

On the Challenges in Evaluating Visually Grounded Stories

Aditya K Surikuchi Raquel Fernández Sandro Pezzelle



UNIVERSITY OF AMSTERDAM
Institute for Logic, Language and Computation

Visually Grounded Story Generation



Visually Grounded Story Generation



S1. *A man ducks for cover. There are bullet holes in the wall near him, and he covers his face in fear. He tries to get away, but another person grabs him from behind. He holds him in place as the man struggles to get away. Suddenly, the bullets start flying. Both men duck or dive for safety. Food and bits of plaster go flying as the bullets fill the room. Two men are hit, and they both go flying back.*

Visually Grounded Story Generation



S1. *A man ducks for cover. There are bullet holes in the wall near him, and he covers his face in fear. He tries to get away, but another person grabs him from behind. He holds him in place as the man struggles to get away. Suddenly, the bullets start flying. Both men duck or dive for safety. Food and bits of plaster go flying as the bullets fill the room. Two men are hit, and they both go flying back.*

S2. *In the opulent room, a man in a tuxedo was playing the accordion, oblivious to the chaos unfolding around him. Suddenly, another man burst into the room, his face contorted with anger as he lunged at the first man. The second man's attack was met with a swift counter from the first, who managed to push him back. As they struggled, a cloud of dust filled the air, obscuring their movements. In the midst of the commotion, a third man appeared, his expression one of shock and fear as he watched the two men grapple.*

Visually Grounded Story Generation



🗣 A man ducks for cover. There are bullet holes in the wall near him, and he covers his face in fear. He tries to get away, but another person grabs him from behind. He holds him in place as the man struggles to get away. Suddenly, the bullets start flying. Both men duck or dive for safety. Food and bits of plaster go flying as the bullets fill the room. Two men are hit, and they both go flying back.

🤖 In the opulent room, a man in a tuxedo was playing the accordion, oblivious to the chaos unfolding around him. Suddenly, another man burst into the room, his face contorted with anger as he lunged at the first man. The second man's attack was met with a swift counter from the first, who managed to push him back. As they struggled, a cloud of dust filled the air, obscuring their movements. In the midst of the commotion, a third man appeared, his expression one of shock and fear as he watched the two men grapple.

Datasets

VIST

- ▶ Sequences constructed using images from Flickr albums.
- ▶ Lacks consistency of entities (e.g., *characters, objects*)
- ▶ Corresponding stories are generally descriptive in nature

VWP

- ▶ Sequences constructed using scenes from movies
- ▶ Semantically well-connected with recurring characters
- ▶ Stories contain diverse entities, are longer, and coherent

Models



Qwen-VL



LLaVA



DeepSeek-VL

General-purpose VLMs
(off-the-shelf)



TAPM (+LLAMA 2)

Specific to Visual Story
Generation

Evaluation

Human evaluation is challenging in terms of:

- 💶 Scalability and costs
- ⚠ Selecting qualified annotators and reliability

Evaluation

Human evaluation is challenging in terms of:

- 💶 Scalability and costs
- ⚠ Selecting qualified annotators and reliability

Automatic Metrics:

- ▶ BLEU, METEOR, CIDEr, SPICE, ROUGE

Evaluation

Human evaluation is challenging in terms of:

- 💶 Scalability and costs
- ⚠ Selecting qualified annotators and reliability

Automatic Metrics:

- ✗ BLEU, METEOR, CIDEr, SPICE, ROUGE
- ✓ Reference-free metrics that assess stories along different aspects—*Coherence, Visual grounding, Repetition*

Evaluation

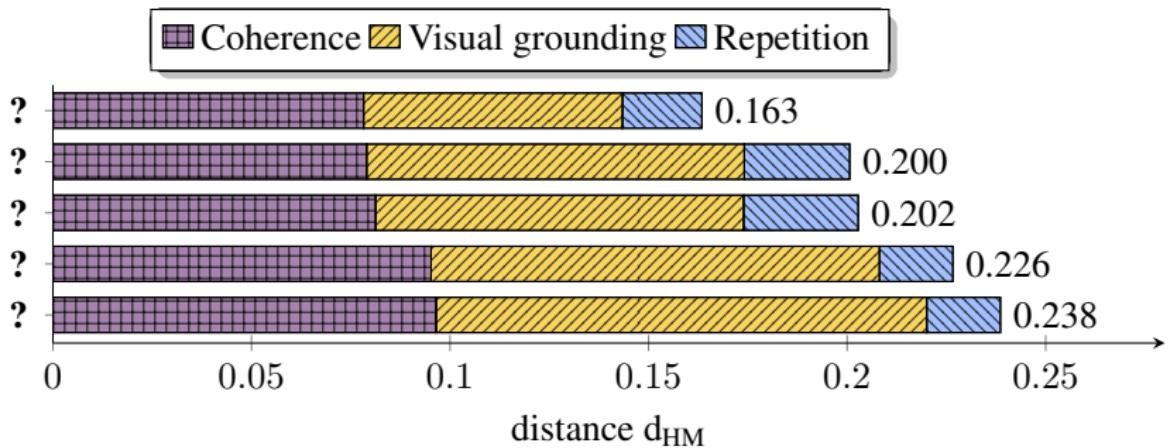
Human evaluation is challenging in terms of:

