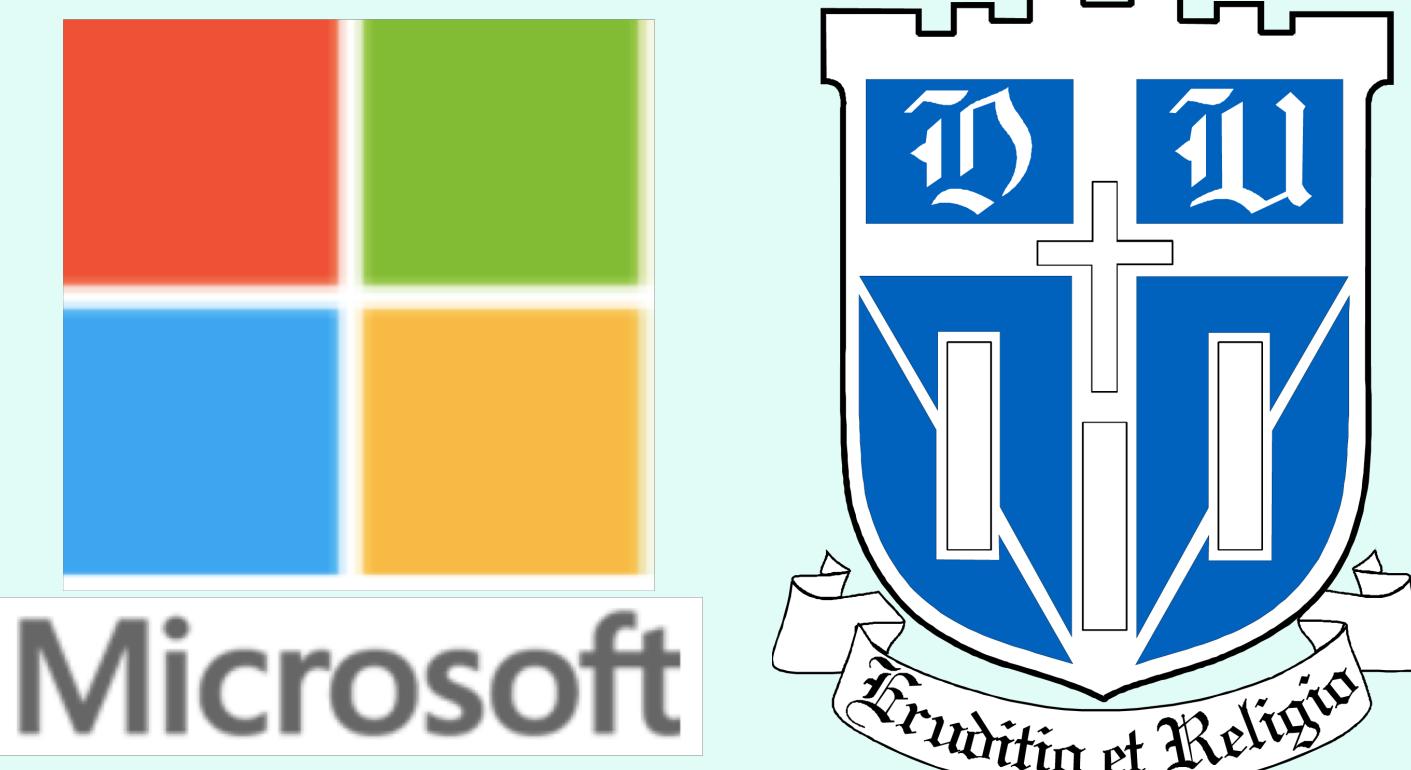




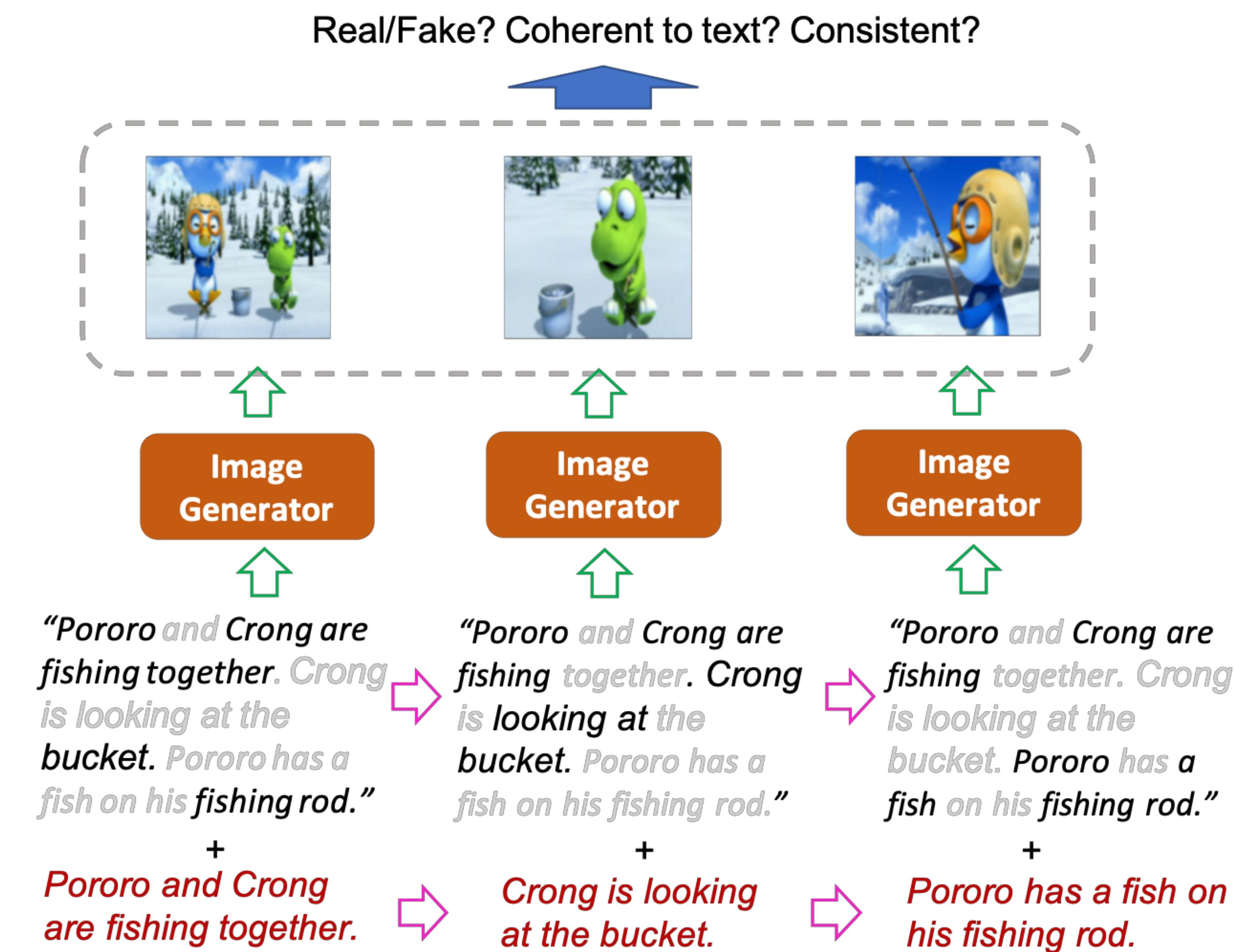
StoryGAN: A Sequential Conditional GAN for Story Visualization

Yitong Li¹, Zhe Gan², Yelong Shen³, Jingjing Liu², Yu Cheng², Yuexin Wu⁴,
Lawrence Carin¹, David Carlson¹, Jianfeng Gao²

