

Visual Writing Prompts: Character-Grounded Story Generation with Curated Image Sequences

Xudong Hong^{1,2}, Asad Sayeed³, Khushboo Mehra²,
Vera Demberg² and Bernt Schiele¹

¹Dept. of Computer Vision and Machine Learning, MPI Informatics

²Dept. of Language Science and Technology, Saarland University

³Dept. of Philosophy, Linguistics, and Theory of Science, University of Gothenburg

{xhong, kmehra, vera}@coli.uni-saarland.de
schiele@mpi-inf.mpg.de, asad.sayeed@gu.se

Abstract

Current work on image-based story generation suffers from the fact that the existing image sequences collections do not have coherent plots behind them. We improve visual story generation by producing a **new image-grounded dataset, Visual Writing Prompts (VWP)**. VWP contains **almost 2K selected sequences of movie shots, each including 5-10 images**. The image sequences are aligned with a total of 12K stories which were collected via crowdsourcing given the image sequences and a set of grounded characters from the corresponding image sequence. Our new image sequence collection and filtering process has allowed us to obtain stories that are more coherent and more diverse compared to previous work. We also propose a character-based story generation model driven by coherence as a strong baseline. Evaluations show that our generated stories are more coherent, visually grounded, and more diverse than stories generated with the current state-of-the-art model. **This is a pre-MIT Press publication version.**

Stories play an important role in natural language understanding and generation because they are the key mechanism for humans to understand meaning and knowledge in the world (Piper et al., 2021). Automatically generating a coherent and interesting story is a complex task requiring various capabilities in language processing, event comprehension, and world knowledge to come together. Previous approaches to story telling have used different kinds of input to guide the story: some use a textual prompt to start the story (Fan et al., 2018). Yet others involve describing a sequence of images to direct the story (Huang et al., 2016). We choose to work inside the latter family of approaches in order to exploit the rich information contained in characters and to prevent suffering from the grounding problem (Harnad, 1990).

Research on visual narratives shows how it would be possible to construct the sort of dataset we propose: image sequences should consist of a series of coherent events centered around one or more main characters (Cohn, 2020). In fact, even Aristotle already points out in *Poetics* that *event* and *character* are the most important elements for a good story.

1 Introduction

In this work, we improve the quality of text stories generated by neural models from image sequences. We do so by improving the curation of the image sequences that form the basis for collecting the story/image pairs used to train the models: we build a dataset in which the images lend themselves better to telling a story. To show the usefulness of our dataset, we train a coherence-driven model where we design a coherence component inspired by entity grid models. Experiments show that our model produces more coherent, visually-grounded stories with more diversity than previous models¹.

¹We will release our code, image features, annotations and collected stories on a website.

To date, several datasets of image sequences for narrative generation exist, such as the Visual Storytelling (VIST; Huang et al., 2016) dataset, which includes sets of images extracted from Flickr albums. However, image sequences generated this way have the drawback that they may not lend themselves well to storytelling. Consider for instance the image sequence shown in the first column of Fig. 1: the people featured across the image sequence are all different, and there is no real development of an event or a plot. This means that the stories that humans were able to write for these types of image sequences are often quite poor from a narrative point of view and therefore lead to low-quality training data for our story gen-

Visual Storytelling	Travel Blogs	Visual Writing Prompts (Ours)
		
Shoppers riding the escalator at the mall.	sorry to be absent lately mes cheris but it was necessary to put myself on a little nyc staycation. with all the running around i have done in the last couple months i finally had the opportunity to rest ...	Jack was on a call with a client, getting stressed over a business deal that wasn't going well.
		
So many people are shopping today.	went shopping in soho. i love passing all the creative storefronts around that nabe. how fun and regal are these doors?	Jack put the phone down after an unsuccessful deal and decided to go get a coffee at the nearby coffee.
		
Two friends going into the mall for the great sales.	you know it's going to be a good day when you start off your morning with magnolia bakery breakfast. raspberry crumb muffin coffee infinity scarf and gaga glasses. done and done.	At the coffee shop, he started talking to the waiter Will about the unfortunate call.
		
Three men in yellow vest outside the mall.	i watched the enterprise space shuttle fly over manhattan as it made its voyage to its new nyc home at the intrepid air and space museum. bonus points for living on the hudson river? ...	Will told him he would convince the client to accept the deal if he could work for Jack.
		
Picture of the old home we will visit on vacation.	had an all day long adventure to ikea on saturday which of course consisted of taking the nyc water taxi out to brooklyn's ikea. the southstreet seaport is always a great photo ...	Will then called the client and successfully struck the deal.

Figure 1: Comparison of Visual Writing Prompts dataset with Visual Storytelling and Travel Blogs datasets. Our dataset has recurring characters across all five images and sub-stories. Each appearance of a character in a sub-story has a bounding box in the corresponding image, which grounds the textual appearance to visual input.

eration algorithms, which in turn, unsurprisingly, generate quite bad stories.

We thus argue that image sequences serving as writing prompts should be comprehensible as visual narratives by humans. Humans (with reasonable writing proficiency) can then “translate” such visual narratives into textual narratives. For an image sequence to qualify as a visual narrative, events and characters must have two properties: *coherence*, meaning that the events are semantically related and centered around recurring characters; and *diversity*, meaning that several different events jointly construct a plot. Psycholinguistic experiments show that missing either of these properties impedes human comprehension of image sequences as visual narratives (Cohn et al., 2012). In addition, the characters should also be easily recognized in the image sequences and can be straightforwardly linked to the stories (*visual groundedness*). Image sequences without these properties are hence not effective writing prompts.

In this work, we define the term *visual tellability* to mean the *tellability* (Hühn et al., 2014) of image sequences, i.e., how likely it is that humans can write a story with an image sequence, which measures whether the image sequences have the

two properties described above. We propose a new dataset, Visual Writing Prompts (VWP), containing curated image sequences and matching user-generated stories, linking the image sequences into coherent stories. Our image selection process allows us to choose optimized image sequences that have high visual tellability, and to encourage our crowdsourced storytellers to produce coherent stories with high diversity.

To obtain coherent and visually grounded stories, we provide cropped images of characters explicitly with image sequences for storytellers. To improve narrativity and diversity, we select images from a data source that is already likely to have a plot: image sequences selected from movie scenes with aligned synopses. To further show the importance of coherence and visual groundedness, we propose a story generation model with a representation of visual coherence focused principally on character continuity as a strong baseline. Experiments show that our model outperforms the current state-of-the-art and generates stories that are more coherent, visually grounded, and have higher diversity.

