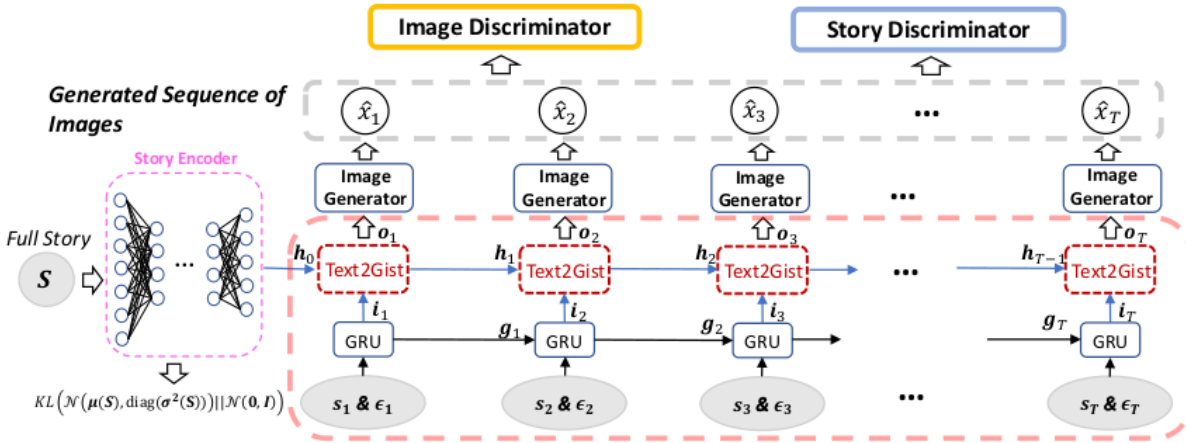


# StoryGAN

GitHub Repo - <https://github.com/yitong91/StoryGAN>



StoryGAN Model Outline

## Components-

- i) **Story Encoder** - encodes story  $S$  into a low-dimensional vector  $h_0$
- ii) **Context Encoder** - encodes input sentence  $s_t$  and its contextual information  $o_t$  (Gist) for each time point  $t$
- iii) **Image Generator** - generates image  $x_t$  based on  $o_t$  for each time step  $t$
- iv) **Image Discriminator & Story Discriminator** - guide the image generation process so as to ensure the generated image sequence  $X$  is locally and globally consistent, respectively.

## Text2Gist-

The lower layer is implemented using standard GRU cells and the upper layer using the proposed Text2Gist cells, which are a variant to GRU cells and are detailed below. At

time step  $t$ , the GRU layer takes as input the concatenation of the sentence  $s_t$  and isometric Gaussian noise  $\epsilon_t$ , and outputs the vector  $i_t$ . The Text2Gist cell combines the GRU's output  $i_t$  with the story context  $h_t$  (initialized by Story Encoder) to generate  $o_t$  that encodes all necessary information for generating an image at time  $t$ .  $h_t$  is updated by the Text2Gist cell to reflect the change of potential context information.

## How Does Context Encoder Work?

Let  $g_t$  and  $h_t$  denote the hidden vectors of the GRU and Text2Gist cells, respectively. The Context Encoder works in two steps to generate its output:

$$i_t, g_t = \text{GRU}(s_t, \epsilon_t, g_{t-1}), \quad (2)$$

$$o_t, h_t = \text{Text2Gist}(i_t, h_{t-1}). \quad (3)$$

## How is input fed in and modified?

We call  $\mathbf{o}_t$  the “*Gist*” vector since it combines all the global and local context information, from  $\mathbf{h}_{t-1}$  and  $\mathbf{i}_t$  respectively, at time step  $t$  (i.e., it captures the “gist” of the information). The Story Encoder initializes  $\mathbf{h}_0$ , while  $\mathbf{g}_0$  is randomly sampled from an isometric Gaussian distribution.

We next give the underlying updates of Text2Gist. Given  $\mathbf{h}_{t-1}$  and  $\mathbf{i}_t$  at time step  $t$ , Text2Gist generates a hidden vector  $\mathbf{h}_t$  and an output vector  $\mathbf{o}_t$  as follows:

$$\mathbf{z}_t = \sigma_z (\mathbf{W}_z \mathbf{i}_t + \mathbf{U}_t \mathbf{h}_{t-1} + \mathbf{b}_z), \quad (4)$$

$$\mathbf{r}_t = \sigma_r (\mathbf{W}_r \mathbf{i}_t + \mathbf{U}_r \mathbf{h}_{t-1} + \mathbf{b}_r), \quad (5)$$

$$\begin{aligned} \mathbf{h}_t = & (\mathbf{1} - \mathbf{z}_t) \odot \mathbf{h}_{t-1} \\ & + \mathbf{z}_t \odot \sigma_h (\mathbf{W}_h \mathbf{i}_t + \mathbf{U}_h (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_h), \end{aligned} \quad (6)$$