- 💶 Scalability and costs
- ⚠ Selecting qualified annotators and reliability

Automatic Metrics:

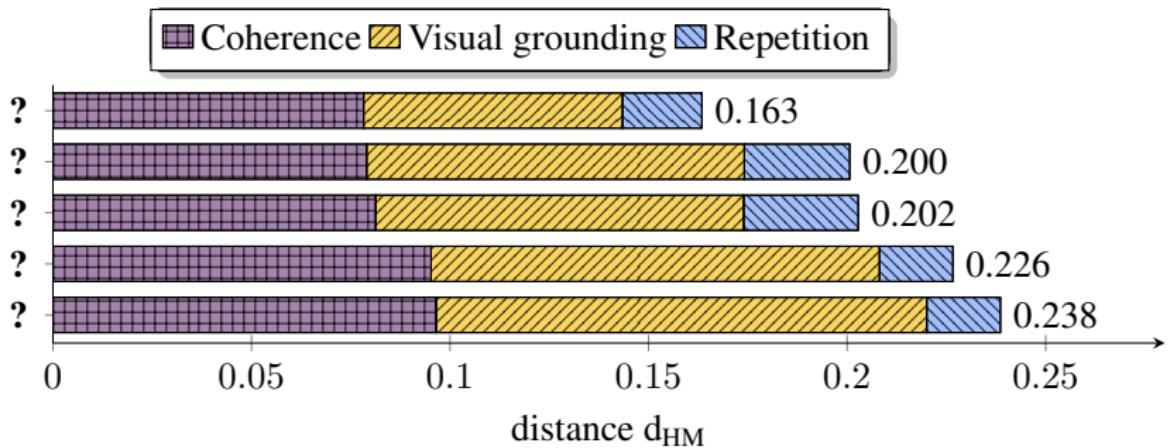
- ✗ BLEU, METEOR, CIDEr, SPICE, ROUGE
- ✓ Reference-free metrics that assess stories along different aspects—*Coherence, Visual grounding, Repetition*
- ✓ d_{HM} measure that computes distances between individual metric scores of model- and corresponding human stories

Results

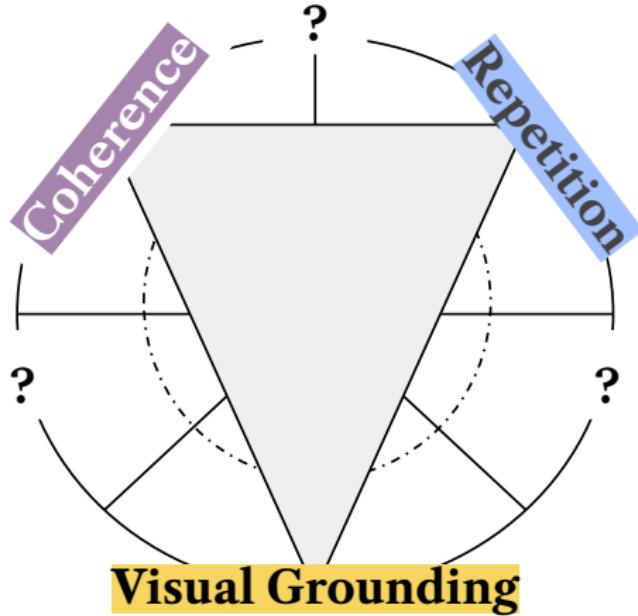


🤔 What is overall best performing model?

Results



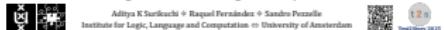
- 🤔 What is overall best performing model?
- 🤔 How did the models fare along each of the 3 dimensions?



🤔 Are there other dimensions relevant for evaluating visually-grounded stories?

To learn more, please stop by our poster:

On the Challenges in Evaluating Visually Grounded Stories



Visually Grounded Story Generation Evaluation Methods

Input: sequence of temporally-ordered images.



Task: To generate a detailed story consistent with the input.

Human-written story: A man walks in the room. There are books hidden in the wall now, and he wants to get them. He goes to get away, but there are people going around. The books are placed in a shelf. The man goes to the shelf to take the books out of the shelf. He takes the book to the safety. First and last of all, flying things are flying in the room. Two men are in it, and they both are flying back.

Annotations

VISIT²

- Sequences are constructed using images from Flickr albums.
- Looks consistency of entities (e.g., colors).
- Corresponding stories are generally described in nature.

