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

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

Shweta Mahajan[1,3]      Leonid Sigal[1,3,4]

[1]University of British Columbia      [2]Snap Inc.

[3]Vector Institute for AI      [4]Canada CIFAR AI Chair

| Method | Reference-text | Char-acc ($\uparrow$) | Char-F1 ($\uparrow$) | BG-acc ($\uparrow$) | BG-F1 ($\uparrow$) | FID ($\downarrow$) |
|---|---|---|---|---|---|---|
| Story-LDM (Ours) | ✓ | 69.19 | 86.59 | 35.21 | **28.80** | 69.49 |
| Story-LDM (Ours) | ✗ | **83.29** | **94.61** | **35.54** | 27.32 | **64.89** |

Table 3. **Experimental results using our proposed Story-LDM with and without the reference text on the FlintstonesSV dataset.**

## 1. Additional Quantitative Results

In Tab. 3, we show the performance of our Story-LDM approach with and without the reference text. We observe that without reference text *i.e.*, with explicit mentions of the characters (by their name), our approach outperforms the prior state-of-the-art VLCStoryGAN [32] (*cf*. Tab. 2, row 1 in the main paper) and the strong LDM [42] baseline on character accuracy by nearly $\sim 3.4\%$.

This demonstrates that introduction of our Memory attention module to the pipeline of the diffusion model with U-Net architecture increases the performance even in the traditional dataset setting. This can be attributed to the fact that even when textual resolution of character names or setting is unnecessary, our memory attention module can still enhance consistency of appearance in the visual domain. When using the descriptions with references, the character accuracy of the model drops by $\sim 12\%$ showing the difficulty of the extended task of generating stories from the co-referenced text. Even the strong LDM baseline which outperforms the prior state-of-the-art for story generation, offers limited performance on the complex task of story generation with co-reference resolution illustrating the hardness and complexity of the new task (*cf*. Tab. 2 in the main paper). Our Story-LDM approach outperforms this strong baseline illustrating the benefits of our autoregressive Story-LDM with memory module.

Table 5 shows quantitative results for consistent story generation with the reference text on the PororoSV [27] dataset. We compare our proposed Story-LDM with DUCO-STORYGAN [32] and VLCStoryGAN [31]. Our method outperforms previous baseline models (including LDM baseline) in character evaluation metrics. To be noted,

PororoSV dataset has no background information, therefore we only perform character level evaluation on this dataset.

Moreover, we conduct human evaluation. We randomly select 10 examples to compare Story-LDM with LDM. Following [31], we select three evaluation criteria: visual quality, consistence, and relevance for each sample. We conducted forced choice experiment with 13 subjects. Preference rate for our Story-LDM is 74.7%, 65.4% & 69.2% in terms of listed criteria.

## 2. Additional Qualitative Results

In Fig. 9 we provide example image frames from our extended MUGEN dataset with three characters *Lisa, Tony* and *Jhon*, and six different backgrounds. Our extended MUGEN dataset is thus more complex than the MUGEN dataset in [16] where only one character is considered for two backgrounds.

Fig. 10 shows examples of positive/negative samples. To evaluate fore- and back-ground consistency we propose two evaluation metrics: character and background classification. The generated story is considered accurate (positive sample) if it appropriately resolves character and background references in generation, otherwise it is considered a negative sample. In the figure, the 1st row (*i.e.* positive sample) has 100% character and background accuracy; the 2nd row (*i.e.* negative) has 75% and 50% character and background accuracy.

