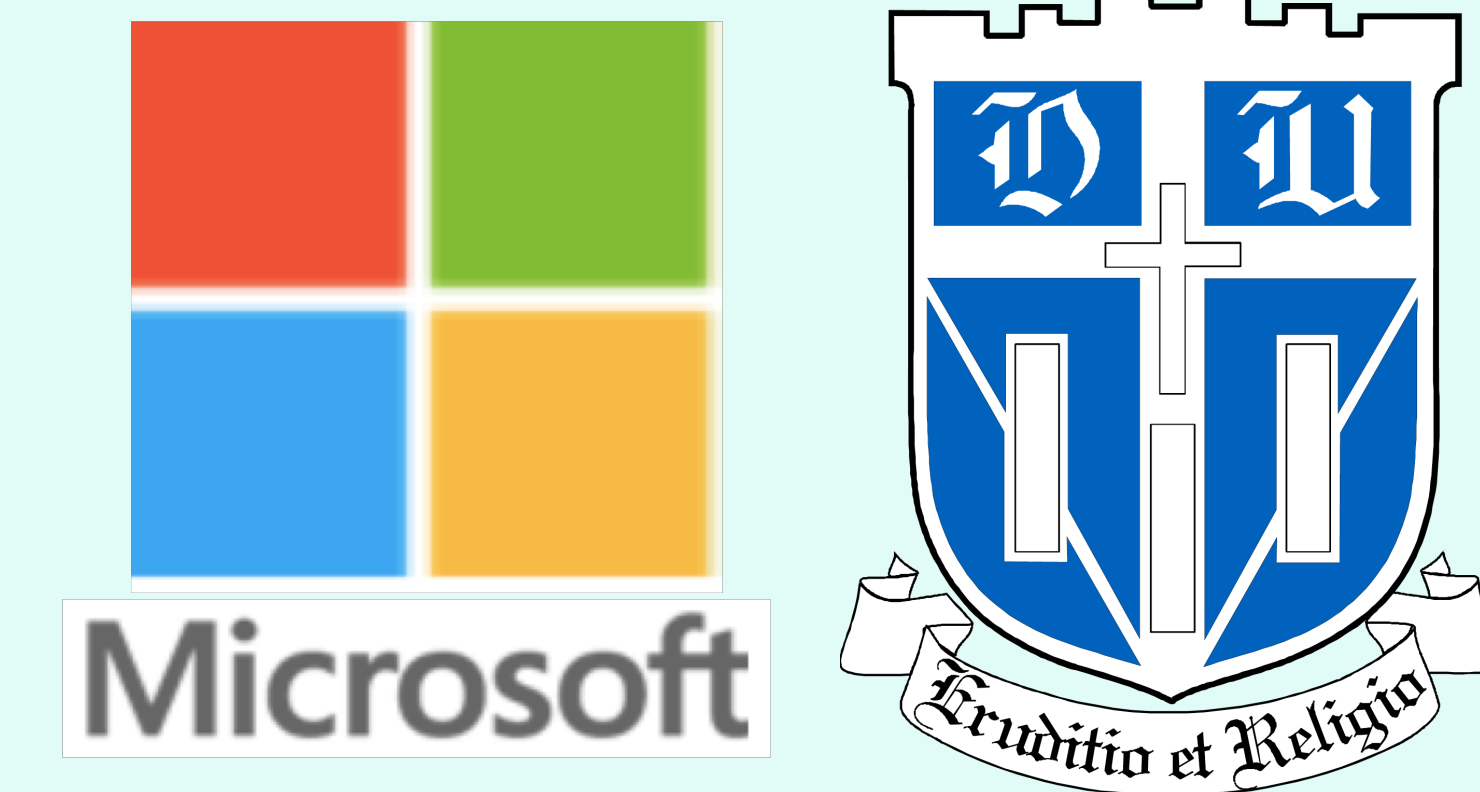




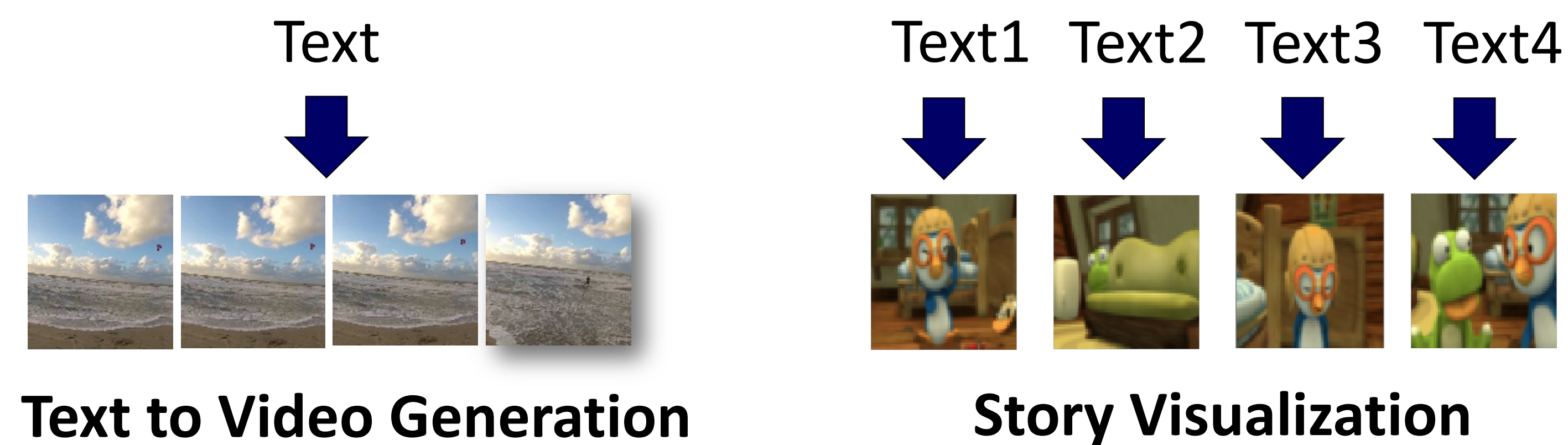
# StoryGAN: A Sequential Conditional GAN for Story Visualization

Yitong Li<sup>1</sup>, Zhe Gan<sup>2</sup>, Yelong Shen<sup>3</sup>, Jingjing Liu<sup>2</sup>, Yu Cheng<sup>2</sup>, Yuexin Wu<sup>4</sup>,  
Lawrence Carin<sup>1</sup>, David Carlson<sup>1</sup>, Jianfeng Gao<sup>2</sup>

<sup>1</sup>Duke University, <sup>2</sup>Microsoft, <sup>3</sup>Tencent, <sup>4</sup>Carnegie Mellon University