We summarize our contributions in this work as follows: (a) We propose a pipeline to extract

images sequences automatically from annotated movies as story writing prompts, which leads to image sequences with higher visual tellability. (b) We collect a new dataset of stories based on curated image sequences with grounded characters, which is more coherent and has better diversity than previous datasets. (c) We propose a character-grounded story generation model driven by visual coherence as a strong baseline for image-based story generation, which generates more coherent, diverse and visually grounded stories than the current state-of-the-art model TAPM (Yu et al., 2021).

2 Related Work

Story generation. Several datasets have been presented for generating a story conditioned on a prompt such as title (Fan et al., 2018), keywords (Yao et al., 2019), cue phrases (Xu et al., 2020) or story plot (Rashkin et al., 2020). The ROCStories corpus (Mostafazadeh et al., 2016) is a collection of short stories with rich causal and temporal relations. In subsequent work, new datasets have also been formed by gathering annotations on subsets of ROCStories for specialized story generation tasks such as modeling character psychology (Rashkin et al., 2018) counterfactual reasoning (Qin et al., 2019), etc. The STORIUM dataset (Akoury et al., 2020) of collaboratively-written long stories contains rich annotations such as narrator prompts, character goals, and other attributes to guide story generation. However, all these datasets relying on textual prompts suffer from the grounding problem that the meanings of textual stories are grounded on textual symbols (Harnad, 1990). In contrast, our dataset contains stories grounded on characters in image sequences, i.e. nonsymbolic prompts from visual perception.

Visually-grounded stories. Early work on the VIST dataset (Huang et al., 2016) identified that language in visually-grounded stories is much more diverse than in image captions. However, most of the previous datasets of visually-grounded stories have limitations because characters are not explicitly annotated (Chandu et al., 2019), the dataset is limited in scale (Xiong et al., 2019), or there is no sequence of events behind the images (Park and Kim, 2015; Huang et al., 2016). Our dataset is the first large-scale dataset that is focused on overcoming these limitations. Unlike

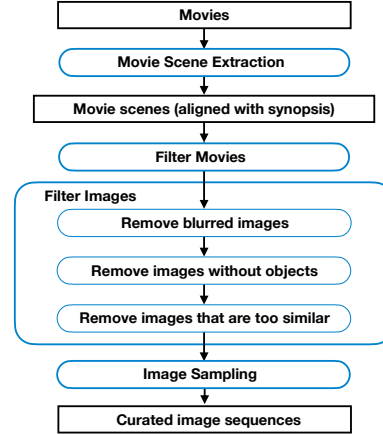


Figure 2: Image processing pipeline. Black squares are input or output. Circles are processing steps.

the VIST images, images in our VWP dataset do not feature people posing for the camera in limited contexts. Instead, they depict a rich range of situations, interactions, and emotions. Furthermore, providing character annotations in VWP ensures that the entities in the narrative are grounded to the image sequence and can be easily tracked across the sequence even when some visual attributes change. We hypothesize that these features will result in more coherent and visually grounded stories while maintaining a high level of diversity.

3 Image Sequence Construction

In this section, we describe how we obtain image sequences and design a pipeline to filter and sample images. Our objective is to construct image sequences that are *visually tellable*, i.e. are coherent and have high diversity. Our full pipeline for image sequence construction is shown in Figure 2.

Movie Scene Extraction. To achieve high coherence and diversity, we choose to select images from movie scenes that have a plot consisting of a series of events around several main characters. We extract movie scenes from MovieNet (Huang et al., 2020) since it is a dataset that contains movie synopses, annotated movie scenes with extracted movie shots and identified main characters. The paragraphs in each movie synopsis describe subplots of the movie plot, which are aligned with one or more movie scenes.

Changing from one paragraph to another in the synopsis indicates scene changes (Xiong et al., 2019). Moreover, events and characters in one movie scene are semantically coherent. We can

View a sequence of images and figure out the content. Then write a story with it.

- View a sequence of images as many times as you wish.
- Figure out who were involved and what happened.
- Then write a story that fits the image sequence.
- After the writing, answer two multiple choice questions.

Requirements:

- You need to be a **native speaker** of English. Please exit this task if you are not a native speaker.
- You should write the story using **at least 5 images**. You need to write **at least 50** but **no more than 300** words. You don't need to write in a text box without a corresponding image unless it is necessary.
- The story should be related to the image sequence. Describe only the most important character(s) and event(s).
- When mentioning the characters, please follow their names on the left. You can use either the first name, a pronoun or a noun phrase according to the context. If the character you want to mention is not there, name the characters as you want, but please be consistent.
- Use punctuation and letter case correctly.
- Don't mention that you're describing an image. Avoid using phrases like "In this image, ...".
- Do **not** write a monologue of a character or a dialogue between characters.

Breaking at least one of the requirements above automatically leads to rejection.

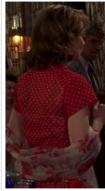
WARNING: This HIT may contain adult content. If you encounter any explicit images, please skip the image or the task.

Thank you for reading and cooperation! If you have any question, feel free to contact Xudong Hong (xhong@coll.uni-saarland.de)

Main characters:



Jack



Helen

- ☐
- ☐
- ☐

Please write the story in the corresponding boxes:



1.

The anxious waiter walked up to the couple and told them that there was a



2.

Jack and Helen



3.



4.



5.



6.



7.

Word Count (50~300): 38

Image Count (at least 5): 2

Full story:

The anxious waiter walked up to the couple and told them that there was a dress code for the restaurant and Jack was breaking it, he would need to leave and return with a jacket. Jack and Helen

Figure 3: Worker interface on Amazon Mechanical Turk. The instructions and the requirements are presented first. The main characters are provided on the left side. On the right side, each image is accompanied by a *textarea*. The full story is presented under the input area. We also show the word count and the number of images used for workers' convenience.

make use of these properties to achieve high diversity by sampling image sequences from movie scenes aligned with only one paragraph, so that image sequences are from one sub-plot with a series of different events.

