

MirrorGAN: Learning Text-to-image Generation by Redescription

Sanghyuck Na

Jan, 15, 2021

Dongguk University

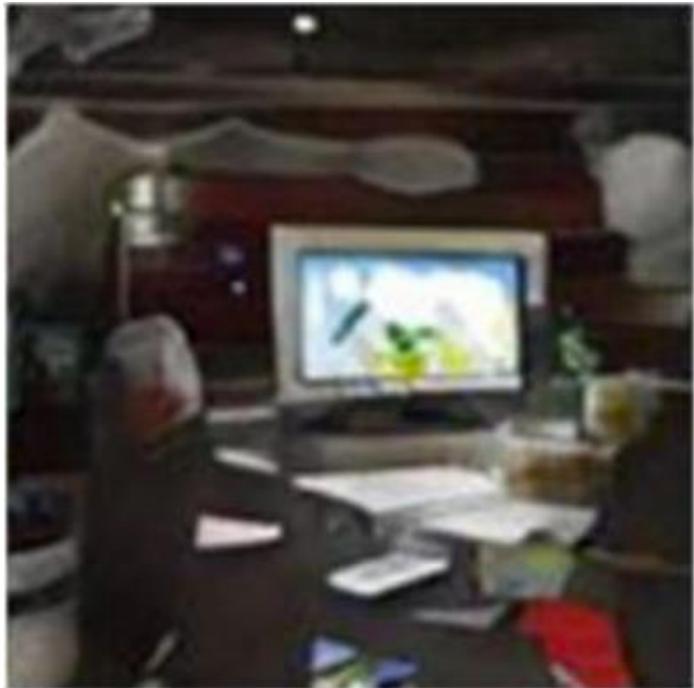
Artificial Intelligence Laboratory

shna@Dongguk.edu

- 1. Introduction**
- 2. StackGAN**
- 3. CycleGAN**
- 4. MirrorGAN**
- 5. Result**
- 6. Reference**

Introduction

Text to image generation



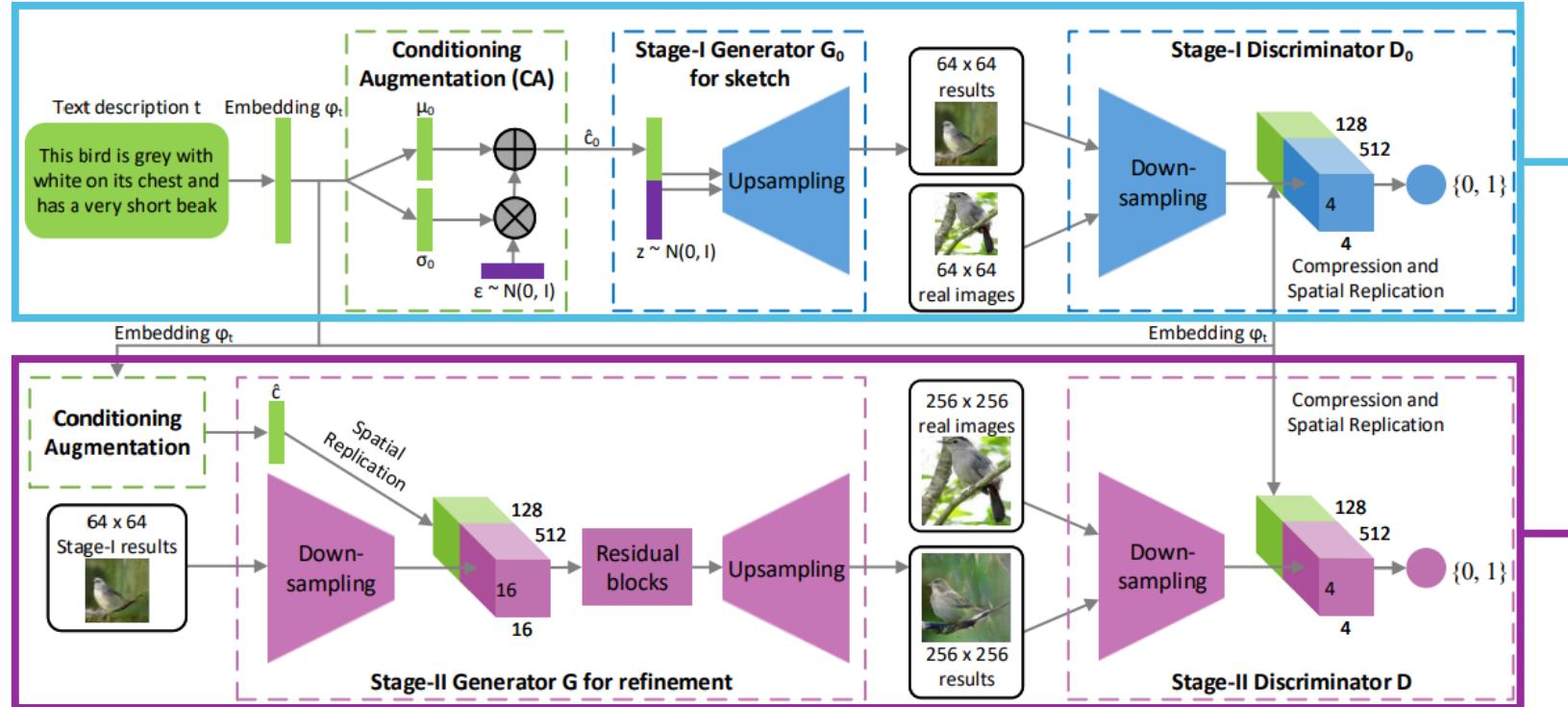
There is a lot of electrical
sitting on the table.

Image to text generation



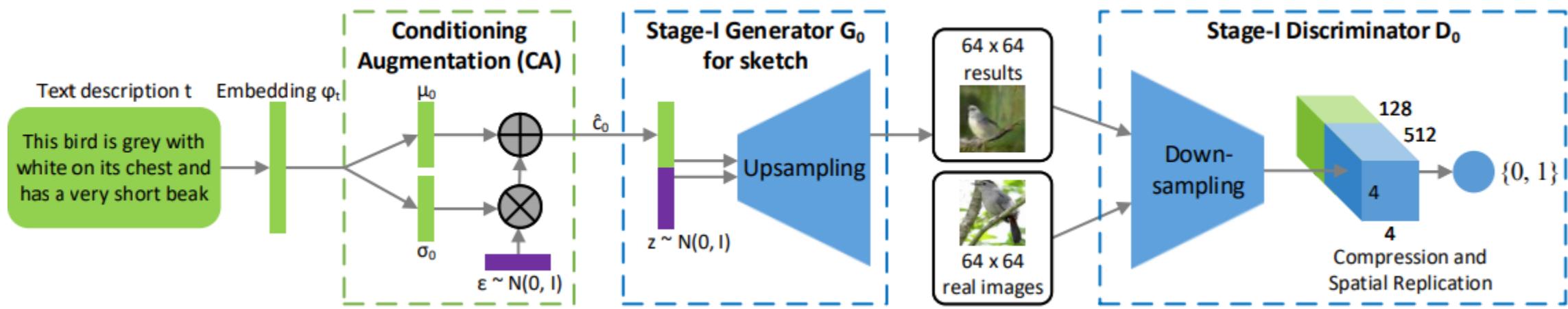
man in graduation robes riding bicycle
cyclist giving thumbs up poses with his bicycle by right
of way sign at park
man riding motorcycle on street

Green: Human ground truth.
Red: Top-scoring sentence from training set.
Blue: Generated sentence.