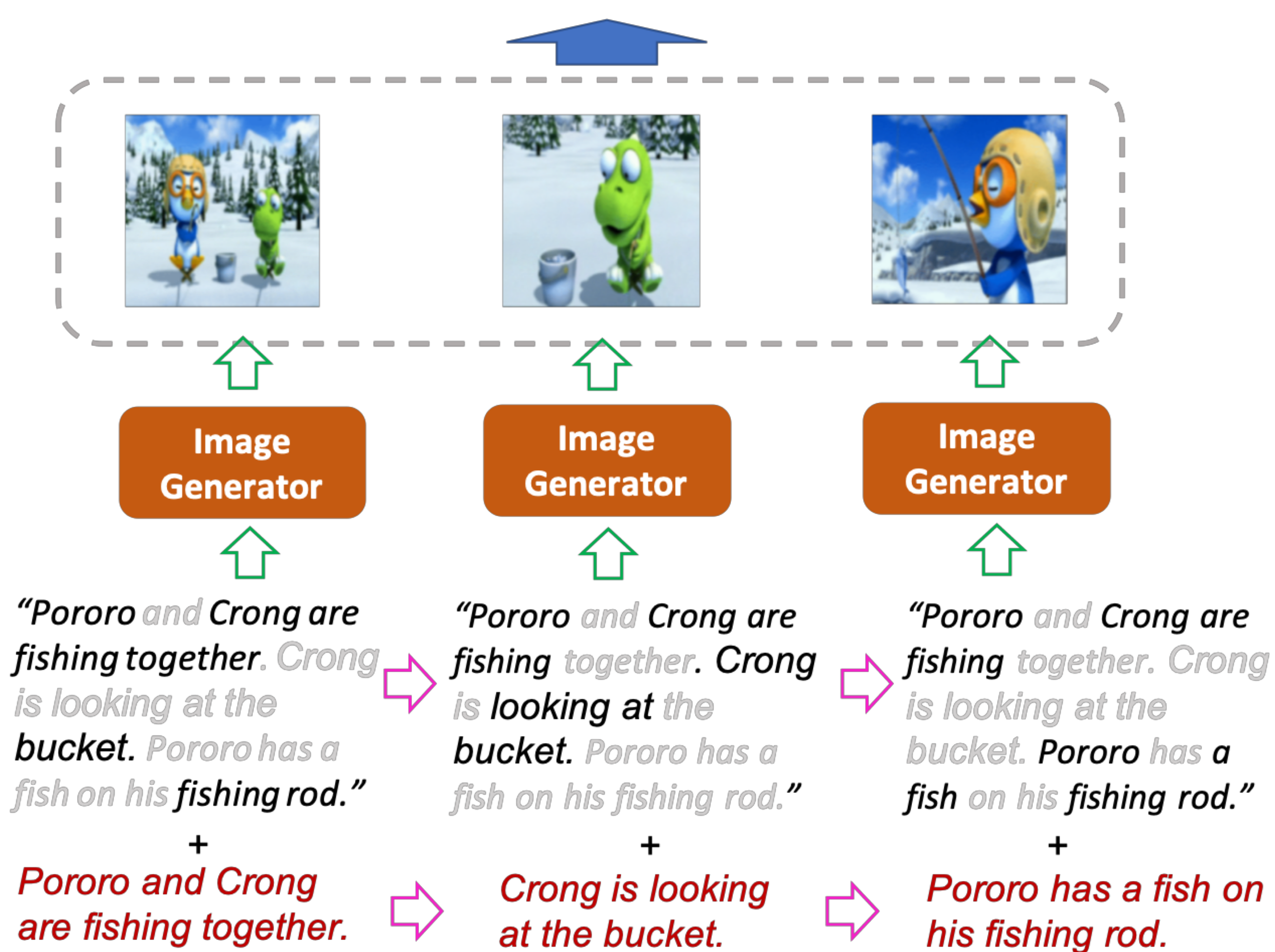
## Story Visualization