Filtering Movies. Since we want to constrain the range of commonsense inferences of storytellers to the real world and help them to produce coherent stories, we first filter out all fantasy, science fiction, and horror movies. We also filter out all animations because their image characteristics are too different from the other movies.

Filtering Images.² To help storytellers to write stories that are visually grounded on characters or objects around them, we discard blurry images and images without any COCO “objects” (Lin et al., 2014)³. We measure the amount of image blur by calculating the variance of the Laplacian (Pech-Pacheco et al., 2000) and remove images with a variance lower than 30. We further apply a MaskRCNN-based object detector and filter out

²Hyper-parameters in this section are determined by a preliminary experiment that optimizes the filter process manually on 50 image sequences.

³A human character is also labeled as an “object” in MSCOCO dataset for object detection.

images without any detected objects – this will help us generate stories with interesting grounding in the image.

To increase the diversity of image sequences, we need to avoid including shots that are very similar (as can happen when a person speaks in a long monologue, for example) to one another. To detect the degree of similarity, we first feed the images to a ResNet-50 pre-trained on ImageNet and extract image features after the $fc7$ layer. Then we compute pairwise cosine similarities of the image features within each image sequence and discard an image if its cosine similarity with any one of the other images is larger than 0.89.

Additionally, we detect adult content by applying a pre-trained classifier⁴ and exclude images that trigger the classifier. We also remove the first or the last image sequence in a movie to avoid images with credits.

Image Sampling. The most intuitive way to collect stories is to use extracted movie scenes directly as writing prompts. Since these movie scenes contain a large number of movie shots, we control the workload by constraining the number of images for each writing task to a lower number K which is obtained through the 2nd pilot studies in Section 4.1. So from each selected movie scene, we sample images consecutively in non-overlapping sliding windows with a size of K and use each set of K images as one writing prompt.

4 Crowdsourcing Experiment Design

In this section, we design a crowdsourcing experiment to collect stories using our collected image sequences as writing prompts. Our objective is to obtain coherent stories that have high diversity from crowdsourced storytellers.

We design and run all of our studies on Amazon Mechanical Turk (AMT). The worker user interface is shown in Figure 3. In each assignment, we ask the worker to select a subset of images from the image sequence and write a short story (50 to 300 words) that fits the image sequence. To ensure the human-written stories are grounded on main characters, we provide names and cropped images of at most five major characters. We retrieve the bounding boxes for each character from

the MovieNet annotations and choose the least blurry appearance of each character in the image sequence. We pose three questions to the workers. The first two questions are used to identify workers who have watched the movie from which the image sequence is taken, as they might exhibit different behaviors during story-writing. The third question is to measure the visual tellability on a 5-point Likert scale, which is used to show the effectiveness of our image sequence construction pipeline.

We also design a review form for story reviewers to judge the quality of collected stories. We ask the reviewers: 1) whether they want to approve the story; 2) if not, which requirement does it break? 3) if yes, judge the statement: *this is a good story*. on a 5-point Likert scale. The first two questions are to assure that the collected stories fulfill the requirements including: the story is grammatical, the story has high diversity and the story is visually grounded. The third question is to get judgments of the quality of the approved stories.

4.1 Pilot studies

We identify the following design questions of the crowdsourcing experiment for data collection:

1. Does the image filtering process improve the tellability of the image sequences?
2. What is the optimal number of images to provide to workers to achieve high visual tellability at a reasonable workload in one writing prompt?

We conducted two pilot studies to investigate these questions. We collect 5 stories per image sequence at most from different writers.

Pilot study 1: Effectiveness of image filtering.

The first study tests whether our image-filtering steps (see Section 3) increase the visual tellability of the extracted image sequences. We extract 180 movie scenes containing 10 images each from selected movies; on half of these, we apply our image filters, while we leave the others as is. All resulting image sequences have 5 to 10 images.

Results show that the average visual tellability score of image sequences with filtering is 3.7, which is significantly higher (unpaired t -test, $t = 4.89$, p -value < 0.001) than the average visual tellability score of image sequences without filtering (3.29). This shows that our image filtering pro-

⁴<https://github.com/notAI-tech/NudeNet/>

cess in the image sequence construction pipeline leads to higher visual tellability and we will apply image filtering in our data collection.

Pilot study 2: Number of images to display.

The second study explores the effect of the number of images K in a writing prompt on workload and visual tellability. We randomly sample 150 movie scenes with 20 images, where writers can choose from 5 to 20 images for their stories. We set the minimum number of images to 5 because the most common narrative structure is *5-part play* that contains five components (Cohn, 2013). In addition, since there are five sentences per story in both ROCStories and VIST datasets, we can make the stories with 5 images in our dataset comparable to theirs. We set the maximum number to 20 because we find in a preliminary experiment that the workload of writing prompts with more than 20 images is too high considering our budget. We then run our study on these scenes.

We find a negative correlation between the actual number of images used by storytellers and the visual tellability scores, $r(500) = -0.17$, $p < 0.001$. This result indicates that showing fewer images can both improve visual tellability and reduce workload. However, we also want to obtain longer stories, we prefer to have a K larger than 5. Since a majority of 89% of the human-written stories use 5 to 10 images out of 20 and achieve a reasonably high average visual tellability (3.75), we set the maximum number of images we display to 10.

5 Data Collection

In this section, we describe how we collect and process the stories in the VWP dataset. Our goal is to obtain narratives given the curated image sequences as writing prompts.

Worker Qualification. In order to improve story quality, we apply a qualification process to workers. We first collect 4.5K stories together with visual tellability judgments and obtain 556 candidate workers. Each story is reviewed by one of five graduate students. To ensure that the reviewers mutually understand the purpose of the task, we let the reviewers judge 100 stories then check the reviews together to agree on the judgment standards. We then select 58 qualified workers with an acceptance rate $\geq 90\%$, average story quality

> 3.1 , and accepted assignments ≥ 5 . We assign a qualification to these workers and invite them to bulk collection.

Bulk Collection. We collect 7.5K stories with the qualified workers in bulk collection. We group about 300 image sequences into a batch and collect 1.5K stories per batch. For each batch, we sample s stories from each worker and review the stories to update the assessment of the worker,

$$s = \begin{cases} 10, & \text{if } n_w < 10 \\ 10 \log n_w, & \text{otherwise} \end{cases}$$

where n_w is the number of stories that worker w wrote in this batch. We run the bulk Collection batch by batch and revoke the qualification if the worker does not satisfy the selection criteria anymore.

