

# StoryGAN:

## A Sequenital Conditional GAN for Story Visualization

Yitong Li, Zhe Gan, Yelong Shen, Jingjing Liu, Yu Cheng, Yuexin Wu, Lawrence Carin, David Carlson and Jianfeng Gao

# Introduction

## Story Visualization task

to **generate a sequence of images** to **describe a story** written in a **multi-sentence paragraph**

## Two Challenge

the sequence of images must **consistently** and **coherently** depict the whole story

- compared with image generation

  - coherent image sequence

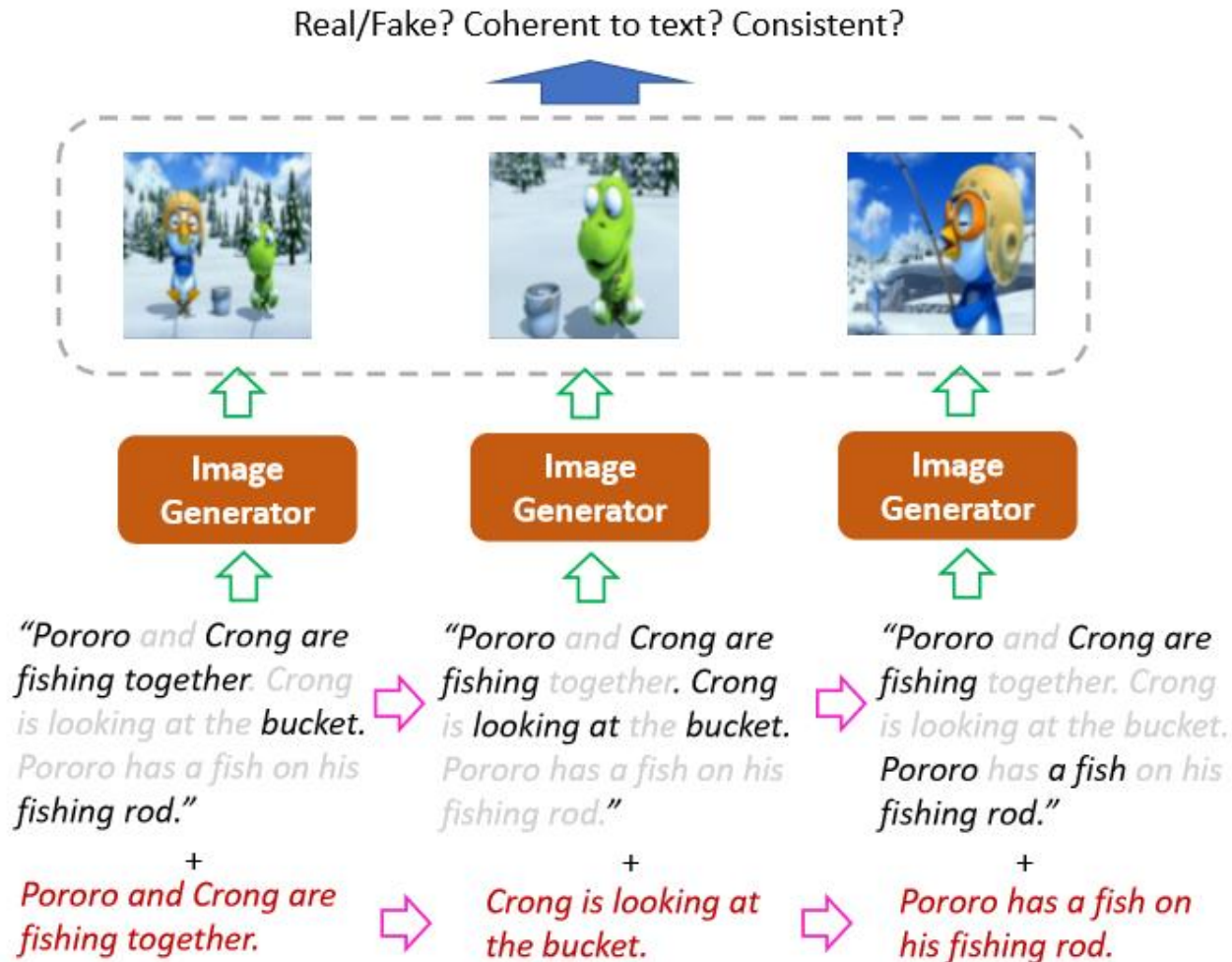