Stage-I GAN: it sketches the primitive shape and basic colors of the object conditioned on the given text description, and draws the background layout from a random noise vector, yielding a low-resolution image.

Stage-II GAN: it corrects defects in the low-resolution image from Stage-I and completes details of the object by reading the text description again, producing a high-resolution photo-realistic image.



Conditioning Augmentation (CA)

$$D_{KL}(\mathcal{N}(\mu(\varphi_t), \Sigma(\varphi_t)) \| \mathcal{N}(0, I))$$

t : *text description*

z : *noise vector from Gaussian Distribution*

φ_t : *text embedding networks (pre – trained)*

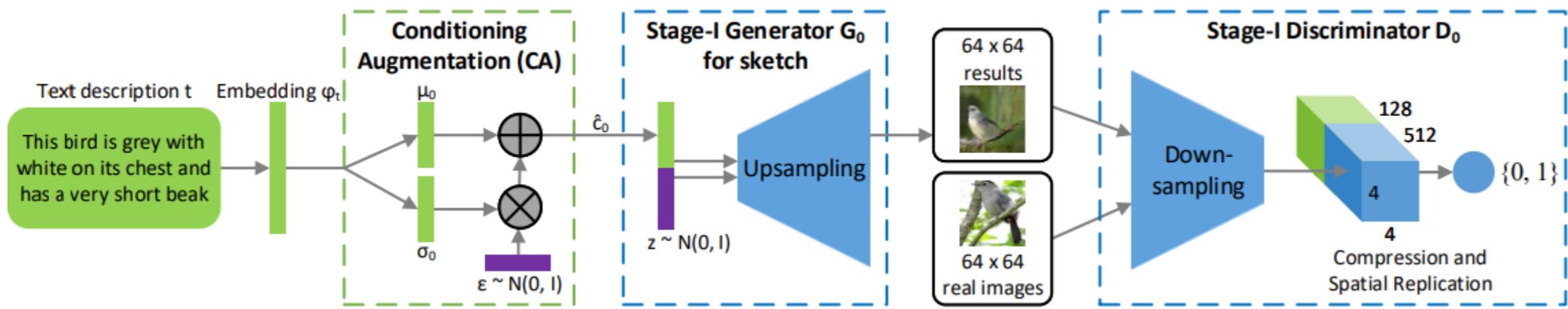
\hat{c}_0 : *conditioning variable*

$\mathcal{N}(\mu(\varphi_t), \Sigma(\varphi_t))$: *conditioning Gaussian distribution*

$\mathcal{N}(0, I)$: *normal distribution*

$\Sigma(\varphi_t)$: *diagonal covariance matrix*

s_0 : *image generated by the Stage-I*



$$\begin{aligned} L_{D_0} &= \mathbb{E}_{(I_0, t) \sim p_{data}} [\log D_0(I_0, \varphi_t)] \\ &+ \mathbb{E}_{z \sim p_z, t \sim p_{data}} [\log(1 - D_0(G_0(z, \hat{c}_0, \varphi_t)))] \end{aligned}$$

$$\begin{aligned} L_{G_0} &= \mathbb{E}_{z \sim p_z, t \sim p_{data}} [\log(1 - D_0(G_0(z, \hat{c}_0, \varphi_t)))] \\ &+ \lambda D_{KL}(\mathcal{N}(\mu_0(\varphi_t), \Sigma_0(\varphi_t)) \| \mathcal{N}(0, I)) \end{aligned}$$

t : text description

z : noise vector from Gaussian Distribution

φ_t : text embedding networks (pre-trained)

\hat{c}_0 : conditioning variable

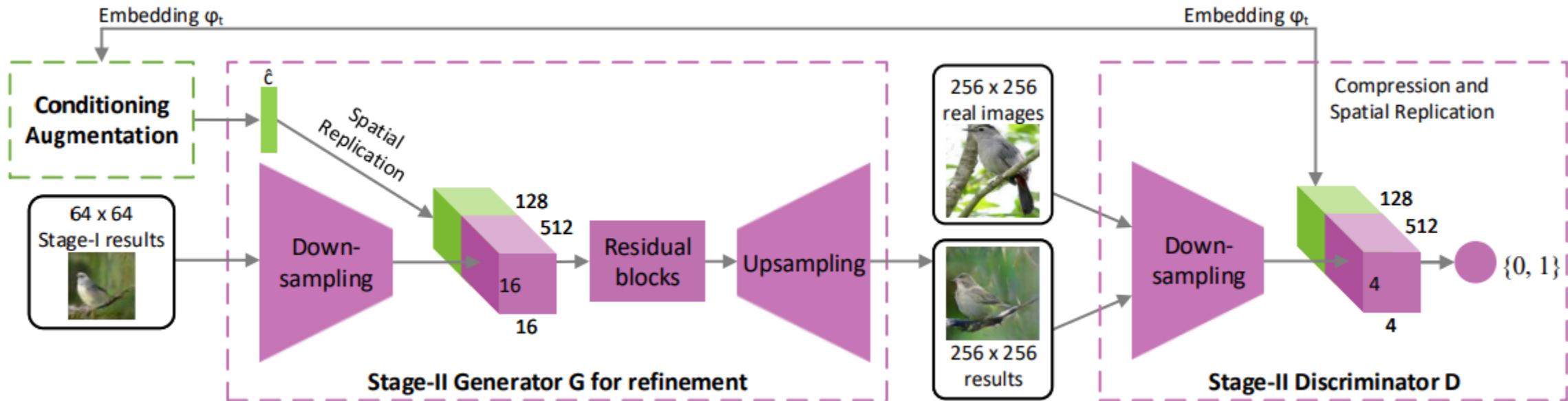
$\mathcal{N}(\mu(\varphi_t), \Sigma(\varphi_t))$: conditioning Gaussian distribution

$\mathcal{N}(0, I)$: normal distribution

$\Sigma(\varphi_t)$: diagonal covariance matrix

s_0 : image generated by the Stage-I

StackGAN



$$L_D = \mathbb{E}_{(I, t) \sim p_{data}} [\log D(I, \varphi_t)] + \\ \mathbb{E}_{s_0 \sim p_{G_0}, t \sim p_{data}} [\log(1 - D(G(s_0, \hat{c}_0), \varphi_t))]$$

$$L_G = \mathbb{E}_{s_0 \sim p_{G_0}, t \sim p_{data}} [\log(1 - D(G(s_0, \hat{c}), \varphi_t))] + \\ \lambda D_{KL}(\mathcal{N}(\mu(\varphi_t), \Sigma(\varphi_t)) \| \mathcal{N}(0, I))$$