Text Processing. We process the raw text to make it easier for training story generation models. We tokenize all stories with the spaCy⁵ English tokenizer. We then recognize all entities using a Name Entity Recognition model (Peters et al., 2017). We change all location names to placeholders and replace all named characters in each story to $[male0], \dots, [maleM], [female0], \dots, [femaleN]$. We obtain the gender of each named person based on a name statistics⁶. Finally, to mark the alignment between images and story sections, we add a special separator token $[sent]$. We randomly sample 849 stories as *validation* split and 586 stories as *test* split.

5.1 Statistics of the Dataset

We present statistics, automatic measures of coherence and diversity of our dataset to show that our collected stories are more coherent and diverse.

Statistics. We compare the properties of our dataset to similar previous datasets including Travel blogs (Park and Kim, 2015) and VIST (Huang et al., 2016) in Table 1. Our VWP dataset has 1965 image sequences with 20763 unique images from 122 movies. Each image sequence has 5 to 10 images. Our stories have 45% more tokens, 103% more events, and 285% more characters per text compared to the VIST dataset. While

⁵<https://spacy.io/>

⁶<https://ssa.gov/oact/babynames/names.zip>

Name	Image Genre	# Text	# Image per Text	# token per Text	# Event per Text	# Char. per Text
VIST	photos	50 K	5	57.6	6.3	3.4
Travel blogs	photos	10 K	1	222.3‡	3.8‡	2.3‡
VWP (Ours)	movie shots	12 K	[5, 10]	83.7	12.8	13.1

Table 1: Comparison of statistics of VWP against previous datasets. Numbers with ‡ are obtained from a small sample of the Disney split of the dataset that is available in their repository.

Dataset	# stories	LL	Avg. LL
VIST	4987	-4017	-0.8055
VWP (Ours)	4680	-3722*	-0.7953*

Table 2: Coherence by log-likelihood (LL) and average log-likelihood (Avg. LL) on validation split of VIST versus a sample split from our VWP dataset with the same number of image sequences. The stories are more coherent if the number is larger.

the Travel blogs dataset has longer stories, it has only one image per story.

Coherence. We first analyze coherence of the stories focusing on the characters and their appearances. According to Centering theory (Grosz et al., 1995), coherent narratives are typically structured such that salient entities often appear in strong grammatical roles like subject or object. As a result, we apply a model based on this theory, Entity Grid (Lapata et al., 2005), to measure the local coherence of our dataset. We apply the generative Entity Grid model implemented in the Cohere toolkit (Smith et al., 2016) on the VIST and our dataset. We calculate the log-likelihood based on entity transitions as the story coherence. The results in Table 2 show that our dataset is significantly more coherent compared to the VIST (unpaired t -test, $t = -5$, p -value < 0.001).

To further check whether event elements are semantically related given the same image sequence, we also compute the average Jaccard similarities between event elements of the stories for each image sequence by main characters, predicates (without auxiliary verbs), and arguments in different semantic roles. We identify the main characters in the raw text using coreference clusters (Lee et al., 2018). To ensure that characters mentioned only once in the story can be detected by the coreference resolution model, we append the stories with one introductory sentence per character. For example, to identify the character *Jack* in Figure 1, we add “*This is Jack.*” before the story. The Jac-

card similarity between story A and B is defined as $J(A, B) = \frac{A \cap B}{A \cup B}$, where A, B are the token sets of predicate/argument in story A and B. The results in Table 3 show that the event elements of stories conditioned on the same image sequence are more semantically related to each other. Our dataset has higher semantic cohesion compared to VIST dataset.

Diversity. We then measure diversity of the stories from two perspectives: 1) If a story has a plot with a series of different events, it must have diverse events instead of just repeating one event; 2) If these events are combined into different n-grams in the plot, then the story must have diverse predicate n-grams. For example, in the last column in Figure 1, the character *Will* has a predicate trigram (*tell, convince, work*), which is different from the next trigram (*convince, work, call*).

For event diversity, we follow Fan et al. (2019) and Goldfarb-Tarrant et al. (2020) to obtain the unique number of verbs, the verb-vocabulary ratio, verb-token ratio and the percentage of diverse verbs (not in the top 5 most frequent verbs). The results in Table 4 show that our dataset has higher event diversity than VIST across all measures. To measure predicate n-gram diversity, we extract and lemmatize verbs obtained from a Semantic Role Labeling model (Shi and Lin, 2019) and calculate the unique:total ratios of predicate unigram, bigram, and trigram (Table 4). We observe that the event sequences in VWP are more diverse than VIST, because VWP has a lower unigram ratio but higher bigram and trigram ratios.

Visual Groundedness. To check the visual groundedness of the stories, we first apply a semantic role labeller to 25 human-written stories each from VWP and VIST. We obtain 299 events and 715 arguments from the VWP samples, 84 events and 196 arguments from the VIST samples. We then manually annotated these events and arguments with 3 labels: 1) *Grounded* means the

Dataset	#	PRD	Characters	Arguments	arg0	arg1	arg2	arg-loc
VIST	998	0.063	0.184	0.055	0.041	0.018	0.018	0.013
VWP (Ours)	1000	0.068	0.21	0.057	0.101	0.048	0.025	0.017

Table 3: Average Jaccard similarity between stories of each image sequence. All numbers are higher the better except the first column which is the number of image sequences.

Dataset	Voc	Verb	Verb : Voc %	Verb : Tok %	Diverse Verb %	unigram	bigram	trigram
VIST	12627	3447	27.3	1.2	73.6	3.39	33.48	75.22
VWP (Ours)	13637	4811	35.28	1.23	79	2.71	34.87	79.10

Table 4: Comparison of diversity. The first column shows the names of the datasets. The next five columns show event diversity for validation split of VIST versus a comparable sample of VWP. We report measures including the vocabulary size (Voc), unique number of verbs (Verb), verb-vocabulary ratio (Verb : Voc %), verb-token ratio (Verb : Tok %) and percentage of diverse verbs (Diverse Verb %). The last three columns show predicate n-grams diversity for VIST versus VWP dataset. We measure diversity using unique:total ratios of predicate unigram, bigram and trigram. All numbers are the higher the better.