**display the logic** of the storyline

- compared with video generation

  - smooth motion transitions is not important

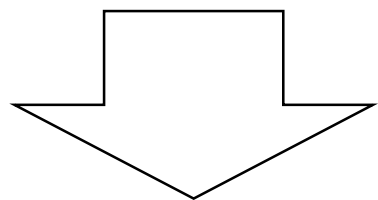
  - sharp change of scenes

# StoryGAN



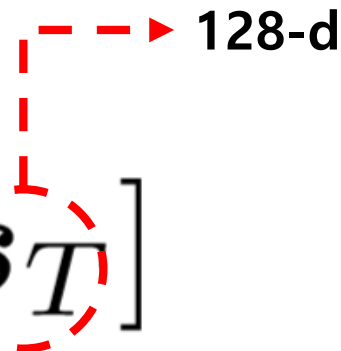
# StoryGAN

$$S = [s_1, s_2, \cdots, s_T]$$

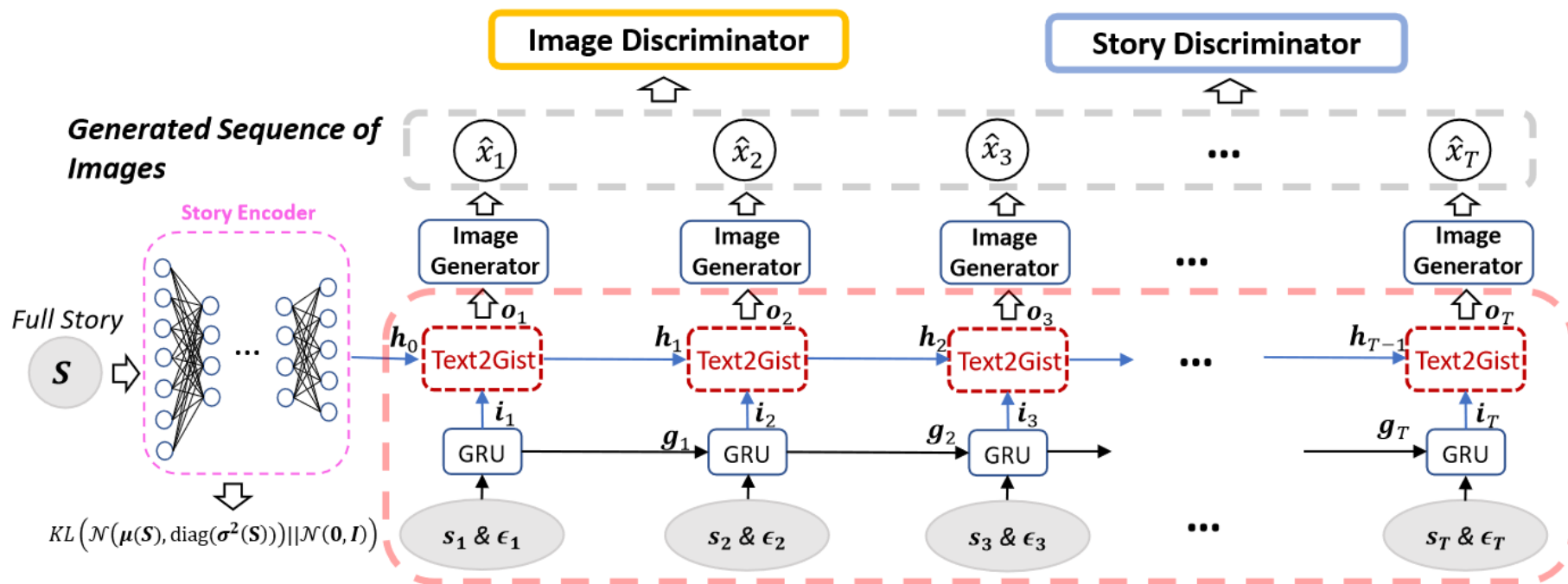


one generated image per sentence  
with **locally** and **globally consistent**  
sentence – image    story – images

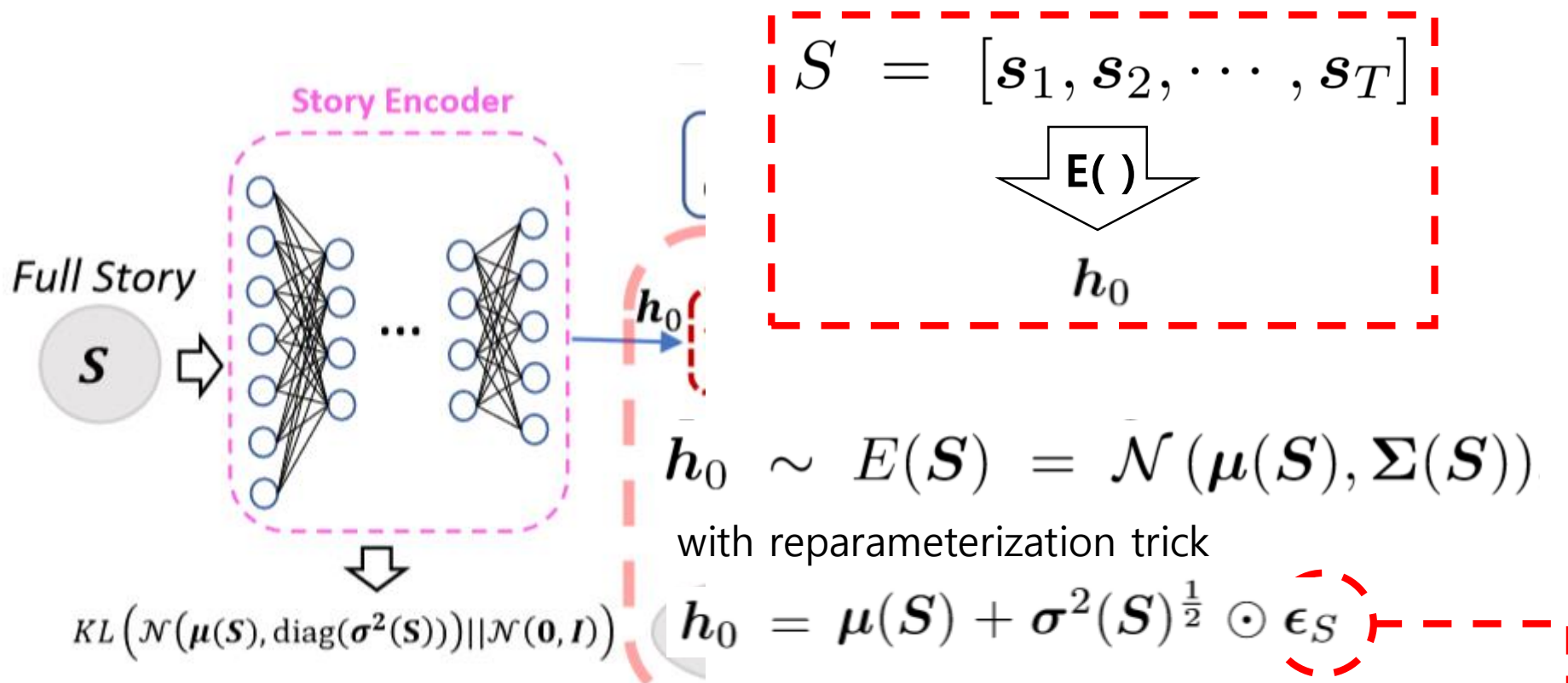
$$\hat{X} = [\hat{x}_1, \hat{x}_2, \cdots, \hat{x}_T]$$