t : *text description*

z : *noise vector from Gaussian Distribution*

φ_t : *text embedding networks (pre-trained)*

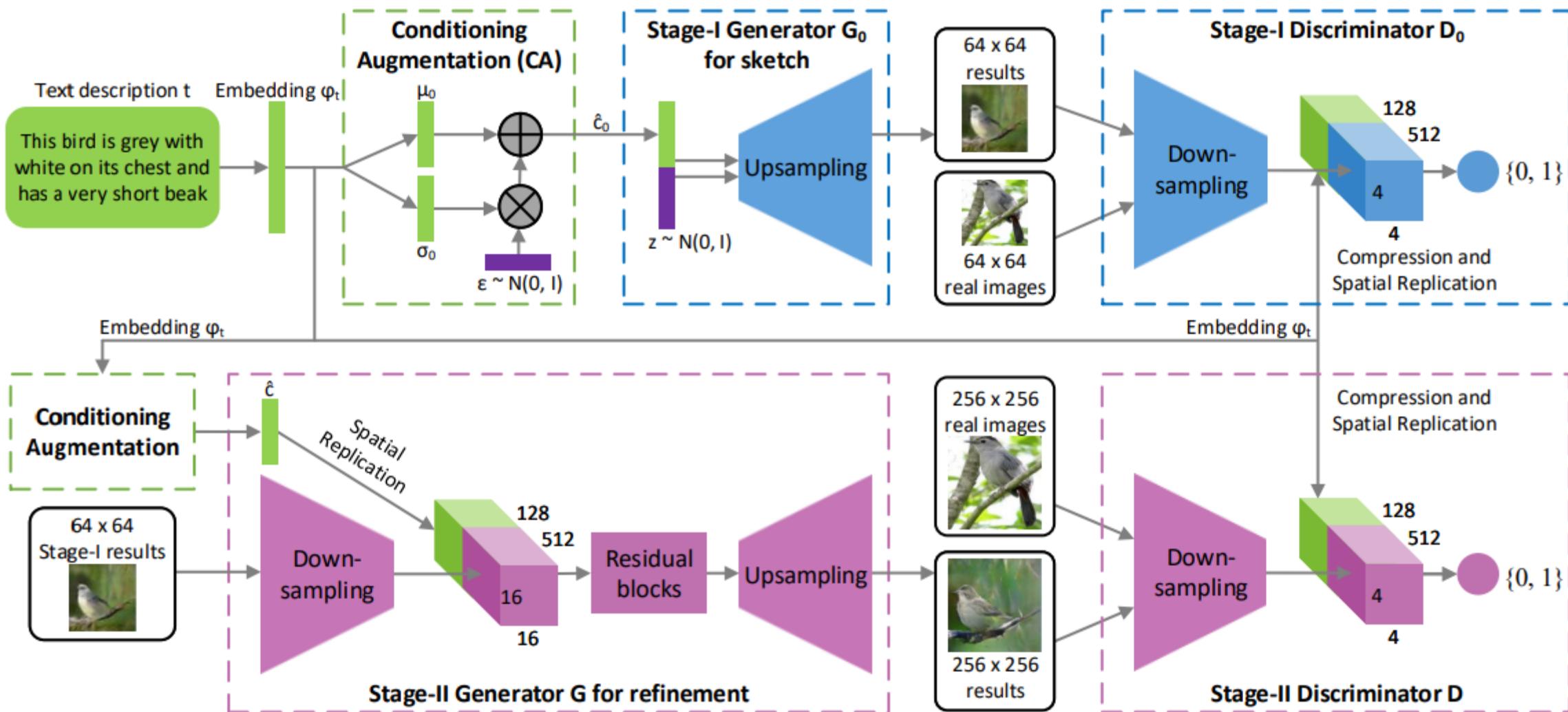
\hat{c}_0 : *conditioning variable*

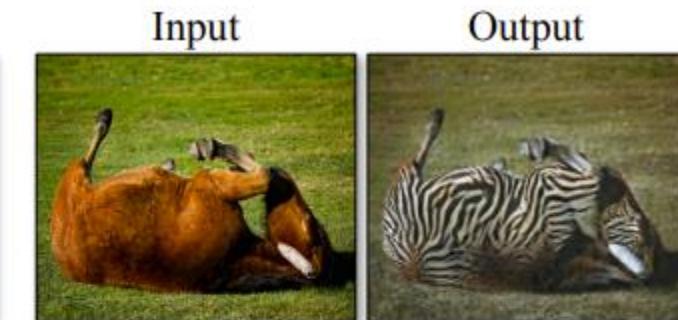
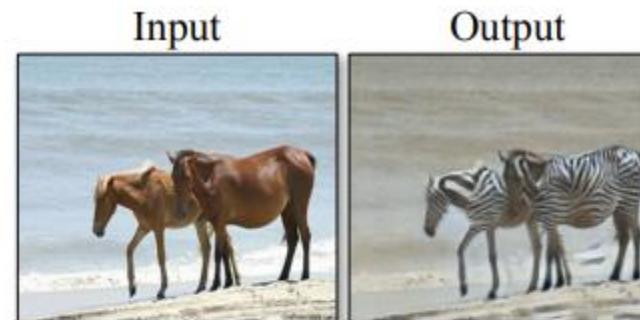
$\mathcal{N}(\mu(\varphi_t), \Sigma(\varphi_t))$: *conditioning Gaussian distribution*