Stories are provided by qualified workers through a crowdsourcing platform.

VISIT²+

- Sequences are constructed using images from Flickr albums.
 - Looks consistency of entities (e.g., colors).
 - Corresponding stories are generally described in nature.
- Stories are provided by qualified workers through a crowdsourcing platform.

VISIT² challenge

- VISIT² dataset with 11770 training, 649 validation, and 565 test samples.
- Stories were evaluated using reference-based metrics such as METEOR and reference-free metrics such as CIDEr (character n-gram).
- The proposed VLMs were evaluated using a metric that was generated from using other visual feature extractors besides Auto-Encoder.

In this work, we evaluate our two models on the VISIT² challenge dataset, using the recently proposed Δ_{GAP} method.

Modeling Approaches

We use 2 different types that are shown to perform well on other visual storytelling datasets: CLIP-ViT and LLaMA.

We also compare two recent general-purpose VLMs that have demonstrated strong performance on various vision-language benchmarks.

Model	Vision Encoder	Language Decoder
General purpose VLM (off-the-shelf)		
LLM(DeBERTa-V3)	CLIP-ViT-150M	Matric-TB
Qwen(1.4B)	CPT-1.4B	Qwen-1.3B
Dolly(7.5B)	T5-7.5B	Dolly-7.5B
Models specific to Visual Story Generation		
TAFM	Auto-Encoder	LLAMA 2
TAFM	ResNet-101, ResNet-101CINN	LLAMA 2

Process: [UNITS] image > Write a story using exactly [Mango] sentence for this image sequence. Do not use any more than [Mango] characters. [UNITS]

The generated VLMs were not directly provided with any of the visual storytelling datasets, including VISIT².

References

- Visual Storytelling Dataset: A Novel 100k Story Generation with Grounded Images. (Song et al., 2020)
- Visual Storytelling Generation Challenge (Song et al., 2020)
- Quality-Style: Improving VQA and visual Storytelling. (Liu et al., 2020)
- VIT2: A Large-Scale Visual Storytelling Dataset. (Xu et al., 2021)
- VIT2+: A Large-Scale Visual Storytelling Dataset. (Xu et al., 2021)
- Qwen: A Large Pre-trained Language Model for Visual Storytelling. (Song et al., 2021)
- Qwen: Learning Reference Metrics for Visual Storytelling. (Song et al., 2021)
- Matric: A Metric for Generating Object- and Visual Storytelling. (Kumar et al., 2021)

Evaluation is challenging. Possibility of several consecutive stories for a single given image sequence, including multiple ground-truths.

Cohen-Kappa-BiLSTM²: average probability with which each sentence follows the preceding sentence (coarse-grained) of the story range $[0, 1]$.

Visual grounding (GIST)²: alignment scores between scene-photographs and corresponding image descriptions, normalized and aggregated in range $[0, 1]$.

Visual grounding (GAP)²: alignment scores between scene-photographs and corresponding image descriptions, normalized and aggregated in range $[0, 1]$.

Human evaluation: human evaluation of the generated stories based on the total number of words in both texts (ground truth), for example average of inter- and intra-baseline repetitions is computed, range $[0, 1]$.

shd[unseen story]: Grounded generative model vs. Ground truth. The model can't find the bottle.

shd[seen story]: Grounded generative model vs. Ground truth. The model can find the bottle.

shd[unseen story]: Grounded generative model vs. Ground truth. The model can't find the bottle.

shd[seen story]: Grounded generative model vs. Ground truth. The model can find the bottle.



Human-written story: It was a normal evening for [Donald]. When she entered her room, she saw her friend [Donald] sitting on the bed. [Donald] didn't know what was happening. She asked him what was wrong. [Donald] told her that he was feeling sick. She asked him if he wanted to go to the hospital. [Donald] said no. She started to open his bag. [Donald] was so happy for her friend as she finally able to see her present for the day. They went to eat and celebrate their happiness.

Qwen²-S: In the cozy bedroom, Donald received her friend, Donald, looking dismally. Donald was lying on the bed, looking pale and weak. Donald's mother came in and asked what was wrong. Donald, still in shock, told Donald about a recent cold he had. Donald had a fever. Donald's mother, worried about Donald's health, decided to call a doctor. Donald's mother got a doctor's appointment for Donald's birthday. Donald, who was still feeling unwell, decided to stay home and rest. Donald's mother checked the body of the temperature, and once again, she found that Donald's temperature was still high. Donald's mother decided to give Donald some medicine. Donald's mother gave Donald some medicine and told him to take care of himself. Donald's mother also told him that Donald should get out of the bed. **It was a large and grand house with grand gardens.**

Takeways

The Δ_{GAP} evaluation metric is used for differentiating "good" stories in terms of how well they align with human evaluations.

- There are several other metrics of a visual story that need to be considered for improving evaluation accuracy.



Thanks for your attention!