# StoryGAN



# StoryGAN – Story Encoder

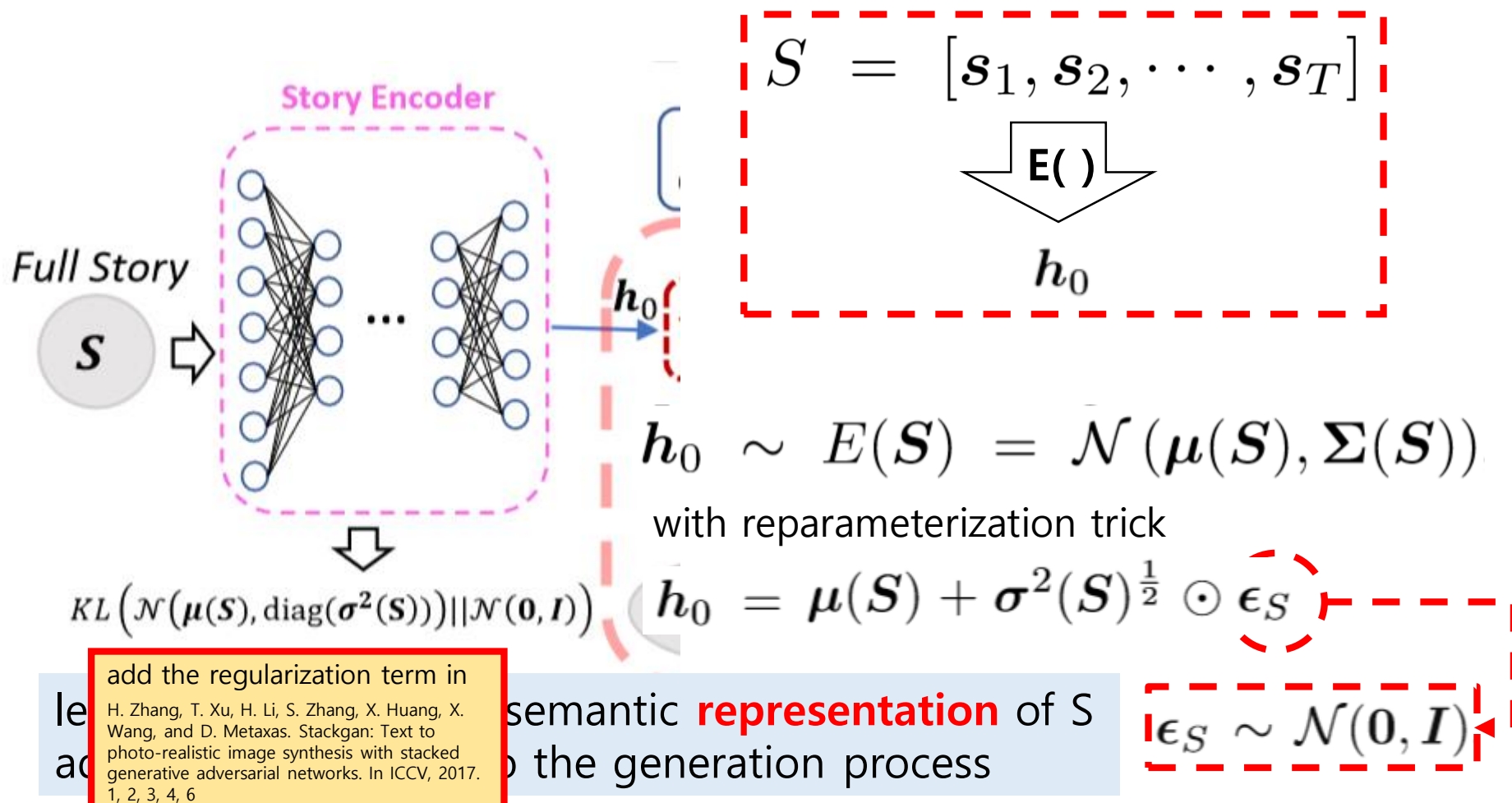


leading to a compact, semantic **representation** of  $S$   
 adding **randomness** to the generation process

$$\epsilon_S \sim \mathcal{N}(\mathbf{0}, I)$$

$$\mathcal{L}_{KL} = KL(\mathcal{N}(\mu(S), \text{diag}(\sigma^2(S))) || \mathcal{N}(\mathbf{0}, I))$$