Label	VWP		VIST	
	#	%	#	%
E Grounded	164	54.9	38	45.2
E Inferred	134	44.8	39	46.4
E Halluciated	1	0.3	7	8.3
A Grounded	447	62.5	105	53.6
A Inferred	254	35.5	64	32.7
A Halluciated	14	2.0	27	13.8

Table 5: Visual Groundedness of stories for VIST versus VWP dataset. We report counts and percentages of each label in each data. E means event and A means argument.

event or argument is in the corresponding image; 2) *Inferred* means not in the image, but can be inferred; 3) *Halluciated* means not in image and cannot be inferred.

The results in Table 5 show that about 55% of the events and 63% of the arguments in VWP stories appear in images which are higher than the 45% of the events and 54% of the arguments in VIST stories that appear in images. About 45% of the events and 35% of the arguments in VWP stories are not in the images but can be inferred from the other images or the previous part of the stories, which are very similar to the results of VIST stories (46% of the events and 33% of the arguments not in images but can be inferred). Only 2% of the arguments in VWP stories are not in the images and cannot be inferred (i.e. not visually grounded). However, there are 8% of the events and 14% of arguments are not visually grounded

in VIST.

6 Experiment and Evaluation

In this section, we propose a strong baseline and show the experimental results on our VWP dataset. Our goal is to demonstrate the usefulness of our dataset.

We extract features for all images with Swin Transformer (Liu et al., 2021), a state-of-the-art computer vision backbone model where all parameters are fixed. We use their official model checkpoint, pre-trained on the ImageNet-21K dataset, to increase domain generality. We extract three different visual features:

1. Global features (global) are most commonly used in image-based language generation. We extract global features from the output of the last feedforward layer.

2. Object features (obj) are widely used in image-based language generation. Since *person* is also a label in object detection (Lin et al., 2014), using object features is a proper baseline for character features. We obtain object features using a Cascade Mask R-CNN object detector (Cai and Vasconcelos, 2019) with the same Swin Transformer backbone. We crop the bounding boxes of the top 20 objects that the detector predicts for each image and extract the features the same way as global features.

3. Character features (char) are extracted by cropping out the least blurry instance of each character

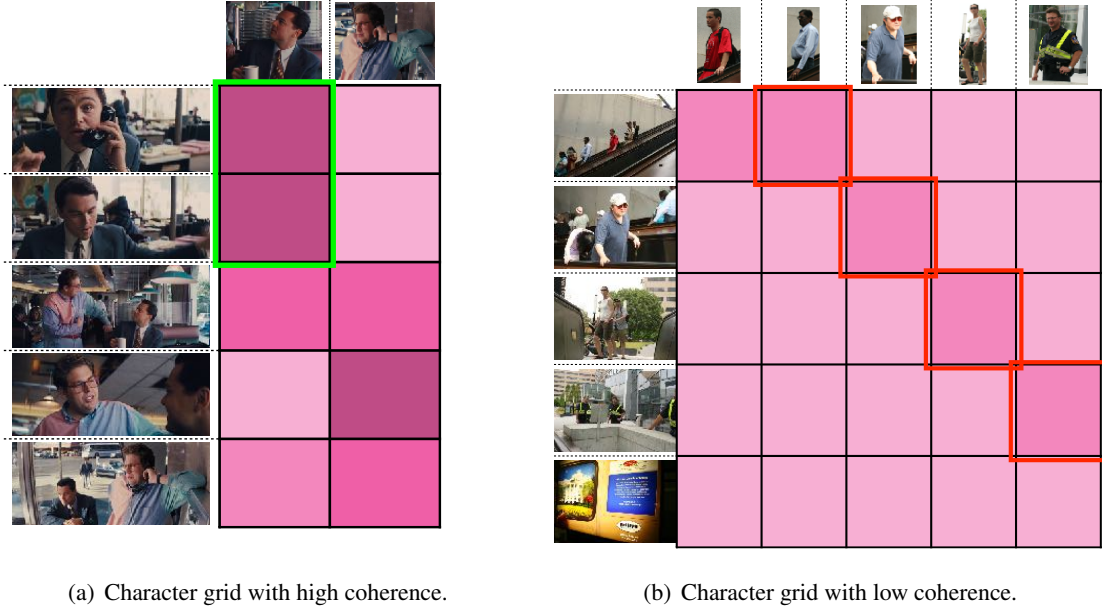


Figure 4: Example of character grid representations. Each row represents an image and each column represents a character. Shades of the cells indicate the similarities between the character features and the image features. The darker colour represents higher similarity. The green square shows the pattern that indicates high coherence and the red square shows the pattern that indicates low coherence.

using bounding boxes from our dataset. We feed the bounding boxes to the same Swin Transformer backbone and get the features from the last feed-forward layer.

We use the following models for visual story generation as baselines:

GPT-2. (GPT-2; Radford et al., 2019) is a Transformer-based language model pre-trained on large-scale text. We use the small version which is widely used in previous work of story generation.

TAPM. (TAPM; Yu et al., 2021) is a Transformer-based model which adapts the visual features with pre-trained GPT-2. This is the current state-of-the-art model for visual story generation.

For each baseline, we consider four different variants with different input: 1) only global image features; 2) global features and object features; 3) global features and character features; 4) all three available features.

6.1 Character-based visual story generation

We propose the character-grid transformer model (CharGrid) as a strong baseline to show the importance of modeling coherence and visual groundedness. We hypothesize that characters and different instances of them in image sequences play

an important role in visual story generation models in two dimensions: firstly, explicit character representations can improve groundedness of generated stories, which has been observed in textual stories (Clark et al., 2018). Secondly, representations that describe different instances of characters across images are beneficial to image-based story generation models.

Character Grid. To represent coherence of image sequences, we proposed a novel visual representation, *character grid*. As we mentioned in Section 5.1, one of the most effective methods to measure text coherence is Entity Grid, a matrix of sentences by entities where the cells are the grammatical roles of the entities in the sentence context (Lapata et al., 2005). The contribution of an entity’s mention to the coherence of a sentence is defined by its within-sentence grammatical role.