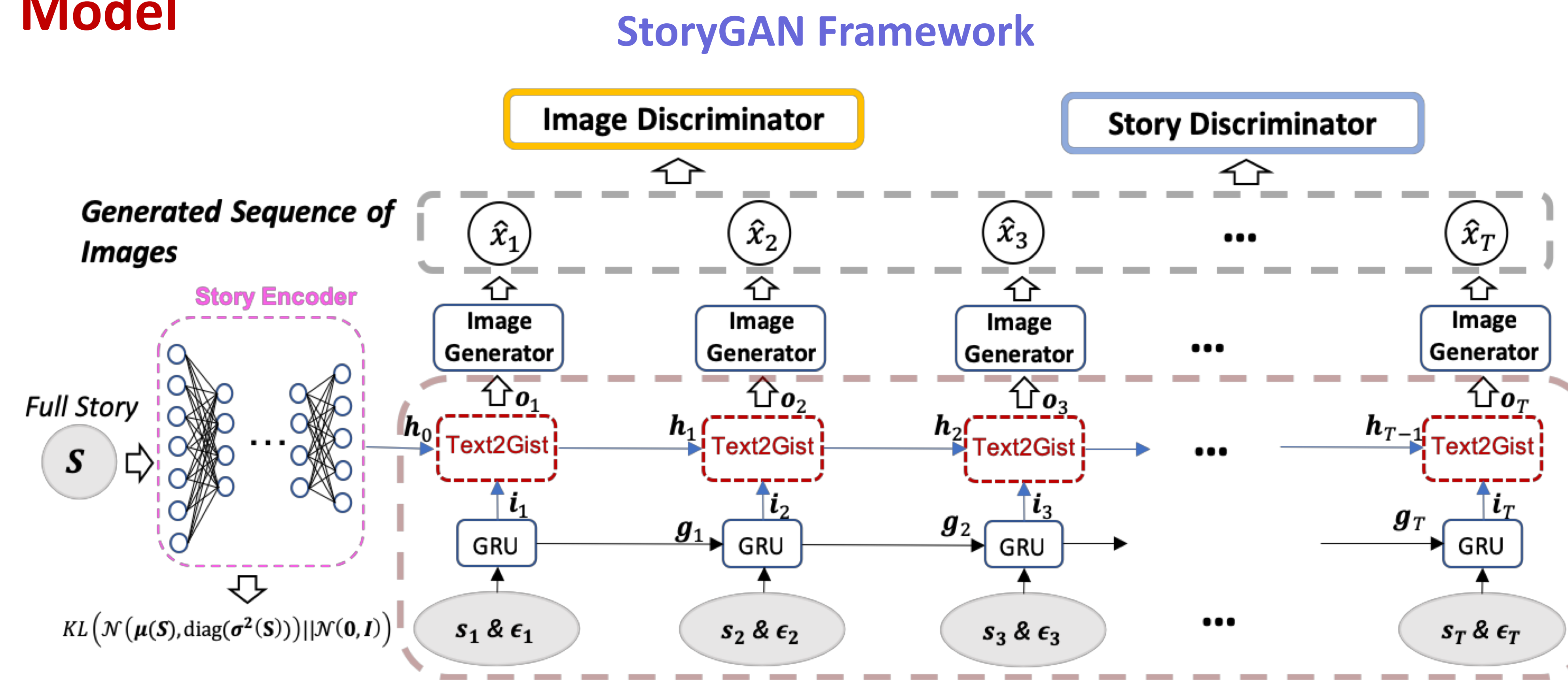
## Text to Video Generation

## Story Visualization