# StoryGAN – Story Encoder



add the regularization term in

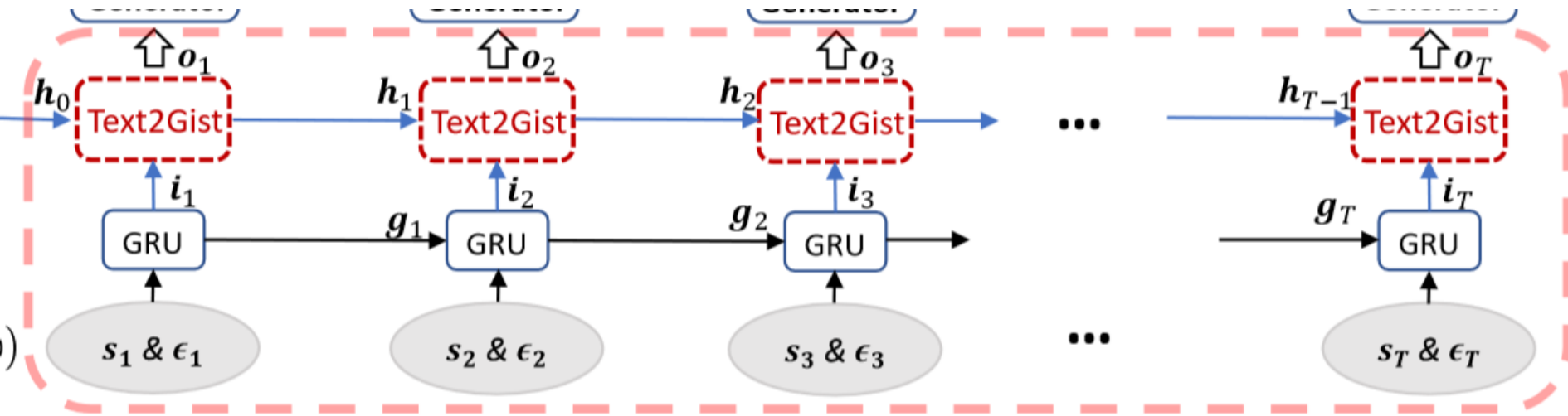
H. Zhang, T. Xu, H. Li, S. Zhang, X. Huang, X. Wang, and D. Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In ICCV, 2017. 1, 2, 3, 4, 6

semantic **representation** of  $S$

to the generation process

$$\mathcal{L}_{KL} = KL(\mathcal{N}(\mu(S), \text{diag}(\sigma^2(S))) || \mathcal{N}(0, I))$$

# StoryGAN – Context Encoder

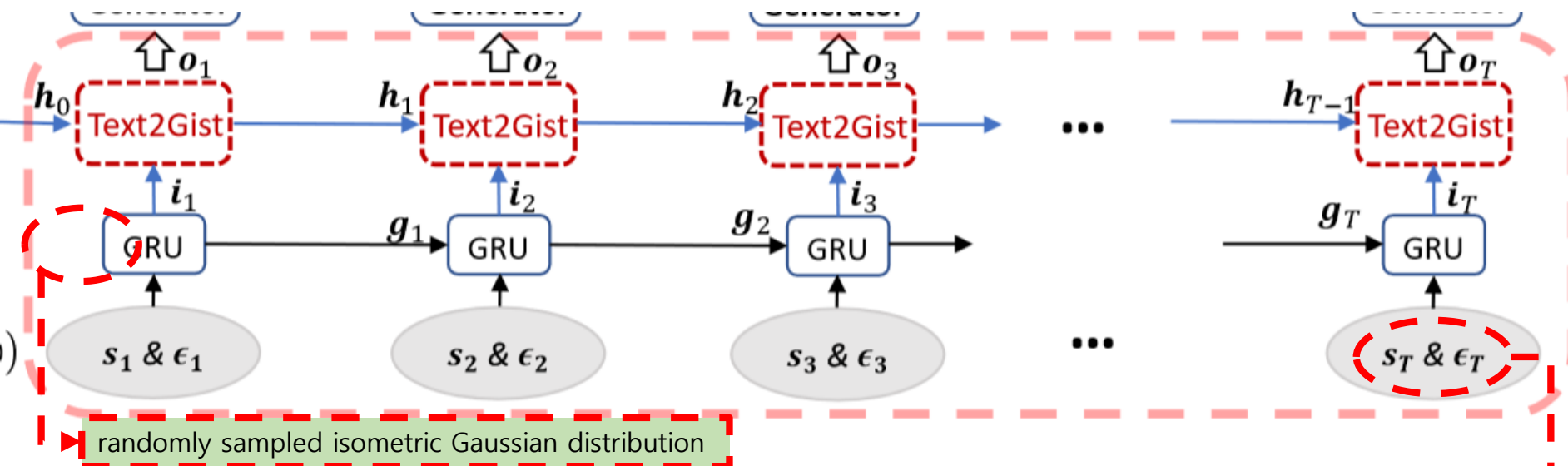


- how to update the **contextual information** to effectively capture background changes
- how to combine **new inputs** and **random noise** when generating each image, to visualize the change of characters, which may shift dramatically

-> **RNN based Context Encoder**



# StoryGAN – Context Encoder



- GRU + Text2Gist

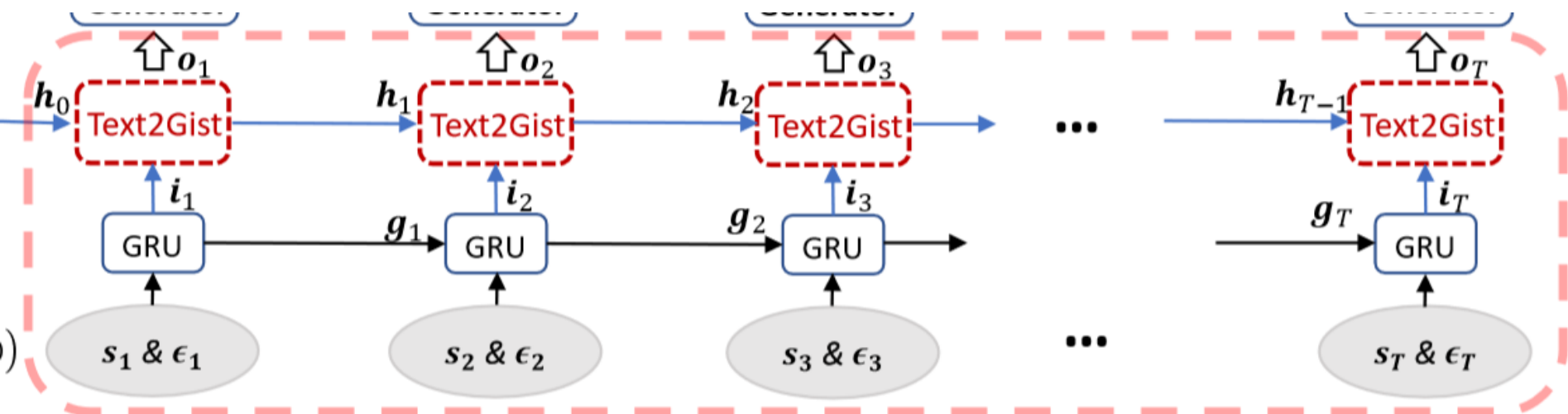
concatenation of sentence and isometric Gaussian noise