Inspired by this, we measure the narrative importance of a character in an image by the similarity between global image features and the character’s features. We thus measure the coherence of an image sequence using a matrix C of images by character instances shown in Figure 4. We obtain the narrative importance of each character instance by computing the dot product of each character’s features and the corresponding global image fea-

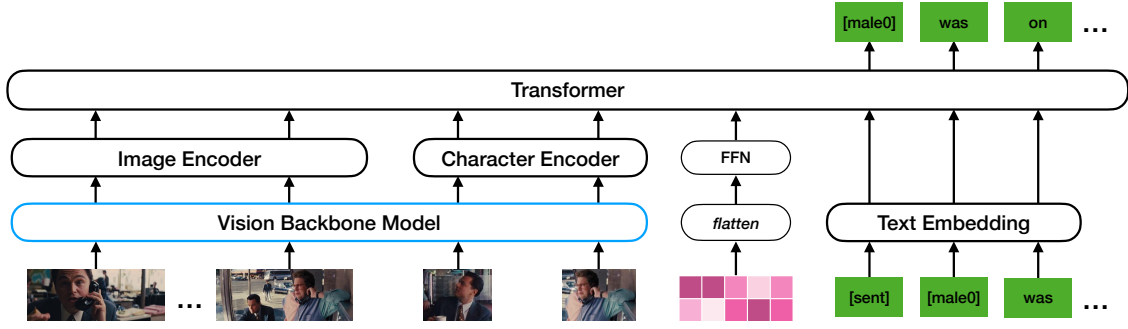


Figure 5: Architecture of character-grid transformer. The blue circles are pre-trained components where the parameters are fixed.

tures. In the character grid \mathbf{C} , each element is computed as $c_{ab} = \mathbf{i}_a \cdot \mathbf{l}_b$, where \mathbf{i}_a is the global features of image a , and \mathbf{l}_b is the features of character b .

Model Architecture. As we show in Figure 5, the architecture is based on the Transformer model. The input to the Transformer is a sequence of tokenized features including: global images features, character features, character grid, and text features. Global images features and character features are the same as the features for baseline models described above, which are first fed to trainable global and character encoders that consists of a feedforward layer. Text features are tokenized representations of the generated context, which are presented to the model incrementally. The character grid is flattened and fed to a feedforward layer. The four inputs then pass through the transformer module. The output obtained at each time step is a probability distribution over all possible output tokens from a pre-trained GPT-2 tokenizer (Wolf et al., 2020).

We also construct two variants of our model to inspect the contributions of each design decision. We replace the character features with object features to obtain the object-grid transformer model (ObjGrid). We use both character features and object features to obtain the entity-grid transformer model (EntiGrid).

Model Training. We randomly initialized the model parameters except for the vision backbone model. We optimize the model by maximizing the likelihood of the image sequence-story pairs in the training set. The parameters are updated via back propagation. We employ Nucleus sampling (Holtzman et al., 2019) to obtain the full output se-

quence for validation. We compute the METEOR score (Banerjee and Lavie, 2005) on the validation set after each training epoch. If the current epoch gets a lower METEOR score, we consider the current epoch as the best epoch and generate stories and run automatic metrics on the test set. We choose the METEOR score following previous work in visual story generation (see Section 2). In addition, Huang et al. (2016) found METEOR correlates better with human judgement than BLEU and Skip-Thoughts similarity on the VIST dataset.

6.2 Reference-based metrics

Our goal is to show the effectiveness of character grid representations. Although it has been shown that reference-based metrics correlate poorly with human judgements in open-ended language generation tasks (Guan and Huang, 2020; Gehrmann et al., 2021), it is still efficient to use them for comparison across many different models. Furthermore, we want to make our results comparable to the original results of the state-of-the-art model TAPM (Yu et al., 2021). They applied greedy search to generate stories with their models for testing and reported reference-based metrics. We thus follow the same setting and compare our proposed CharGrid model against several previous baselines.

We train all the models for at most 15 epochs with 3 different random seeds. We apply the reference-based metrics including unigram (B-1), bigram (B-2), trigram (B-3), and 4-gram (B-4) BLEU scores (B; Papineni et al., 2002), METEOR (M; Banerjee and Lavie, 2005), ROUGE-L (R; Lin, 2004), and CIDEr (C; Vedantam et al., 2015), which were used in the visual storytelling shared task (Mitchell et al., 2018). We then report the mean and standard deviation of 3 runs.

Model	Features	B-1	B-2	B-3	B-4	M	R-L	C
GPT-2	global	38.65**	20.28**	9.78**	4.68*	31.64**	24.24+	1.66**
GPT-2 + obj	global, obj	40.65**	21.35**	10.2**	4.87*	31.69**	24.05+	1.85**
GPT-2 + char	global, char	39.95**	21.04**	10.11**	4.92+	31.85*	24.19+	1.57**
GPT-2 + obj,char	global, obj, char	40.41**	21.44**	10.56**	5.06+	32.03*	24.38	1.87**
TAPM	global	39.85**	21.7**	10.72**	5.19	32.38+	25.09	1.48**
TAPM + obj	global, obj	40.86**	22.13**	10.83**	5.25	32.34+	24.91	1.82**
TAPM + char	global, char	40.03**	21.68**	10.66**	5.18	32.42+	24.88	1.4**
TAPM + obj,char	global, obj, char	40.87**	21.99**	10.72**	5.06+	32.48+	24.87	1.59**
<i>Ours</i>								
ObjGrid	global, obj	47.66	25.26	11.95	5.42	32.83	24.42	4.68
EntityGrid	global, obj, char	45.83	24.85	12.11	5.7	32.68	24.89	3.53+
CharGrid	global, char	47.71	25.33	11.95	5.42	33.03	25.01	4.83

Table 6: Results of all models using different input features on the test set of VWP using reference-based metrics including BLEU (B), METEOR (M), ROUGE-L (R-L), and CIDEr (C). All numbers are average of three runs with different random seeds. (pretrain) indicates models initialised the Transformer with GPT-2 pre-trained weights. +, * and ** represent that the number is one, two or three standard deviations away from the mean of CharGrid model.