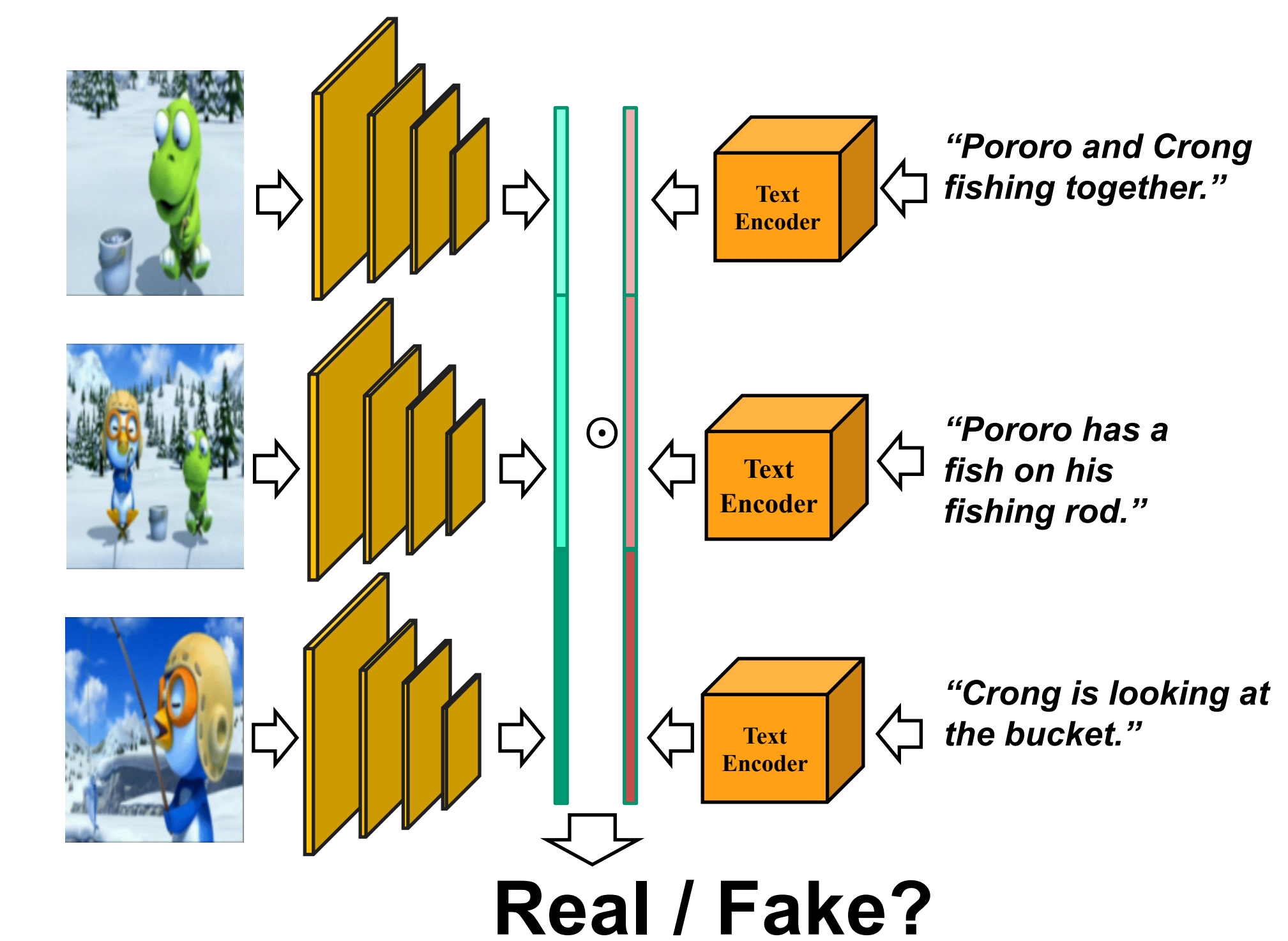
Real/Fake? Coherent to text? Consistent?