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

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

"Gist" vector: combines all the global and local context information

# StoryGAN – Context Encoder



## update gate

decide how much information from the previous step should be kept

## reset gate

determine what to forget from  $h$

$$z_t = \sigma_z (W_z i_t + U_z h_{t-1} + b_z),$$

$$r_t = \sigma_r (W_r i_t + U_r h_{t-1} + b_r),$$

$$h_t = (1 - z_t) \odot h_{t-1}$$

$$+ z_t \odot \sigma_h (W_h i_t + U_h (r_t \odot h_{t-1}) + b_h),$$

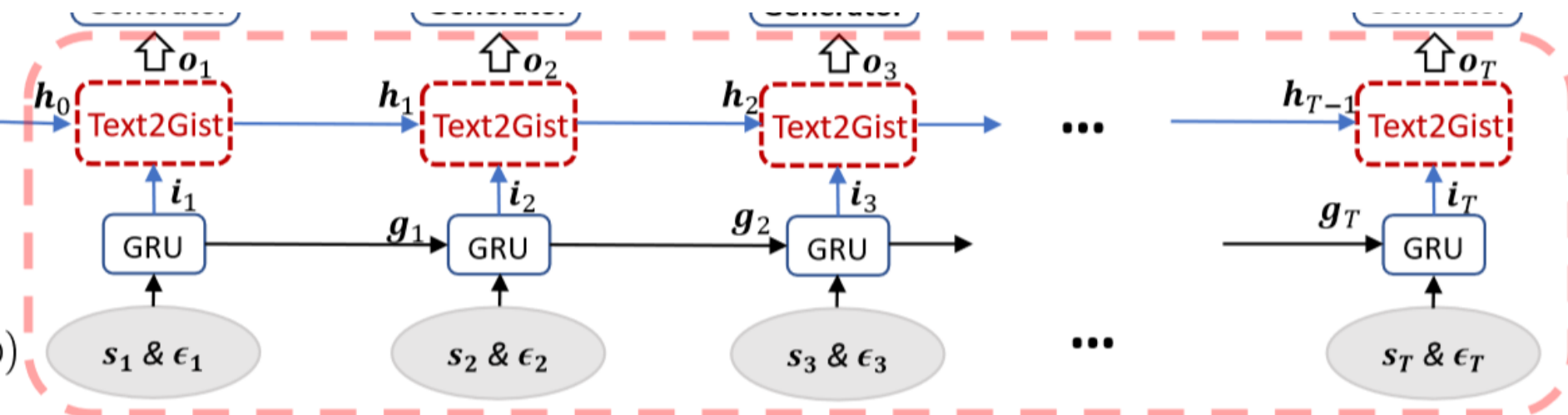
$$o_t = \text{Filter}(i_t) * h_t,$$

## Convolution

Filter(): transforms vector  $i$  to a multi-channel filter 1D filter

$$C_{out} \times 1 \times 1 \times \text{len}(h_t)$$

# StoryGAN – Context Encoder



## update gate

decide how much information from the previous step should be kept

## reset gate

determine what to forget from  $h$

## Convolution

Filter(): transforms vector  $i_t$  to a multi-channel filter  
1D filter

$i_t$ : from sentence  
 $h_t$ : from story

infuse the global contextual information from  $h_t$  and local formation from  $i_t$   
help  $s_t$  to pick out the important part from the story