Results in Table 6 show that the character-grid transformer model (CharGrid) driven by visual coherence outperforms TAPM with character features (TAPM + char) significantly on BLEU-1/2/3 and CIDEr and marginally on METEOR. CharGrid model also outperforms GPT-2 with character features (GPT-2 + char) significantly on most metrics except marginally on BLEU-4 and METEOR. The object-grid transformer model (ObjGrid) outperforms TAPM with object features (TAPM + obj) significantly on BLEU-1/2/3 and CIDEr and marginally on METEOR. ObjGrid model also outperforms GPT-2 with object features (GPT-2 + obj) significantly on most metrics except marginally on BLEU-4. The entity-grid transformer model (EntiGrid) outperforms TAPM with all features (TAPM + obj,char) significantly on most metrics except marginally on METEOR and ROUGE-L. EntiGrid model also outperforms GPT-2 with all features (GPT-2 + obj,char) on most metrics except BLEU-4. These results show the effectiveness of character/object/entity grid representations for coherence of image sequences.

6.3 Human evaluation

Because story generation is an open-domain task, reference-based metrics can only show how output stories match with the references. In order to measure the quality of generated stories directly, we conduct a crowdsourcing experiment to obtain human binary judgments between two systems. We design the first question for *Grammaticality*, which measures whether the textual outputs are at least grammatical and sets a foundation for other

metrics. We then design questions for two properties that we identified for good textual stories: *Coherence* and *Diversity*. Finally, we ask a question to compare the *Visual Groundedness* in order to make sure that the stories are relevant to the input image sequence.

We conduct the experiment with 28 crowd workers over 50 pairs of stories and report the percentage of the judgments for each system that annotators are in favor of. To make the stories more readable, we change the generated character placeholders to randomly sampled names. The results in table 7 show that TAPM with character features (TAPM + char) outperforms TAPM in Visual Groundedness significantly. CharGrid outperforms TAPM + char on all metrics significantly. We use two-sided binomial tests. This indicates that our character grid representation yields better stories. These results confirm the findings in the evaluation with reference-based metrics.

6.4 Qualitative evaluation

We also conduct a qualitative evaluation to show that stories generated by TAPM with character features are more visually grounded than without character features and character grid representation further improves the coherence and visual groundedness. To obtain more diverse text, we use Nucleus Sampling (Holtzman et al., 2019) with $p = 0.1$ on all models to generate the stories. As in Figure 6, TAPM generates unreasonable noun phrases *the train*. With character features, TAPM + char is able to explore character-object interac-

Model	Grammatical	Coherence	Visual Groundedness	Diversity
TAPM + char vs. TAPM	+2.45	+1.99	+3.99*	+1.69
CharGrid vs. TAPM + char	+6.49**	+8.41**	+6.25*	+11.06**

Table 7: Human binary judgments in percentage of generated stories between TAPM and TAPM with character features (TAPM + char), TAPM + char and our model (CharGrid) on test set of VWP across four criteria: *Grammaticality*, *Coherence*, *Visually Groundedness* and *Diversity*. The numbers are percentages. * means p -value < 0.05 . ** means p -value < 0.01 .


Image sequence	
	
Human	It was a sunny day in the city. Inside the large car parking garage, a group of robbers approached a car and the driver. They were out to steal the driver's car and money. The driver noticed the men approaching them. He started to show fear since he didn't know what they wanted.
TAPM	the train was on its way to the city . the police were looking for the suspect . Adam was a bit nervous . Adam told him to get out of his car . Andrew was shocked and told Adam that he was going to have to go and kill the man .
TAPM + Char (Ours)	Greyson was walking down a street . they were looking for a suspect . Greyson was surprised to see that the car was a gun and he was a gun expert . he had never seen such a thing before . Greyson shot him in the chest and he fell to the ground . Greyson decided that he was going to shoot the man .
CharGrid (Ours)	the building is a huge skyscraper that is being built in [location] . Henry and Charles were in the car waiting for the train to arrive . Henry was very scared and did n't know where he was . Henry is driving the car and is ready to shoot at any moment . Charles was in a lot of pain and was trying to get away .

Figure 6: Qualitative results of generated and human-written stories. Red colour represents errors made by models and green colour indicates better output.

tion and reason that there is no train in the image. So it generates more reasonable terms *a street*.

However, TAPM + char model fails to represent the relations between characters, TAPM + char generates pronoun *they* without introducing characters in the second image. In contrast, CharGrid introduces two new characters correctly.

7 Conclusions and Future Work

We show that curated image sequences with characters are effective as writing prompts for visual story generation in both data collection and model design. By filtering images without any objects and removing highly similar images to boost diversity, we can improve visual tellability of image sequences. Presenting selected characters during the story-writing yields stories with characters grounded in images, which are more coherent and have more narrativity. Correspondingly, using character features as input to the story generation model can improve the quality of generated stories. Adding the character grid representation can bring further improvements in coherence, grammaticality, visual groundedness, and narrativity.

Future work. One important property of visual narratives not covered in this work is *narrativity* (Piper et al., 2021), i.e. whether an image sequence contains necessary narrative structures to be a nar-

rative. A narrative structure can be achieved by events following a typical order with roles like *Establisher*, *Initial*, *Initial Peak* and *Release* (Cohn, 2013). We observe that these roles of events emerge in our collected stories. Our annotation of different instances of the same character across a story allow us to construct event chains for each character. Future work should investigate how to annotate roles of these events, measure narrativity and build a model to generate stories with higher narrativity.

A major assumptions of all previous work in storytelling is that all humans are equally and reasonably proficient in story-writing and can translate visual narratives into textual narratives. However, individual differences in the writing proficiency of humans must have an impact on story quality. How to explore this from the perspective of both data selection and model design would be an interesting future direction to take.

References

Nader Akoury, Shufan Wang, Josh Whiting, Stephen Hood, Nanyun Peng, and Mohit Iyyer. 2020. Storium: A dataset and evaluation platform for machine-in-the-loop story generation. In *Proceedings of the 2020 Conference on Em-*