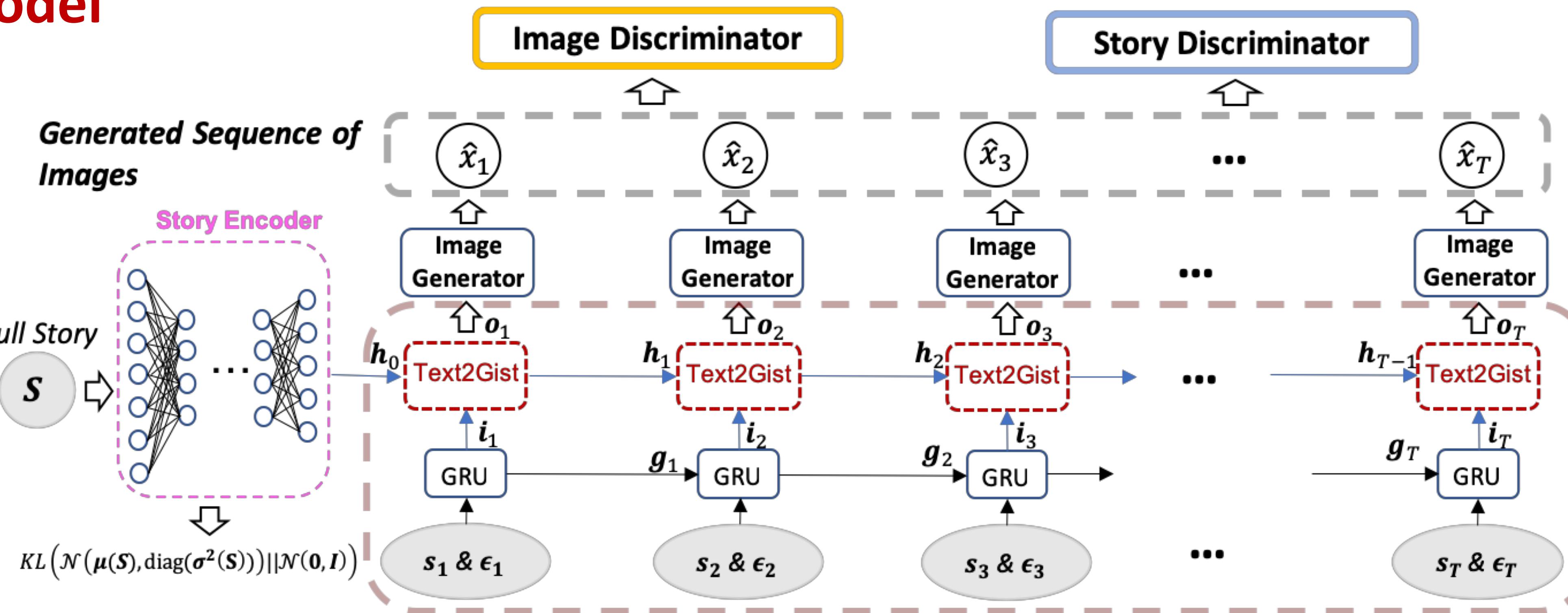
¹Duke University, ²Microsoft, ³Tencent, ⁴Carnegie Mellon University