$$C_{out} \times 1 \times 1 \times \text{len}(h_t)$$

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

# StoryGAN – Discriminators

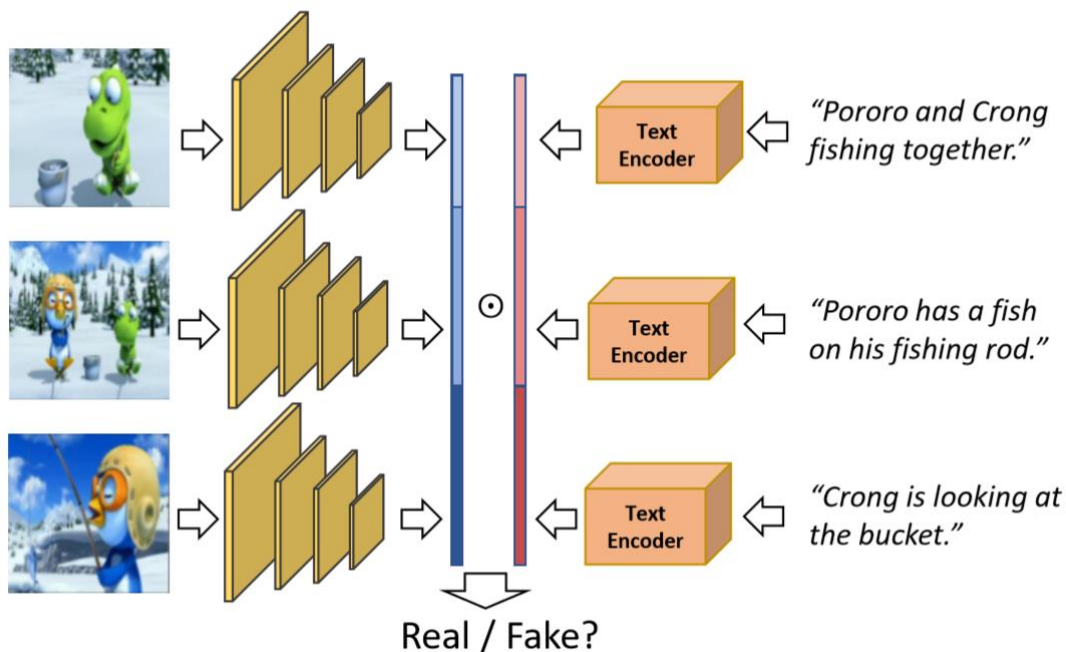
## Image Discriminator

$$\{s_t, h_0, \hat{x}_t\} \longleftrightarrow \{s_t, h_0, x_t\}$$

- the **same sentence** can have a significantly **different generated image** depending on the **context**,
- so it is important to give the **encoded context information** to the **discriminator**

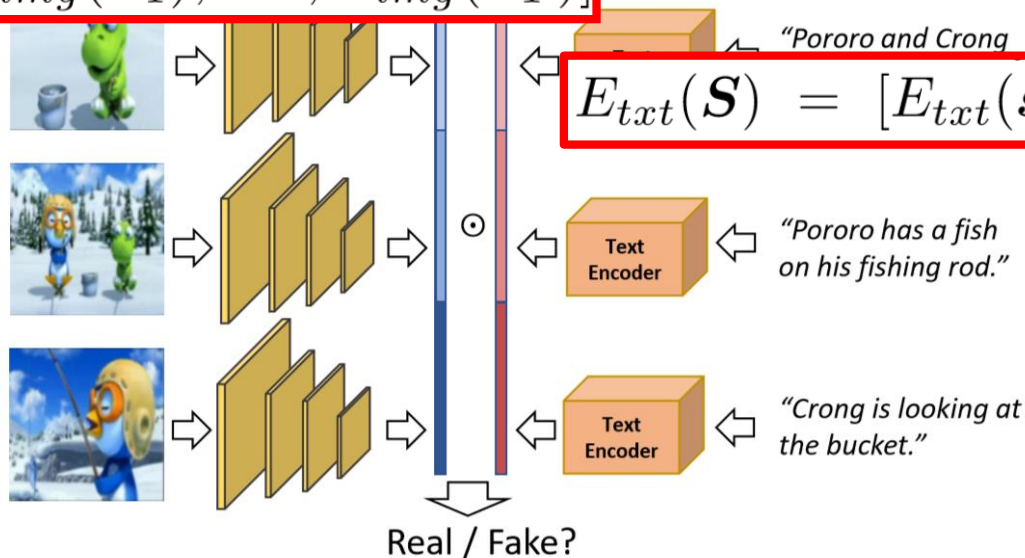
## Story Discriminator

- enforce the global consistency of the generated image sequence given story S



# StoryGAN – Discriminators

$$E_{img}(X) = [E_{img}(x_1), \dots, E_{img}(x_T)]$$



$$E_{txt}(S) = [E_{txt}(s_1), \dots, E_{txt}(s_T)]$$

$$D_S = \sigma(w^\top (E_{img}(X) \odot E_{txt}(S)) + b)$$

By **pairing each sentence and image**, the story discriminator can consider both **local matching** and **global consistency jointly**

# StoryGAN – Objective function

$$\min_{\boldsymbol{\theta}} \max_{\boldsymbol{\psi}_I, \boldsymbol{\psi}_S} \alpha \mathcal{L}_{Image} + \beta \mathcal{L}_{Story} + \mathcal{L}_{KL}$$

$$\begin{aligned} \mathcal{L}_{Image} = & \sum_{t=1}^T (\mathbb{E}_{(\mathbf{x}_t, \mathbf{s}_t)} [\log \textcolor{red}{D}_I(\mathbf{x}_t, \mathbf{s}_t, \mathbf{h}_0; \boldsymbol{\psi}_I)] \\ & + \mathbb{E}_{(\boldsymbol{\epsilon}_t, \mathbf{s}_t)} [\log(1 - D_I(G(\boldsymbol{\epsilon}_t, \mathbf{s}_t; \boldsymbol{\theta}), \mathbf{s}_t, \mathbf{h}_0; \boldsymbol{\psi}_I))]) \end{aligned}$$

$$\begin{aligned} \mathcal{L}_{Story} = & \mathbb{E}_{(\mathbf{X}, \mathbf{S})} [\log \textcolor{red}{D}_S(\mathbf{X}, \mathbf{S}; \boldsymbol{\psi}_S)] \\ & + \mathbb{E}_{(\boldsymbol{\epsilon}, \mathbf{S})} \left[ \log(1 - D_S([G(\boldsymbol{\epsilon}_t, \mathbf{s}_t; \boldsymbol{\theta})]_{t=1}^T), \mathbf{S}; \boldsymbol{\psi}_S)) \right] \end{aligned}$$

# StoryGAN – Algorithm

---

**Algorithm 1** Training Procedure of StoryGAN

---

**Input:** Encoded sentence vectors  $\mathbf{S}_n = [s_{n1}, s_{n2}, \dots, s_{nT}]$  and corresponding images  $\mathbf{X}_n = [x_{n1}, \dots, x_{nT}]$  for  $n = 1, \dots, N$ .

**Output:** Generator parameters  $\theta$  and discriminator parameters  $\psi_I$  and  $\psi_S$ .

---

**for**  $iter = 1$  to  $max\_iter$  **do**

**for**  $iter_I = 1$  to  $k_I$  **do**

        Sample a mini-batch of story-sentence pairs  $\{(s_t, S, x_t)\}$  from the training set.

        Compute  $h_0$  as the initialization of the Text2Gist layer and the KL regularization term as Eq. (1).

        Generate a single output image  $\hat{x}$ .