$\mathcal{N}(0, I)$: *normal distribution*

$\Sigma(\varphi_t)$: *diagonal covariance matrix*

s_0 : *image generated by the Stage-I*



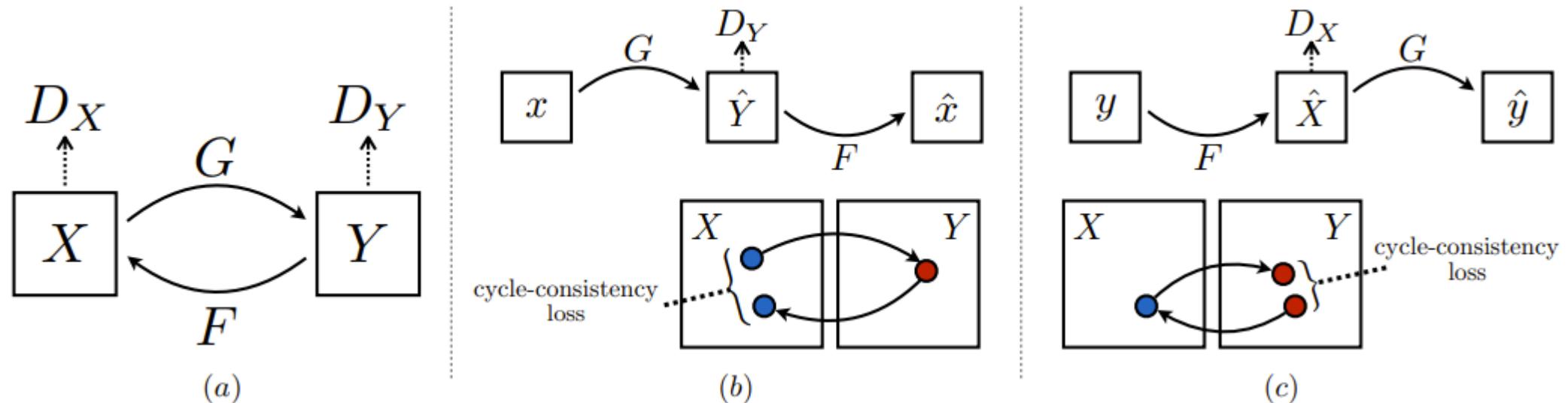


horse → zebra



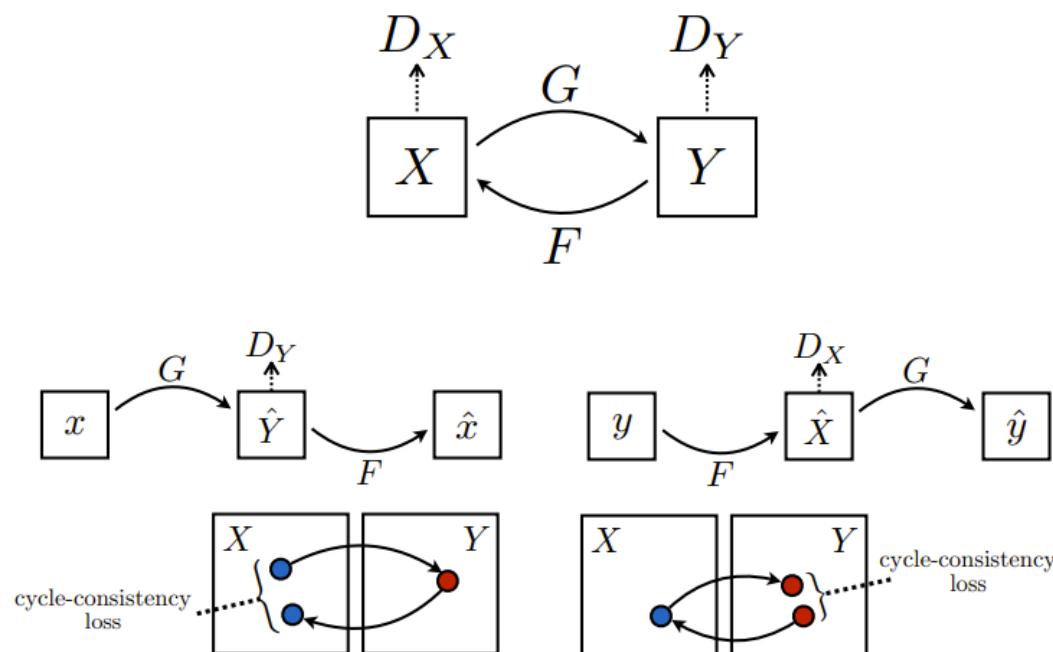
zebra → horse

Given any two unordered image collections, CycleGAN translates image to image.



$$G^*, F^* = \arg \min_{G,F} \max_{D_X D_Y} L(G, F, D_X, D_Y)$$

CycleGAN contains two mapping functions and Associated adversarial discriminator D_X and D_Y .



Adversarial loss

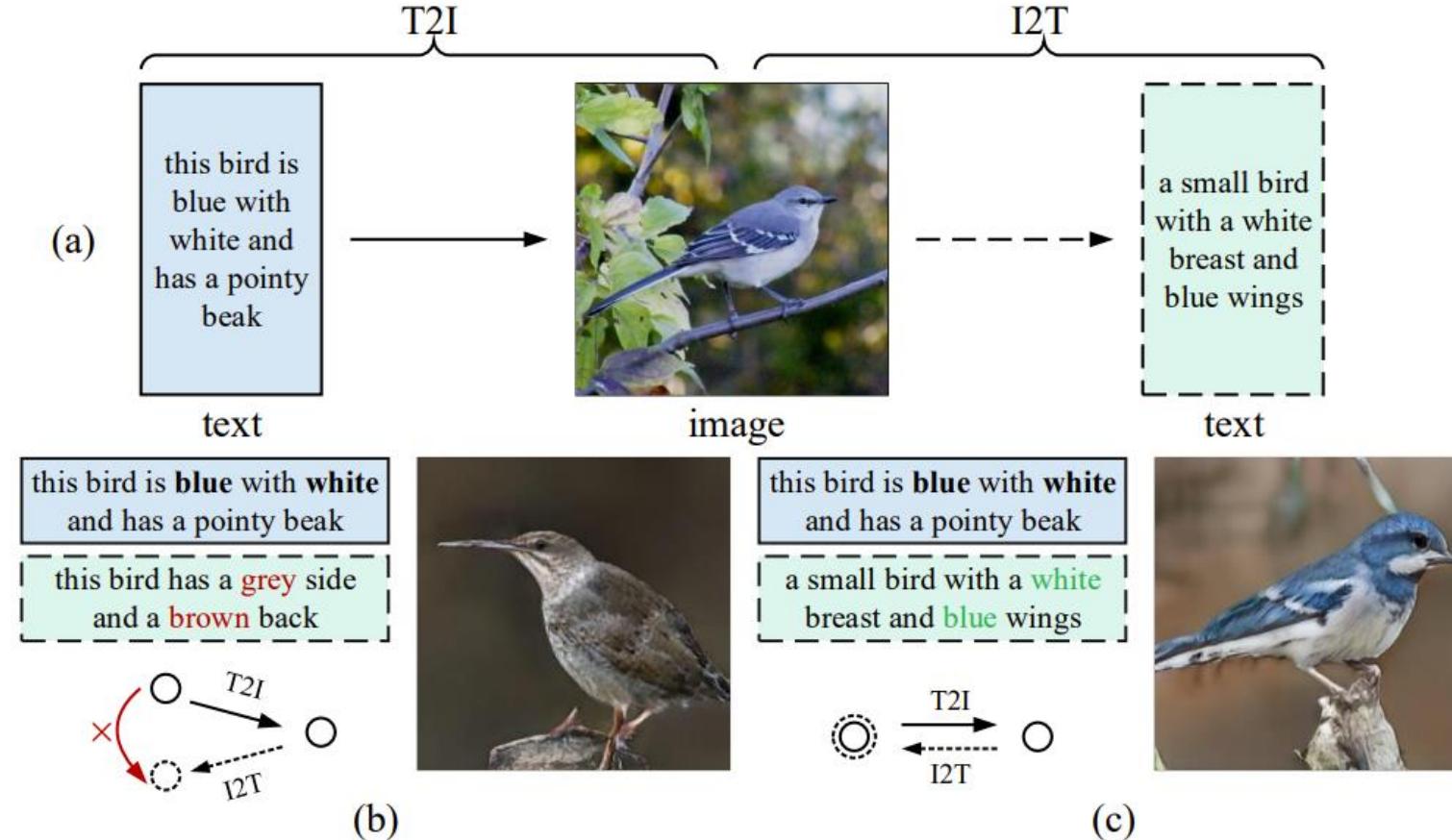
$$\begin{aligned} L_{GAN}(G, D_Y, X, Y) \\ = \mathbb{E}_{y \sim p_{data}(y)} [\log D_Y(y)] + \mathbb{E}_{x \sim p_{data}(x)} [\log D_Y(G(x))] \end{aligned}$$

Cycle consistency loss

$$\begin{aligned} L_{cyc}(G, F) \\ = \mathbb{E}_{x \sim p_{data}(x)} [\|F(G(x)) - x\|_1] + \mathbb{E}_{y \sim p_{data}(y)} [\|G(F(y)) - y\|_1] \end{aligned}$$