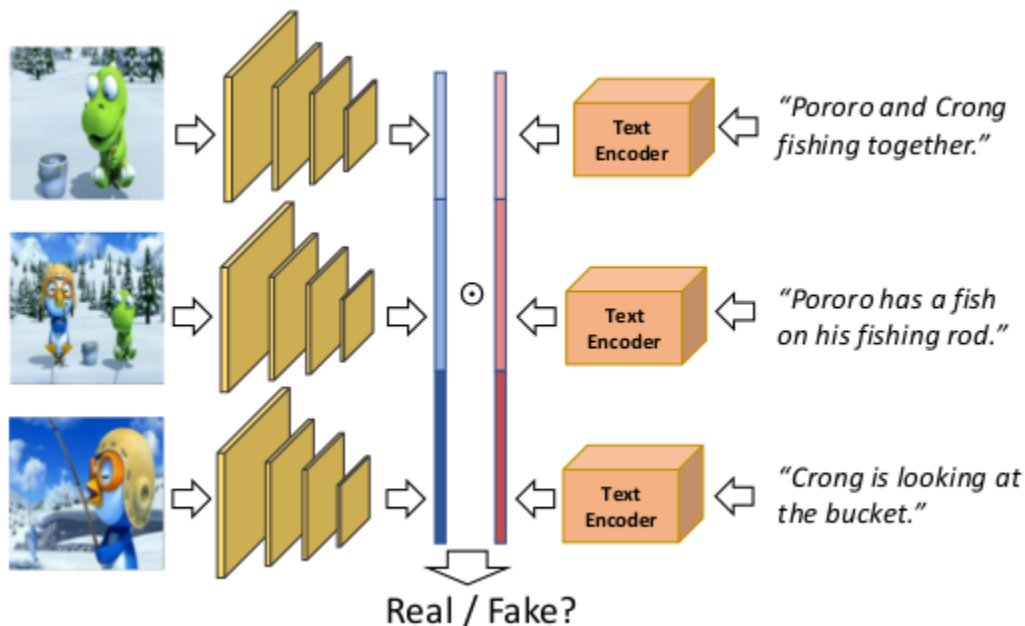
$$\mathbf{o}_t = \text{Filter}(\mathbf{i}_t) * \mathbf{h}_t, \quad (7)$$

where  $\mathbf{z}_t$  and  $\mathbf{r}_t$  are the outputs from the update and reset gates, respectively. The update gate decides how much information from the previous step should be kept, and the reset gate determines what to forget from  $\mathbf{h}_{t-1}$ .  $\sigma_z$ ,  $\sigma_r$  and  $\sigma_h$  are sigmoid non-linearity functions. In contrast to standard GRU cells, output  $\mathbf{o}_t$  is the convolution between  $\text{Filter}(\mathbf{i}_t)$  and  $\mathbf{h}_t$ . The filter  $\mathbf{i}_t$  is learned to adapt to  $\mathbf{h}_t$ . Specifically,  $\text{Filter}(\cdot)$  transforms vector  $\mathbf{i}_t$  to a multi-channel filter of size

$C_{out} \times 1 \times 1 \times \text{len}(\mathbf{h}_t)$  using a neural network, where  $C_{out}$  is the number of output channels. Since  $\mathbf{h}_t$  is a vector, this filter is used as a 1D filter as in a standard convolutional layer.

The convolution operator in Eq. (7) infuses the global contextual information from  $\mathbf{h}_t$  and local information from  $\mathbf{i}_t$ .  $\mathbf{o}_t$  is the output of the Text2Gist cell at time step  $t$ . Since  $\mathbf{i}_t$  encodes information from  $\mathbf{s}_t$  and  $\mathbf{h}_t$  from  $\mathbf{S}$ , which reflects the whole picture of the story, the convolutional operation in Eq. (7) can be seen as helping  $\mathbf{s}_t$  to pick out the important part from the story in the process of generation. Empirically, we find that Text2Gist is more effective than traditional RNNs for story visualization.

## How does Image and Story Discriminator Work?



Compares generated triplets  $\{s_t, h_0, x_t^*\}$  (generated) and  $\{s_t, h_0, x_t\}$  (ground truth)

The story discriminator helps enforce the global consistency of the generated image sequence given story  $S$ . It differs from the discriminators used for video generation, which often use 3D convolution to smooth the changes b/w the frames.

the changes between frames. The overall architecture of the story discriminator is illustrated in Figure 3. The left part is an image encoder, which encodes an image sequence into a sequence of feature vectors  $E_{img}(X) = [E_{img}(x_1), \dots, E_{img}(x_T)]$ , where  $X$  are either real or generated images (which are denoted by  $\hat{X}$ ). These vectors are concatenated into a single vector, shown as the blue rectangular in Fig. 3. Similarly, the right part is a text encoder, which encodes the multi-sentence story  $S$  into a sequence of feature vectors  $E_{txt}(S) = [E_{txt}(s_1), \dots, E_{txt}(s_T)]$ . Likewise, these are concatenated into one big vector, shown as the red rectangle in Fig. 3. The image encoder is implemented as a deep convolutional network and the text encoder as a multi-layer perceptron. Both output a same dimensional vector.

The global consistency score is computed as

$$D_S = \sigma(w^T (E_{img}(X) \odot E_{txt}(S)) + b), \quad (8)$$

## Datasets

### i) CLEVR-SV

## ii) Cartoon Dataset (Pororo)

About 40 video clips forms a complete story. Each story has several QA pairs. In total, the Pororo dataset contains 16K clips of one second videos about 13 distinct characters. The

manually written description has an average length of 13.6 words that describes what is happening and which characters are in each video clip. These 16K video clips are sorted

into 408 movie stories.

We modified the Pororo dataset to fit story visualization task by considering the description for each video clip as the story's text input. For each video clip, we randomly pick out one frame (sampling rate is 30Hz) during training as the real image sample.

Five continuous images form a single story. Finally, we end up with 15, 336 description-story pairs, where 13, 000 pairs are used as training, the remaining 2, 336 pairs for testing. We call this dataset Pororo-SV to differ it from the original Pororo QA dataset.