        Update  $\psi_I$  and  $\theta$ .

**end for**

**for**  $iter_S = 1$  to  $k_S$  **do**

        Sample a mini-batch of story-image pair  $\{(S, X)\}$  from training set.

        Compute  $h_0$  and update  $h_t$  at each time step  $t$

        Generate image sequence  $\hat{X}$ .

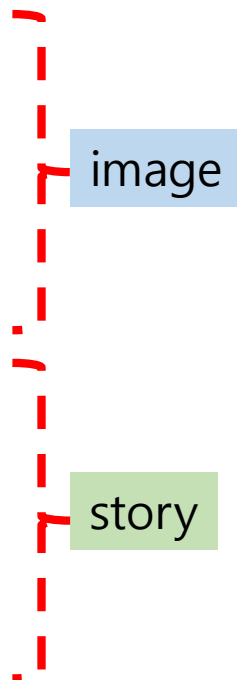
        Update  $\psi_S$  and  $\theta$ .

**end for**

**end for**

---

- using **Adam** optimizer
- **different mini-batch sizes** for image and story discriminators



# StoryGAN – Experiments

Video generation result is too blurry and not comparable to StoryGAN

- our comparisons are mainly to ablated versions of our proposed model
- for a fair comparison, use same image generator, context encoder and discriminators

**ImageGAN:** ImageGAN follows the work in [28, 36] and does not use the story discriminator, story encoder and Context Encoder. Each image is generated independently. However, for a reasonable comparison, we concatenate  $s_t$ , encoded story  $S$  and a noise term as input. Otherwise, the model fails on the task. This is the simplest version of StoryGAN.

**SVFN:** In “Story Visualization by Filter Network” (SVFN), the concatenation in SVC is replaced by a filter network. Sentence  $s_t$  is transformed into a filter and convolved with the encoded story. Specifically, the image generator input is  $o_t = \text{Filter}(i_t) * h_0$  instead of Eq. 7.

**SVC:** In “Story Visualization by Concatenation” (SVC), the Text2Gist cell in StoryGAN is replaced by simple concatenation of the encoded story and description feature vectors [31]. Compared to ImageGAN, SVC includes the additional story discriminator, and is visualized in Figure 4.

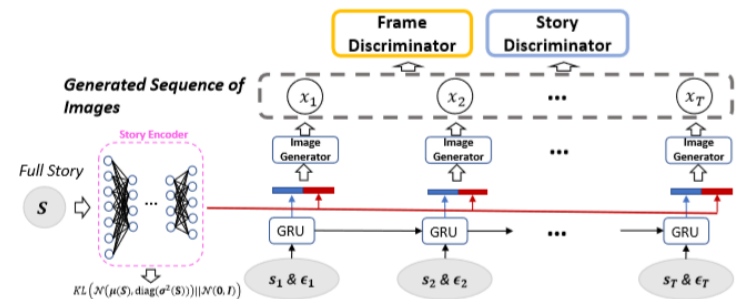


Figure 4: The framework of the baseline model SVC, where the story and individual sentence are concatenated to form the input.



# StoryGAN – Experiments with CLEVR-SV

## CLEVR -> CLEVR-SV

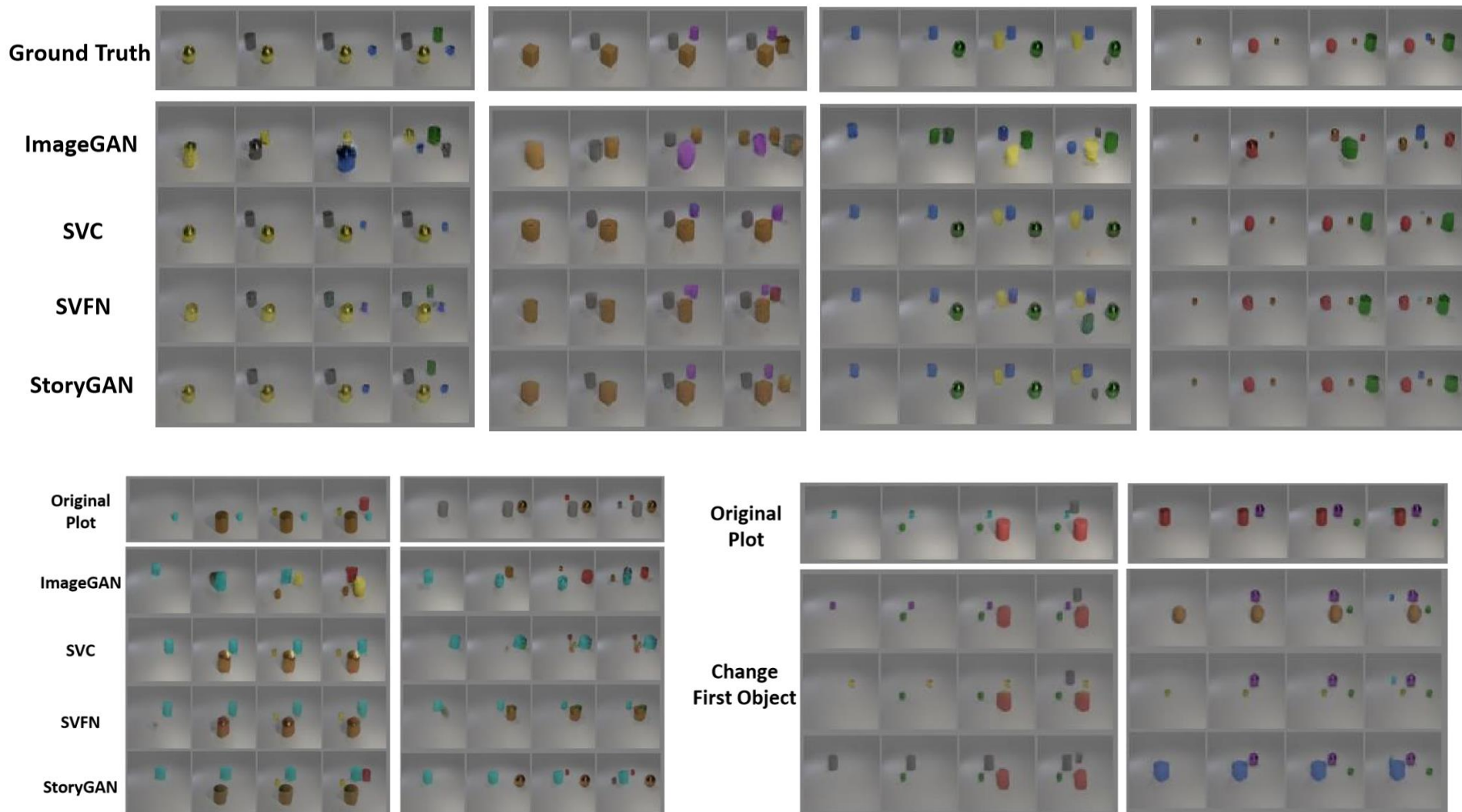
- the **maximum number of objects** in one story is limited to four.
- objects are made of **metallic/rubber** with **eight different colors** and **two different sizes**.
- the **object shape** can be cylinder, cube or sphere.
- the object is **added one at a time**, resulting in a **four-image** sequence per story.

