# Make-A-Story: Visual Memory Conditioned Consistent Story Generation

Tanzila Rahman[1,3]     Hsin-Ying Lee[2]     Jian Ren[2]     Sergey Tulyakov[2]

Shweta Mahajan[1,3]     Leonid Sigal[1,3]

1 UBC     2 (Snapchat logo)     3 (logo)

TUE-AM-238

# Problem Formulation

Imagine if you can just tell story to your children….

Wilma and Betty are in a room. Wilma is talking to Betty and standing on a pedestal while Betty hems her dress. **She** turns to grab quills from a porcupine.

. . .

Wilma is in a room. **She** is speaking to someone while looking over her right shoulder.

. . .

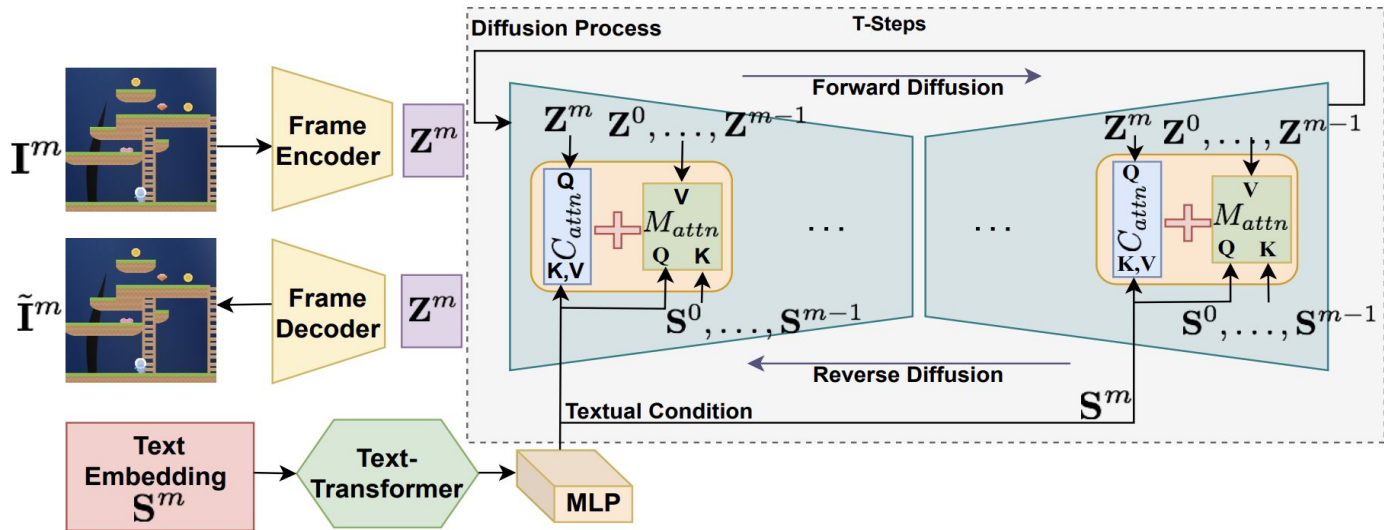Wilma and Betty are in the living room. Betty talks and then Wilma moves her head.

. . .

Wilma is standing in the room talking.

# Problem Formulation

And the story can be generated automatically by maintaining resolution of references and consistency in subject/background appearance:

Wilma and Betty are in a room. Wilma is talking to Betty and standing on a pedestal while Betty hems her dress. **She** turns to grab quills from a porcupine.

. . .

Wilma is in a room. **She** is speaking to someone while looking over her right shoulder.

. . .

Wilma and Betty are in the living room. Betty talks and then Wilma moves her head.

. . .

Wilma is standing in the room talking.

# Our Contributions

- Story-LDM: A novel autoregressive deep generative framework for the task of story generation.

- A novel memory-attention mechanism to generate consistent stories.

- Extend existing datasets to include more complex scenarios to validate co-reference resolution for character and background consistency in the visual domain.

- Novel evaluation metrics to evaluate for foreground (character) as well as background consistency.