Full objective

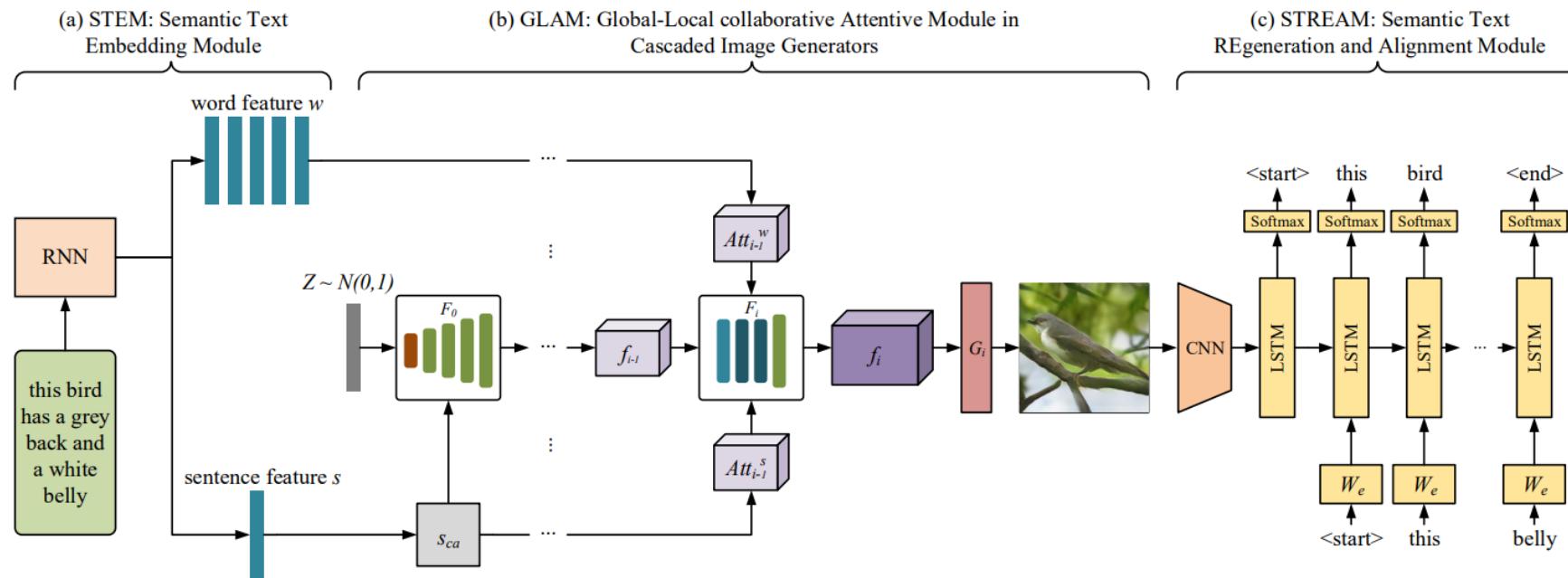
$$\begin{aligned} L(G, F, D_X, D_Y) \\ = L_{GAN}(G, D_Y, X, Y) + L_{GAN}(F, D_X, Y, X) + \lambda L_{cyc}(G, F) \end{aligned}$$



MirrorGAN has two goals : T2I and I2T.

First, T2I : An image generates from a text description.

Second, I2T : A re-description is generated from the image created by T2I.



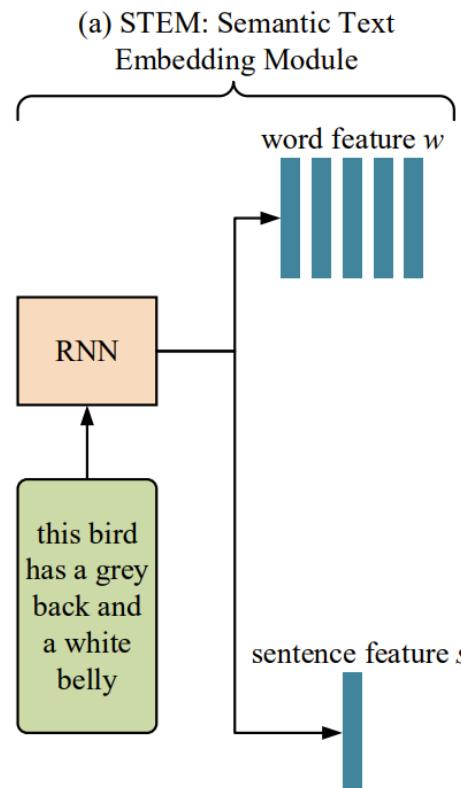
MirrorGAN has three modules.

First, STEM : Semantic Text Embedding Module.

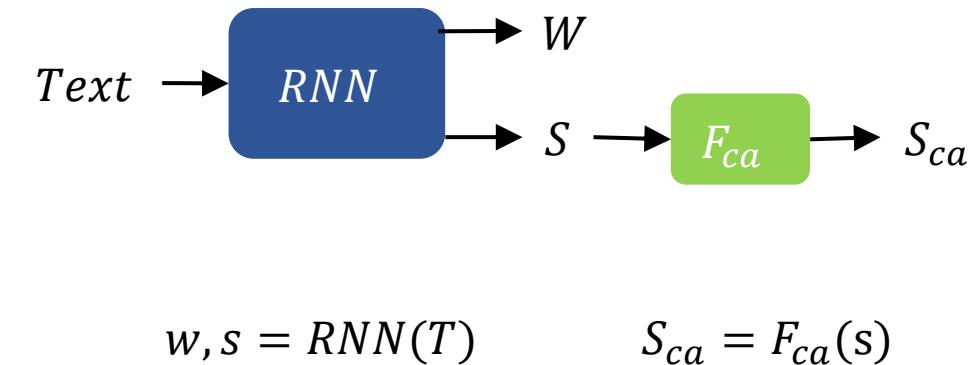
Second, GLAM : Global-Local collaborative Attentive Module in Cascaded Image Generators.

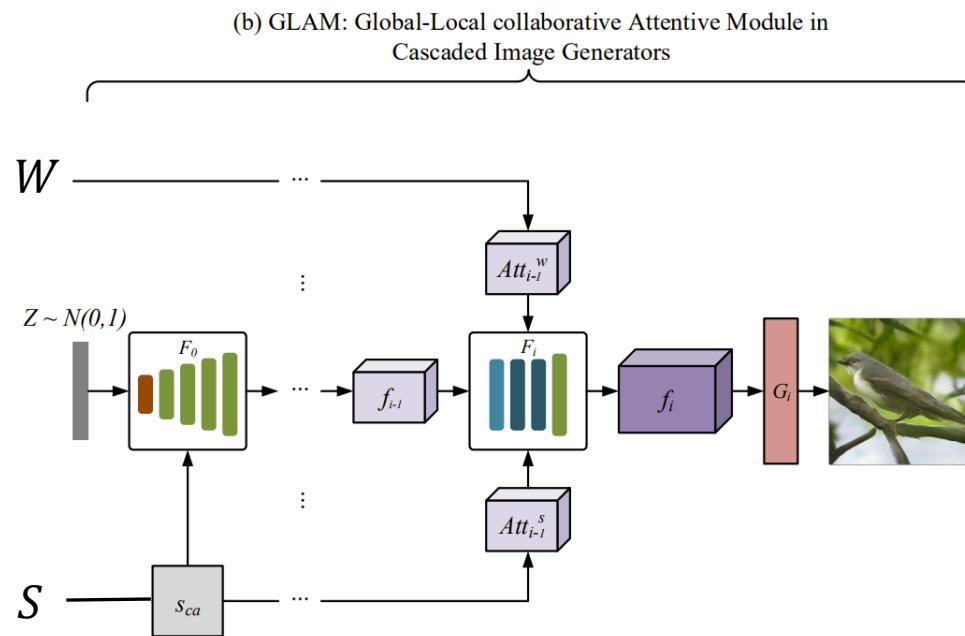
Third, STREAM : Semantic Text REgeneration and Alignment Module.