- pirical Methods in Natural Language Processing (EMNLP)*, pages 6470–6484.
- Satanjeev Banerjee and Alon Lavie. 2005. [ME-TEOR: An automatic metric for MT evaluation with improved correlation with human judgments](#). In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Zhaowei Cai and Nuno Vasconcelos. 2019. Cascade r-cnn: high quality object detection and instance segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 43(5):1483–1498.
- Khyathi Chandu, Eric Nyberg, and Alan W Black. 2019. Storyboarding of recipes: Grounded contextual generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6040–6046.
- Elizabeth Clark, Yangfeng Ji, and Noah A Smith. 2018. Neural text generation in stories using entity representations as context. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2250–2260.
- Neil Cohn. 2013. Visual narrative structure. *Cognitive science*, 37(3):413–452.
- Neil Cohn. 2020. Visual narrative comprehension: Universal or not? *Psychonomic Bulletin & Review*, 27(2):266–285.
- Neil Cohn, Martin Paczynski, Ray Jackendoff, Phillip J Holcomb, and Gina R Kuperberg. 2012. (pea) nuts and bolts of visual narrative: Structure and meaning in sequential image comprehension. *Cognitive psychology*, 65(1):1–38.
- 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.
- Sebastian Gehrmann, Tosin Adewumi, Karmanya Aggarwal, Pawan Sasanka Ammanamanchi, Anuoluwapo Aremu, Antoine Bosselut, Khyathi Raghavi Chandu, Miruna-Adriana Clinciu, Dipanjan Das, Kaustubh Dhole, Wanyu Du, Esin Durmus, Ondřej Dušek, Chris Chinenye Emezue, Varun Gangal, Cristina Garbacea, Tatsunori Hashimoto, Yufang Hou, Yacine Jernite, Harsh Jhamtani, Yangfeng Ji, Shailza Jolly, Mihir Kale, Dhruv Kumar, Faisal Ladhak, Aman Madaan, Mounica Maddela, Khyati Mahajan, Saad Mahamood, Bodhisattwa Prasad Majumder, Pedro Henrique Martins, Angelina McMillan-Major, Simon Mille, Emiel van Miltenburg, Moin Nadeem, Shashi Narayan, Vitaly Nikolaev, Andre Niyongabo Rubungo, Salomey Osei, Ankur Parikh, Laura Perez-Beltrachini, Niranjan Ramesh Rao, Vikas Rautnak, Juan Diego Rodriguez, Sashank Santhanam, João Sedoc, Thibault Sellam, Samira Shaikh, Anastasia Shimorina, Marco Antonio Sobrevilla Cabezudo, Hendrik Strobelt, Nishant Subramani, Wei Xu, Diyi Yang, Akhila Yerukola, and Jiawei Zhou. 2021. [The GEM benchmark: Natural language generation, its evaluation and metrics](#). In *Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)*, pages 96–120, Online. Association for Computational Linguistics.
- Seraphina Goldfarb-Tarrant, Tuhin Chakrabarty, Ralph Weischedel, and Nanyun Peng. 2020. Content planning for neural story generation with aristotelian rescoring. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4319–4338.
- Barbara J Grosz, Aravind K Joshi, and Scott Weinstein. 1995. Centering: A framework for modelling the local coherence of discourse.
- Jian Guan and Minlie Huang. 2020. Union: An unreferenced metric for evaluating open-ended story generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9157–9166.

- Stevan Harnad. 1990. The symbol grounding problem. *Physica D: Nonlinear Phenomena*, 42(1-3):335–346.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. In *International Conference on Learning Representations*.
- Qingqiu Huang, Yu Xiong, Anyi Rao, Jiaze Wang, and Dahua Lin. 2020. Movienet: A holistic dataset for movie understanding. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 709–727. Springer.
- Ting-Hao Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. 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.
- Peter Hühn, Jan Christoph Meister, John Pier, and Wolf Schmid. 2014. *Handbook of narratology*. Walter de Gruyter GmbH & Co KG.
- Mirella Lapata, Regina Barzilay, et al. 2005. Automatic evaluation of text coherence: Models and representations. In *IJCAI*, volume 5, pages 1085–1090. Citeseer.
- Kenton Lee, Luheng He, and L. Zettlemoyer. 2018. Higher-order coreference resolution with coarse-to-fine inference. In *NAACL-HLT*.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint arXiv:2103.14030*.
- Margaret Mitchell, Ting-Hao ‘Kenneth’ Huang, Francis Ferraro, and Ishan Misra, editors. 2018. [Proceedings of the First Workshop on Storytelling](#). Association for Computational Linguistics, New Orleans, Louisiana.
- Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. A corpus and cloze evaluation for deeper understanding of commonsense stories. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 839–849.
- 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, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Cesc C Park and Gunhee Kim. 2015. Expressing an image stream with a sequence of natural sentences. *Advances in neural information processing systems*, 28:73–81.
- José Luis Pech-Pacheco, Gabriel Cristóbal, Jesús Chamorro-Martínez, and Joaquín Fernández-Valdivia. 2000. Diatom autofocusing in bright-field microscopy: a comparative study. In *Proceedings 15th International Conference on Pattern Recognition. ICPR-2000*, volume 3, pages 314–317. IEEE.
- Matthew E. Peters, Waleed Ammar, Chandra Bhagavatula, and R. Power. 2017. Semi-supervised sequence tagging with bidirectional language models. In *ACL*.
- Andrew Piper, Richard Jean So, and David Bamman. 2021. [Narrative theory for computational narrative understanding](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 298–311, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Lianhui Qin, Antoine Bosselut, Ari Holtzman,

- Chandra Bhagavatula, Elizabeth Clark, and Yejin Choi. 2019. Counterfactual story reasoning and generation. 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 5043–5053.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Hannah Rashkin, Antoine Bosselut, Maarten Sap, Kevin Knight, and Yejin Choi. 2018. 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, Asli Celikyilmaz, Yejin Choi, and Jianfeng Gao. 2020. Plotmachines: Outline-conditioned generation with dynamic plot state tracking. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4274–4295.
- Peng Shi and Jimmy Lin. 2019. Simple bert models for relation extraction and semantic role labeling. *arXiv preprint arXiv:1904.05255*.
- Karin Sim Smith, Wilker Aziz, and Lucia Specia. 2016. Cohere: A toolkit for local coherence. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 4111–4114.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4566–4575.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45.
- Yu Xiong, Qingqiu Huang, Lingfeng Guo, Hang Zhou, Bolei Zhou, and Dahua Lin. 2019. A graph-based framework to bridge movies and synopses. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4592–4601.
- Peng Xu, Mostofa Patwary, Mohammad Shoeybi, Raul Puri, Pascale Fung, Anima Anandkumar, and Bryan Catanzaro. 2020. [MEGATRON-CNTRL: Controllable story generation with external knowledge using large-scale language models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2831–2845, Online. 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.
- Youngjae Yu, Jiwan Chung, Heeseung Yun, Jongseok Kim, and Gunhee Kim. 2021. Transitional adaptation of pretrained models for visual storytelling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12658–12668.