# Story-LDM Architecture



For each frame cross-attention: $C_{attn} = \sum_i \hat{f}(\mathbf{Z}^m)_i f(\mathbf{S}^m)_i$

Memory attention module: $M_{attn} = \sum_{k=1}^{m-1} \sum_i \hat{f}(\mathbf{Z}^k)_i f(\mathbf{S}^k)_i f(\mathbf{S}^m)_i$

Final Output of the attention-module: $C_{attn} + M_{attn}$

# Dataset Statistics

| Dataset | #Ref (avg.) | #Chars | #Backgrounds |
|---|---|---|---|
| MUGEN [1] | None | 1 | 2 |
| Extended MUGEN | 3 | 3 | 6 |
| FlintstonesSV [2] | 3.58 | 7 | 323 |
| Extended FlintstonesSV | 4.61 | 7 | 323 |
| PororoSV [3] | 1.01 | 9 | None |
| Extended PororoSV | 1.16 | 9 | None |

# Evaluation Metrics

- Character Classification : Measure frame accuracy (*i.e.* exact matching of characters) and F1-score by using fine-tuned Inception-v3.

- Background Classification: Measure the correspondence of the background to the ground-truth and consider F1-score as a measure of quality.

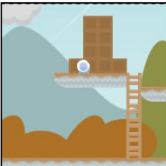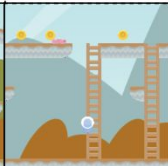- Frechet Inception Distance (FID): Measure the distance between feature vectors from real and generated images.

# Quantitative Results

| Dataset | Method | w/ref. text | Char-acc (↑) | Char-F1 (↑) | BG-acc (↑) | BG-F1 (↑) | FID (↓) |
|---|---|---|---|---|---|---|---|
| FlintstonesSV | VLCStoryGAN [4] | x | 27.73 | 42.01 | 4.83 | 16.49 | 120.85 |
| | LDM [5] | x | 79.86 | 92.33 | 48.02 | 37.86 | 61.40 |
| | LDM [5] | ✓ | 57.38 | 78.68 | 44.19 | 28.25 | 87.39 |
| | Story-LDM (Ours) | ✓ | 69.19 | 86.59 | 35.21 | 28.80 | 69.49 |
| PororoSV | VLCStoryGAN [4] | ✓ | 17.36 | 43.02 | - | - | 84.96 |
| | LDM [5] | ✓ | 16.59 **+22** | 56.30 | - | - | 60.23 |
| | Story-LDM (Ours) | ✓ | **20.26** | **57.95** | - | - | **36.64** |
| MUGEN | LDM [5] | ✓ | 31.39 | 21.28 | 15.74 | 18.66 | 120.99 |
| | Story-LDM (Ours) | ✓ | **93.40** **+62** | **95.60** | **92.19** | **92.37** | **62.16** |

Improvement of **22%** and **62%** in terms of character accuracy on the PororoSV and the Mugen dataset.

# Story Generation Results on MUGEN

# Story Generation Result on FlintstonesSV

# Story Generation Result on PororoSV

# Qualitative Comparison



(a) Comparison on the FlintstonesSV dataset visual **story generation**.

(b) Comparison on the FlintstonesSV dataset for **story continuation.**

# Additional qualitative results



(a) Generating different yet consistent stories by branching the storyline.

(b) Diverse outputs for a single storyline obtained with our Story-LDM.

# Conclusion

- We formulate consistent story generation in a more realistic way by co-referencing actors/backgrounds in the story descriptions.

- We introduce an autoregressive Story-LDM approach with memory attention capable of maintaining consistency across the frames.

- We expect our proposed formulation and models to be conductive to the real-world use cases and further the research.

# References

[1] Hayes, Thomas, et al. "Mugen: A playground for video-audio-text multimodal understanding and generation." ECCV 2022.

[2] Gupta, Tanmay, et al. "Imagine this! scripts to compositions to videos." ECCV 2018.

[3] Li, Yitong, et al. "Storygan: A sequential conditional gan for story visualization." CVPR 2019.

[4] Maharana, Adyasha, and Mohit Bansal. "Integrating visuospatial, linguistic and commonsense structure into story visualization." arXiv preprint arXiv:2110.10834 (2021).

[5] Rombach, Robin, et al. "High-resolution image synthesis with latent diffusion models." CVPR 2022.