STEM : Semantic Text Embedding Module



- RNN is used to extract word and sentence features from the given text.
- Sentence features apply conditioning augmentation.



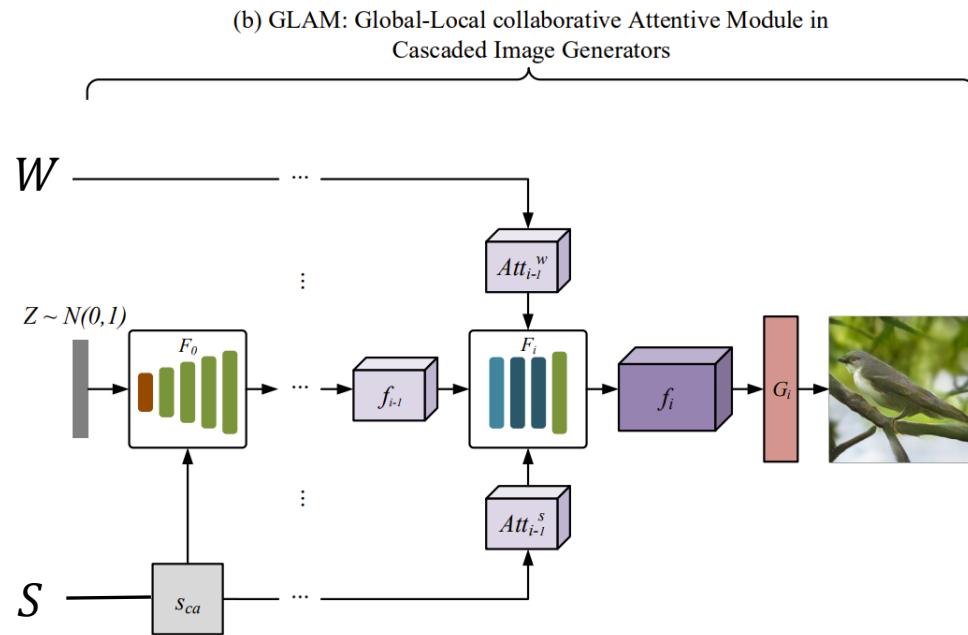


GLAM : Global-Local collaborative Attentive Module in Cascaded Image Generators

- Image generate through attention of *word* \leftrightarrow *local image* and *sentence* \leftrightarrow *global image*
- There is construct a multi-stage cascaded generator by stacking three image generation networks sequentially.

$$\text{Visual feature transormers} = \{F_0, F_1, \dots F_{m-1}\}$$

$$\text{Image generator} = \{G_0, G_1, \dots G_{m-1}\}$$



GLAM : Global-Local collaborative Attentive Module in Cascaded Image Generators

Att_{i-1}^w is context attn score between words and images.
 Att_{i-1}^s is context attn score between sentences and images.
 V_{i-1} is the term for s_{ca} mapping to a common semantic space.

$$f_0 = F_0(z, s_{ca})$$

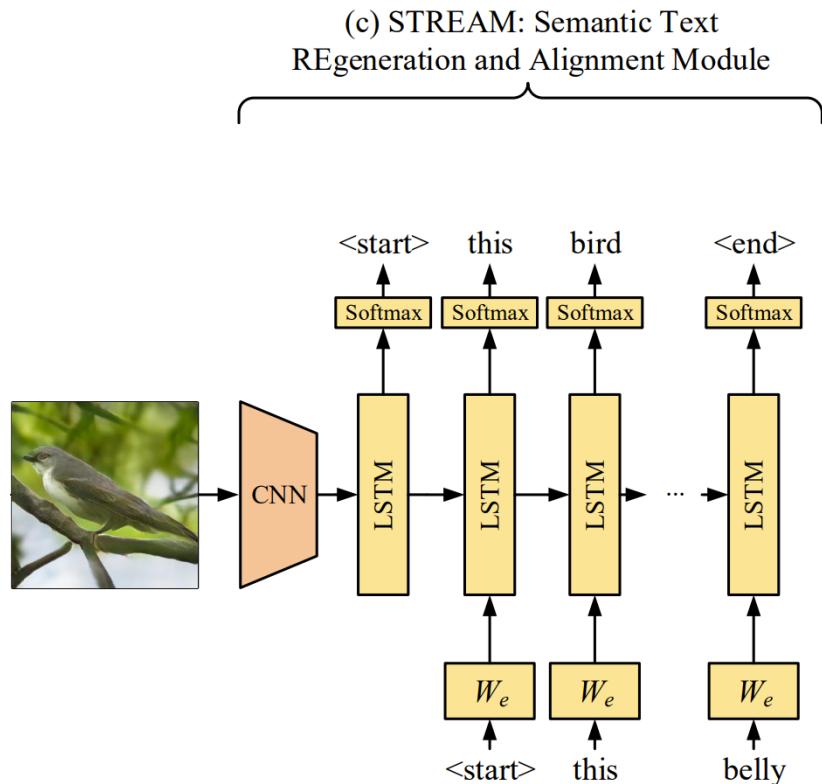
$$f_i = F_i(f_{i-1}, F_{att_i}(f_{i-1}, w, s_{ca})), i \in \{1, 2, \dots, m-1\}$$

$$I_i = G_i(f_i), i \in \{1, 2, \dots, m-1\}$$

$$Att_{i-1}^w = \sum_{l=0}^{L-1} (U_{i-1} w^l) (\text{softmax} (f_{i-1}^T (U_{i-1} w^l)))^T$$

$$Att_{i-1}^s = (V_{i-1} s_{ca}) \circ (\text{softmax} (f_{i-1} \circ (V_{i-1} s_{ca})))$$