Story Visualization



Model

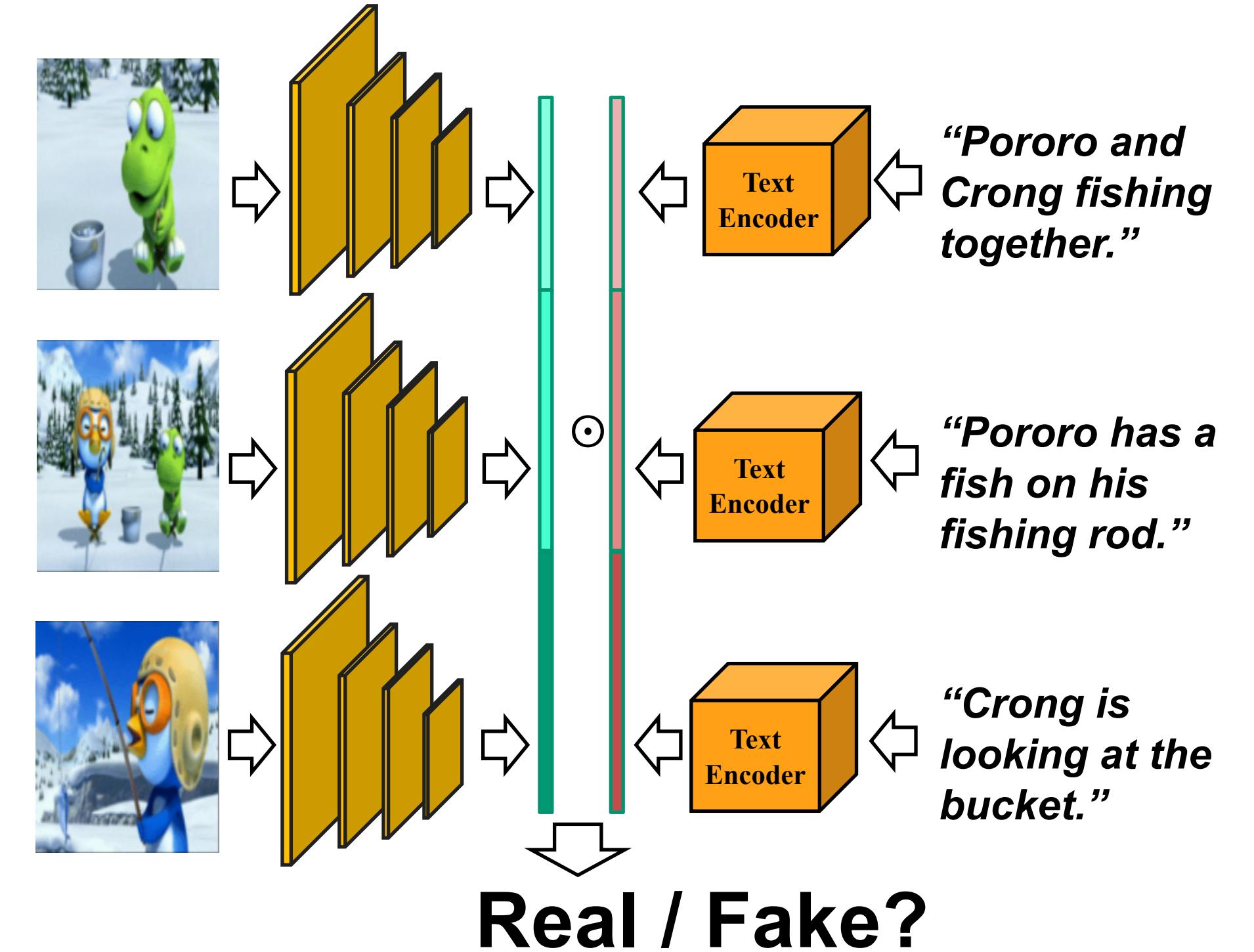


Model Framework

- 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 ϵ_t is 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 give sequence consistency. The input i_t is transformed to a filter, then convolve with the hidden state h_t as

$$o_t = \text{Filter}(i_t) * h_t$$

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



Story Discriminator

Motivation and Contribution

1) New task (Story Visualization): Describe a story (multi-sentence paragraph) by generating a sequence of images.

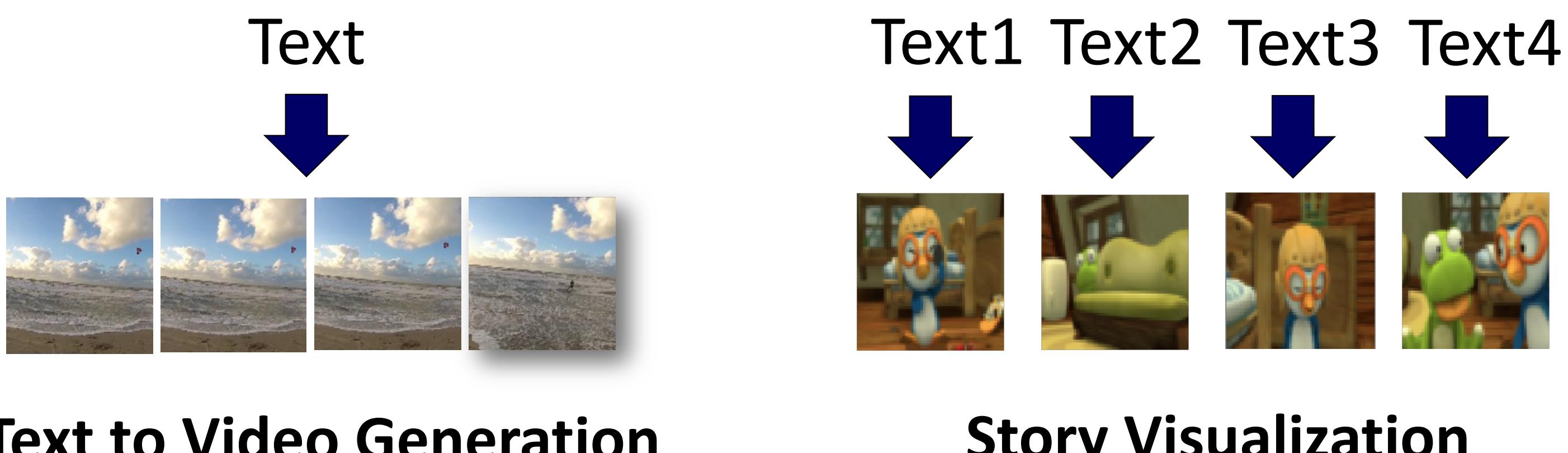
2) Challenge: The generated sequence of images must consistently and coherently depict the whole story and display the logic of the storyline.

3) New model (StoryGAN): Consists of a deep Context Encoder that dynamically tracks the story flow and two discriminators: one to enhance the image quality and the other to enforce consistency of the generated sequences.

4) New datasets: CLEVR-SV and Pororo-SV. Both of them are modified to have text sequence as input and image sequence as output.

5) Potential applications to interactive image editing

6) Code: <https://github.com/yitong91/StoryGAN>



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.