Figs. 11, 12 and 13 show random samples obtained for story generation with our Story-LDM approach on the MUGEN, FlintstonesSV and PororoSV datasets respectively. Clearly, our approach yields high quality frames with character and background consistency.

| Dataset | Method | w/ ref. text | Char-acc (↑) | Char-F1 (↑) | BG-acc (↑) | BG-F1 (↑) | FID (↓) |
|---|---|---|---|---|---|---|---|
| PororoSV | DUCO-STORYGAN [32] | ✓ | 13.97 | 38.01 | - | - | 96.51 |
| | VLCStoryGAN [31] | ✓ | 17.36 | 43.02 | - | - | 84.96 |
| | LDM [42] | ✓ | 16.59 | 56.30 | - | - | 60.23 |
| | Story-LDM (Ours) | ✓ | **20.26** | **57.95** | - | - | **36.64** |

Table 4. **Quantitative results.** Experimental results on the PororoSV datasets.



Figure 9. **Combination of different characters and backgrounds from the MUGEN dataset [16].**

| Dataset | Method | w/ ref. text | Char-acc (↑) | Char-F1 (↑) | BG-acc (↑) | BG-F1 (↑) | FID (↓) |
|---|---|---|---|---|---|---|---|
| Flintstones | VLCStoryGAN [31] | × | 27.73 | 42.01 | 4.83 | 16.49 | 120.85 |
| | LDM [42] | × | 79.86 | 92.33 | **48.02** | **37.86** | **61.40** |
| | Story-LDM (Ours) | × | **83.29** | **94.61** | 35.54 | 27.32 | 64.89 |
| MUGEN | LDM [42] | × | 95.25 | 97.04 | 21.10 | 23.98 | 123.69 |
| | Story-LDM (Ours) | × | **97.60** | **98.44** | **74.72** | **80.51** | **79.41** |

Table 5. **Experimental results using non-referential text on the Flintstones and Mugen dataset.**

Figure 10. **Positive (1st row) and negative (2nd row) examples.**

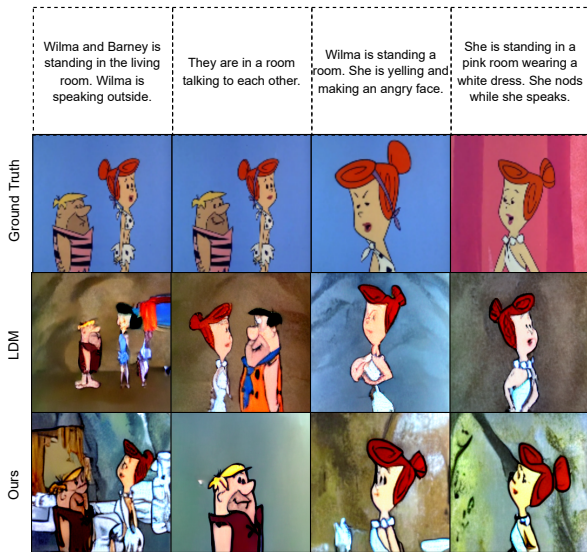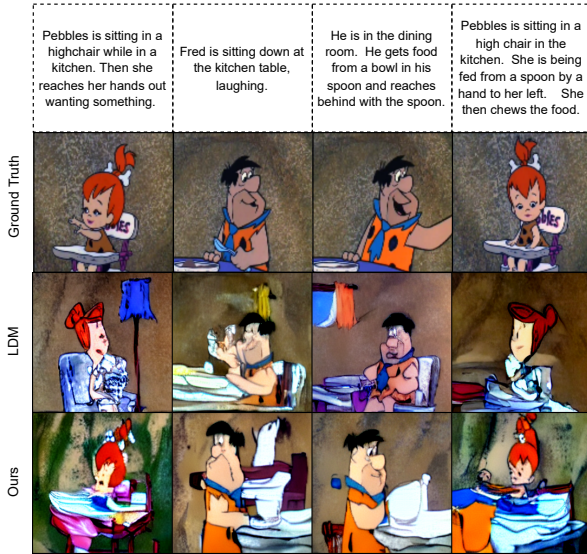Figure 11. **Qualitative results on the MUGEN dataset [16].**

Figure 12. **Qualitative results on the FlintstonesSV dataset [14].**

Figure 13. **Qualitative results on the PororoSV dataset [27].**