StoryGAN generates more feasible images than others

- **Text2Gist** cell tracks the progress of story
- **story and image discriminators** keep the consistency of objects in the generation process
- using the **Story Encoder** to initialize the Text2Gist cell gives better result on first generated image

	ImageGAN [28]	SVC	SVFN	StoryGAN
SSIM	0.596	0.641	0.654	0.672

# StoryGAN – Experiments with CLEVR-SV



# StoryGAN – Experiments with Pororo-SV

Pororo -> Pororo-SV

- randomly pick out one frame during training as the real image sample
- five continuous images form a single story

StoryGAN's **first image** has a much **higher quality** than other baselines -> **Story Encoder**

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.

Ground Truth



ImageGAN



SVC



SVFN



StoryGAN



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.

Ground Truth



ImageGAN



SVC



SVFN



StoryGAN

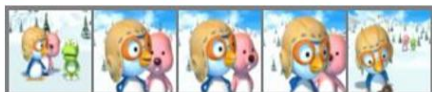


# StoryGAN – Experiments

Input Story: *c1 and c2 are standing in the snow. c1 tells a story to c3. c3 wants to joint c1 and c2. c1 continuous to talk. c1 looks down. They suddenly noticed that there is something lying on the snow.*

C1 = Pororo, C2 = Lopy, C3 = Crong

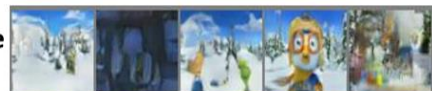
GT



C1 = Pororo, C2 = Eddy, C3 = Rody



Image  
GAN



SVC



SVFN



Story  
GAN



Method	ImageGAN	SVC	SVFN	StoryGAN
Rank	$2.91 \pm 0.05$	$2.42 \pm 0.04$	$2.77 \pm 0.04$	$1.94 \pm 0.05$

	StoryGAN vs ImageGAN		
Choice (%)	StoryGAN	ImageGAN	Tie
Visual Quality	$74.17 \pm 1.38$	$18.60 \pm 1.38$	7.23
Consistence	$79.15 \pm 1.27$	$15.28 \pm 1.27$	5.57
Relevance	$78.08 \pm 1.34$	$17.65 \pm 1.34$	4.27

# Network Structure

Layer	Story Encoder
1	LINEAR-( $128 \times T$ , 128), BN, RELU
Layer	Context Encoder
1	LINEAR-(NOISEDIM + TEXTDIM, 128), BN, RELU
2	GRU-(128, 128)
3	Text2Gist-(128, 128)
Layer	Filter Network
1	LINEAR-(128, 1024), BN, TANH
2	RESHAPE(16, 1, 1, 64)
Layer	Image Generator
1	CONV-(C512, K3, S1, P1), BN, RELU
2	UPSAMPLE-(2,2)
3	CONV-(C256, K3, S1, P1), BN, RELU
4	UPSAMPLE-(2,2)
5	CONV-(C128, K3, S1, P1), BN, RELU
6	UPSAMPLE-(2,2)
7	CONV-(C64, K3, S1, P1), BN, RELU
8	UPSAMPLE-(2,2)
9	CONV-(C3, K3, S1, P1), BN, TANH
Layer	Image Discriminator
1	CONV-(C64, K4, S2, P1), BN, LEAKY RELU
2	CONV-(C128, K4, S2, P1), BN, LEAKY RELU
3	CONV-(C256, K4, S2, P1), BN, LEAKY RELU
4	CONV-(C512, K4, S2, P1), BN, LEAKY RELU
5*	CONV-(C512, K3, S1, P1), BN, LEAKY RELU
6	CONV-(C1, K4, S4, P0), SIGMOID
Layer	Story Discriminator (Image Encoder)
1	CONV-(C64, K4, S2, P1), BN, LEAKY RELU
2	CONV-(C128, K4, S2, P1), BN, LEAKY RELU
3	CONV-(C256, K4, S2, P1), BN, LEAKY RELU
4	CONV-(C512, K4, S2, P1), BN, LEAKY RELU
5	CONV-(C32, K4, S2, P1), BN, CONCAT
6	RESHAPE-(1, $32 \times 4 \times T$ )
Layer	Story Discriminator (Text Encoder)
1	LINEAR-( $128 \times T$ , $32 \times 4 \times T$ ), BN