## Model



## Story Discriminator



- The Story Encoder learns a stochastic mapping from story  $S$  to a low-dimensional embedding vector  $h_0$ , where  $S = [s_1, \dots, s_T]$
- At each stage, a sentence  $s_t$  and a noise term  $\epsilon_t$  are input
- Text2Gist is built on a GRU cell, which combines the current sentence  $s_t$  with the encoded story  $S$  and the encoded hidden state  $h_{t-1}$  to maintain sequence consistency. The input  $i_t$  is transformed to a filter, then convolved with the hidden state  $h_t$  as  $o_t = Filter(i_t) * h_t$

- The Image Discriminator ensures individual image quality. Note that full story information is incorporated to encourage global consistency
- The Story Discriminator helps enforce the global consistency of the generated image sequence given story  $S$ . It can be written as  $D = \sigma(w^T Encoder(S) \odot Encoder(X) + bias)$ , where  $X = [x_1, \dots, x_T]$  (the image sequence)
- Final loss is  $L_{image} + L_{story}$  from the two-level discriminators

## Experiments

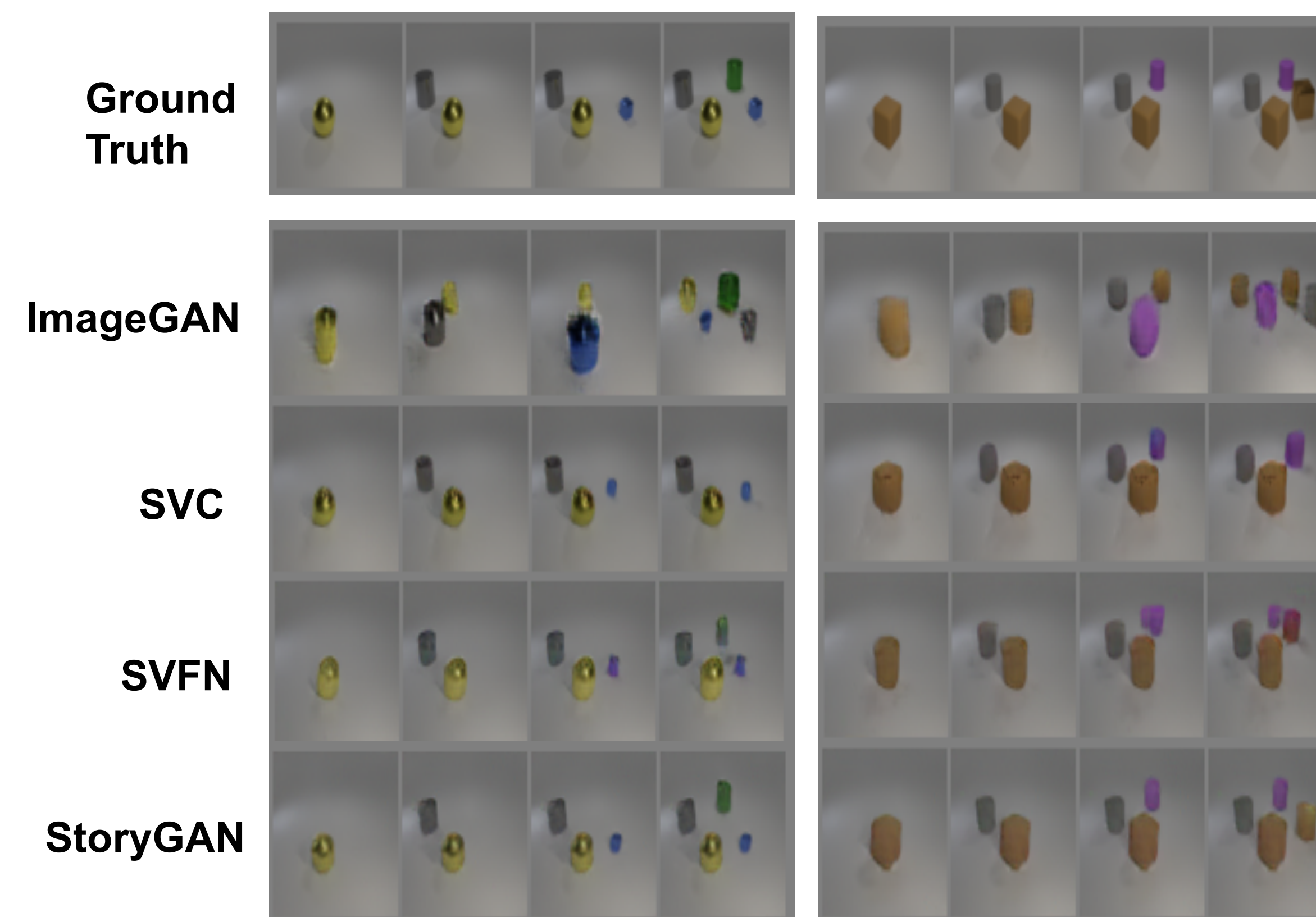
- CLEVR-SV contains 13,000 samples. Each sample is a sequence of four images
- Pororo-SV contains 13,556 samples. Each sample is a sequence of five images

Loopy laughs but tends to be angry. Pororo is singing and dancing and loopy is angry. Loopy says stop to Pororo. Pororo stops. Loopy asks reason to Pororo. Pororo is startled. Pororo is making an excuse to loopy.

Eddy is shocked at what happened now. Pororo tells Eddy that Crong was cloned. Pororo tells Eddy that Crong got into the machine. Eddy says it is not a problem. Eddy tells them that Eddy made a machine to reverse the cloning.

## Motivations and Contributions

- Challenge:** The generated image sequence must consistently and coherently depict the whole story and maintain the logic of the storyline
- New task (Story Visualization):** Visualize a textual story (multi-sentence paragraph) by generating a sequence of images
- New model (StoryGAN):** Consist of a deep Context Encoder that dynamically tracks the story flow and two discriminators: one to enhance the image quality (Image Discriminator) and the other (Story Discriminator) to enforce consistency of the generated sequence
- New datasets:** CLEVR-SV and Pororo-SV. Both have text sequences as input and image sequence as output
- Potential application:** interactive image editing
- Code:** <https://github.com/yitong91/StoryGAN>



CLEVR-SV Dataset Results



Pororo-SV Dataset Results