$$F_{att_i}(f_{i-1}, w, s_{ca}) = concat(Att_{i-1}^w, Att_{i-1}^s)$$



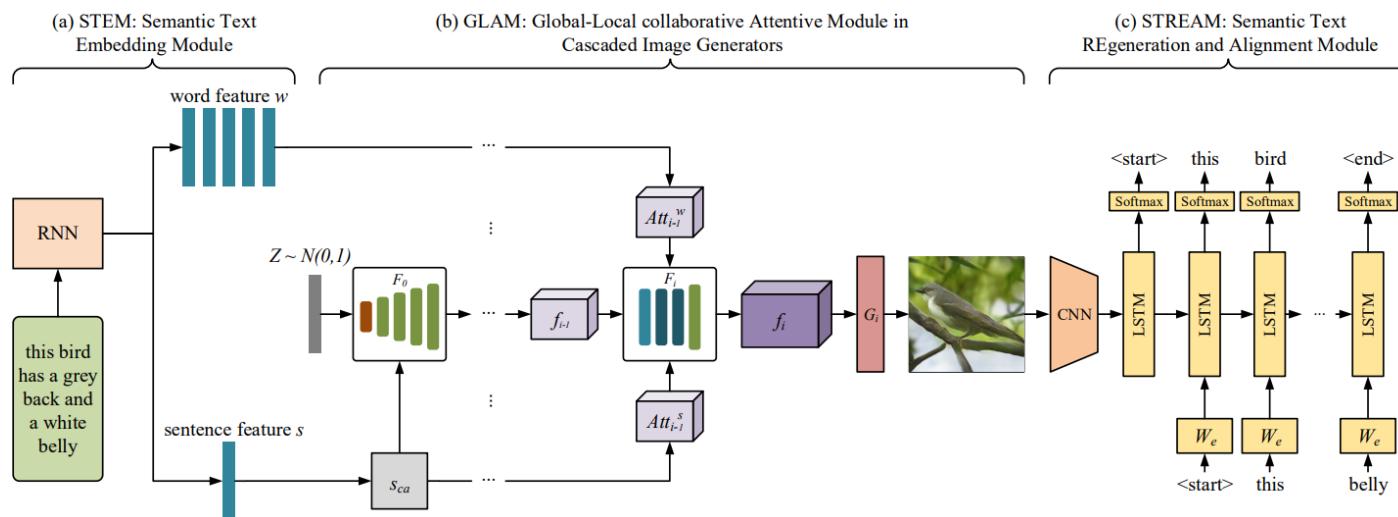
STREAM : Semantic Text Regeneration and Alignment Module

- Regenerate text description from generation image.
- CNN was pretrained on imageNet.

$$x_{-1} = CNN(I_{m-1})$$

$$x_t = W_e T_t, t \in \{0, \dots L - 1\}$$

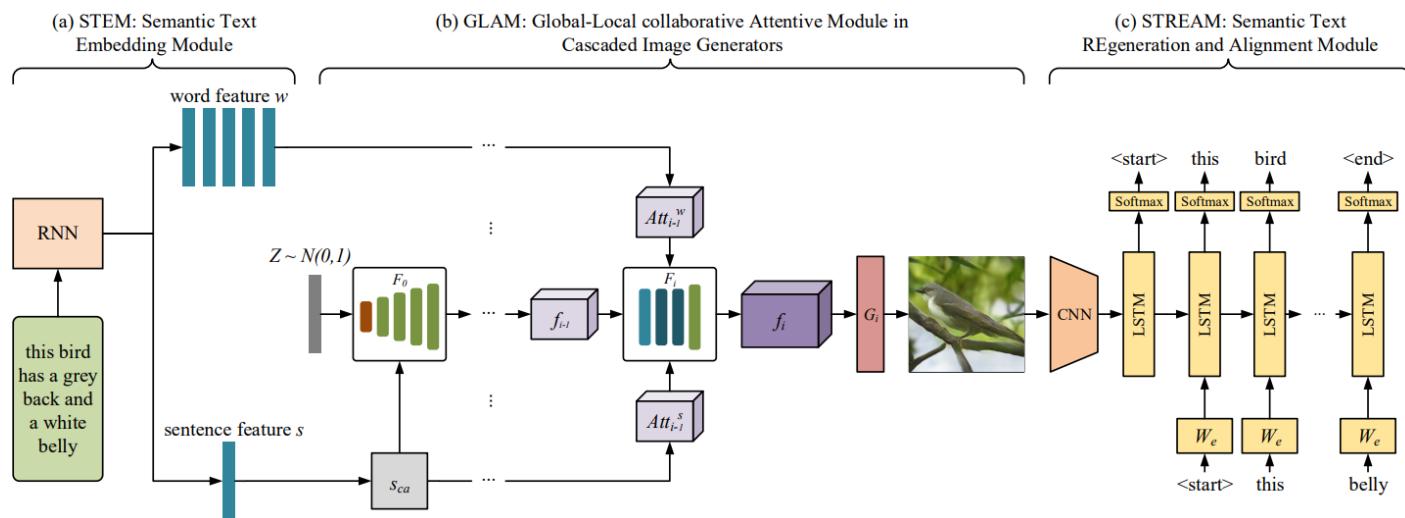
$$p_{t+1} = RNN(x_t), t \in \{0, \dots L - 1\}$$



$$\begin{aligned} L_{G_i} = & -\frac{1}{2} \mathbb{E}_{I_i \sim P_{I_i}} [\log(D_i(I_i))] \\ & -\frac{1}{2} \mathbb{E}_{I_i \sim P_{I_i}} [\log(D_i(I_i, s))] \end{aligned}$$

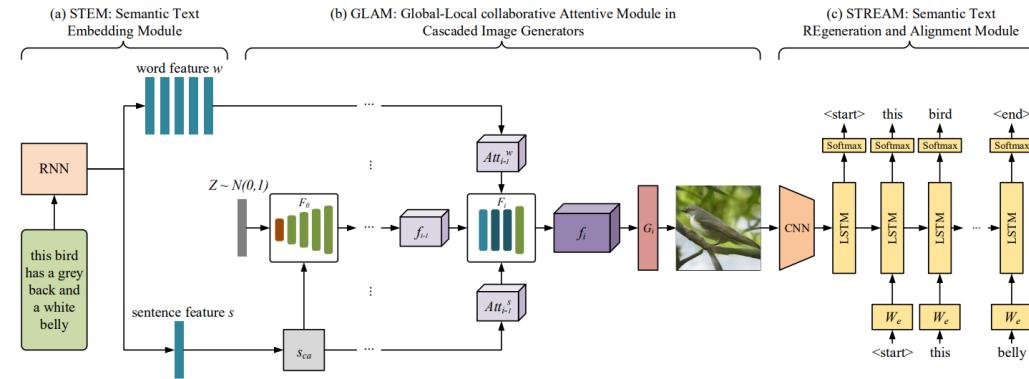
$$L_{stream} = - \sum_{t=0}^{L-1} \log p_t(T_t)$$

$$L_G = \sum_{i=0}^{m-1} L_{G_i} + \lambda L_{stream}$$



$$\begin{aligned}
 L_{D_i} = & -\frac{1}{2} \mathbb{E}_{I_i^{GT} \sim P_{I_i^{GT}}} [\log(D_i(I_i^{GT}))] \\
 & -\frac{1}{2} \mathbb{E}_{I_i \sim P_{I_i}} [\log(1 - D_i(I_i))] \\
 & -\frac{1}{2} \mathbb{E}_{I_i^{GT} \sim P_{I_i^{GT}}} [\log(D_i(I_i^{GT}, s))] \\
 & -\frac{1}{2} \mathbb{E}_{I_i \sim P_{I_i}} [\log(1 - D_i(I_i, s))]
 \end{aligned}$$

$$L_D = \sum_{i=0}^{m-1} L_{D_i}$$



$$\begin{aligned} L_{D_i} = & -\frac{1}{2} \mathbb{E}_{I_i^{GT} \sim P_{I_i^{GT}}} [\log(D_i(I_i^{GT}))] \\ & -\frac{1}{2} \mathbb{E}_{I_i \sim P_{I_i}} [\log(1 - D_i(I_i))] \\ & -\frac{1}{2} \mathbb{E}_{I_i^{GT} \sim P_{I_i^{GT}}} [\log(D_i(I_i^{GT}, s))] \\ & -\frac{1}{2} \mathbb{E}_{I_i \sim P_{I_i}} [\log(1 - D_i(I_i, s))] \end{aligned}$$

