- Which story uses more interesting language?

We recruited 40 participants and each participant reads a randomly selected subset of 10 story pairs, comprised of one story from StoRM and one from baselines—GPT-outline or Yao et al.. For the above four questions, participants answered which story best met the criteria. Our study was approved by our Institutional Review Board, and we payed participants the equivalent of \$15/hr. To generate the stories, we randomly selected 25 stories from the ROCStories corpus.

For each question, please rate which story best fits.

	Eric loved ice cream He got a cone He went to eat it He got distracted and ate too much Eric got fat because of it Eric began to lose weight Eric is no longer ice cream obsessed Eric now says he needs to start using other foods more Eric is glad he stopped eating ice cream	Eric loved ice cream Eric really fat He didn't want to burn his fat Eric's doctor told him to start exercising to burn fat Eric wanted to be fit and fat again Eric started exercising and losing weight After 10 months of exercise jim no longer felt fat and felt great
1. Which story makes better Logical Sense?	<input type="radio"/>	<input type="radio"/>
2. Which story Follows A Single Topic better?	<input type="radio"/>	<input type="radio"/>
3. Which story Avoids Repetition more?	<input type="radio"/>	<input type="radio"/>
4. Which story uses more Interesting Language?	<input type="radio"/>	<input type="radio"/>

Figure 6: Screenshot of the human study on evaluating coherence.