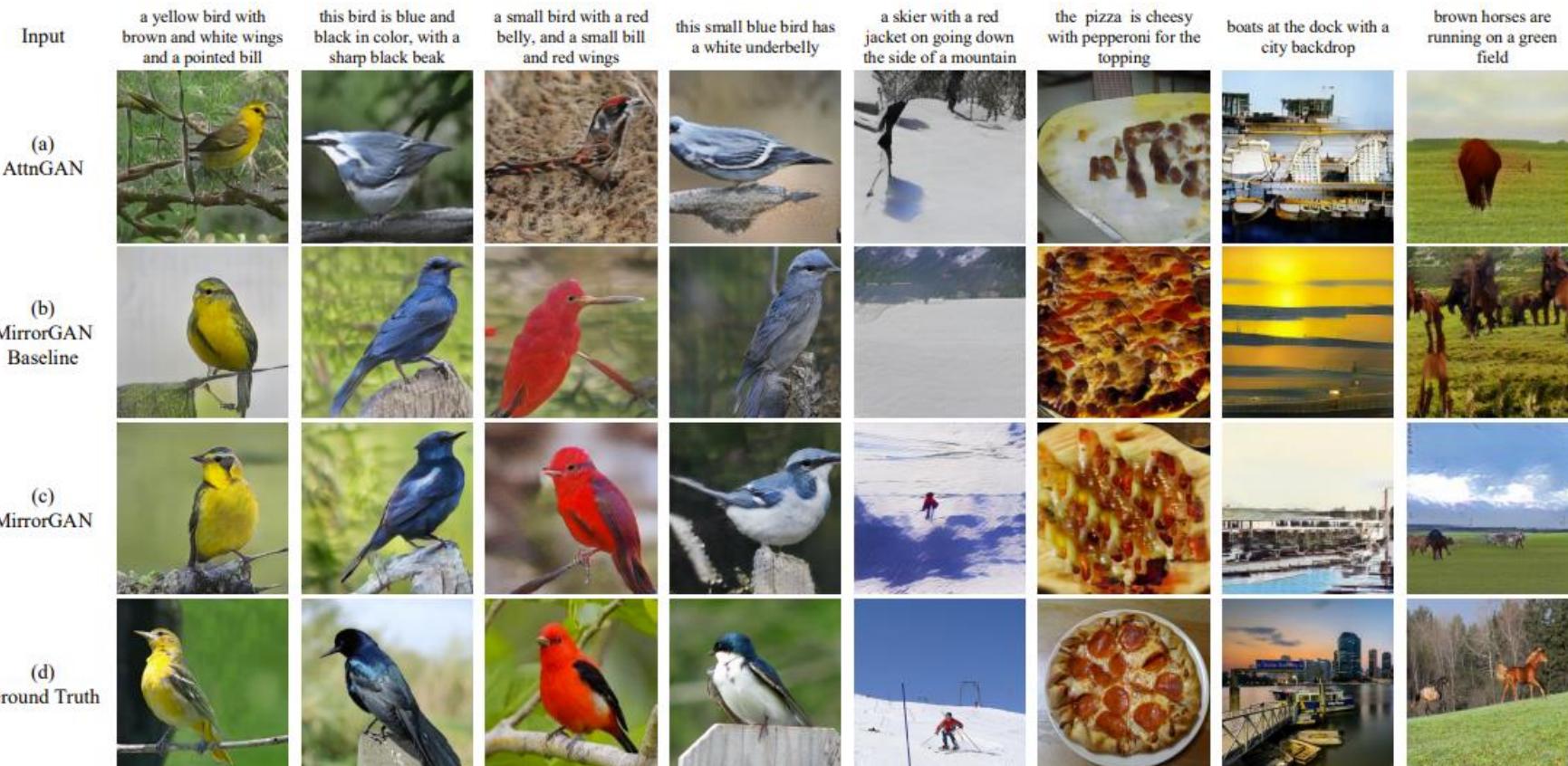
$$L_D = \sum_{i=0}^{m-1} L_{D_i}$$

$$L_{G_i} = -\frac{1}{2} \mathbb{E}_{I_i \sim P_{I_i}} [\log(D_i(I_i))] - \frac{1}{2} \mathbb{E}_{I_i \sim P_{I_i}} [\log(D_i(I_i, s))]$$

$$L_{stream} = -\sum_{t=0}^{L-1} \log p_t(T_t)$$

$$L_G = \sum_{i=0}^{m-1} L_{G_i} + \lambda L_{stream}$$

Result



5

Result

Statistics of datasets

Dataset	CUB [28]		COCO [14]	
	train	test	train	test
#samples	8,855	2,933	80k	40k
caption/image	10	10	5	5

Inception Score

Model	CUB	COCO
GAN-INT-CLS [24]	2.88 ± 0.04	7.88 ± 0.07
GAWWN [25]	3.62 ± 0.07	-
StackGAN [38]	3.70 ± 0.04	8.45 ± 0.03
StackGAN++ [39]	3.82 ± 0.06	-
PPGN [20]	-	9.58 ± 0.21
AttnGAN [34]	4.36 ± 0.03	25.89 ± 0.47
MirrorGAN	4.56 ± 0.05	26.47 ± 0.41

Inception Score evaluate quality and diversity from generated image.

R-precision

Dataset	CUB			COCO			
	top-k	k=1	k=2	k=3	k=1	k=2	k=3
AttnGAN [34]		53.31	54.11	54.36	72.13	73.21	76.53
MirrorGAN		57.67	58.52	60.42	74.52	76.87	80.21

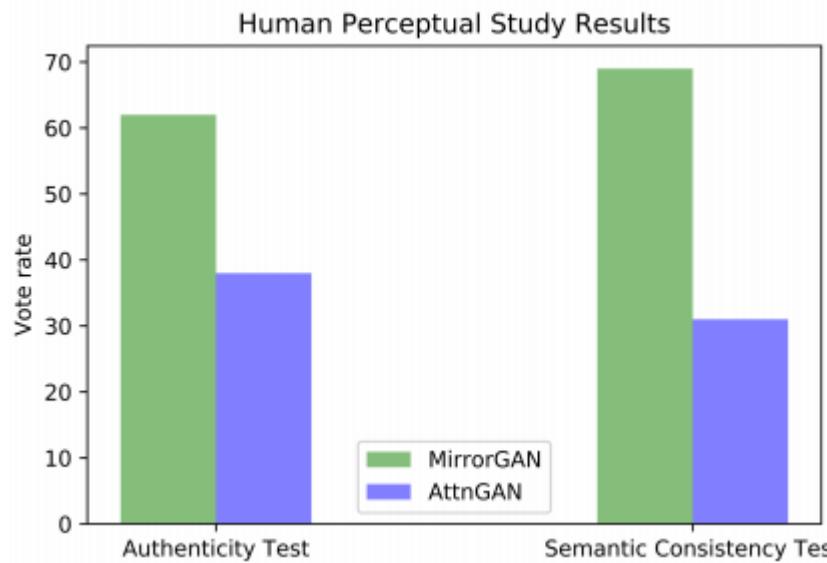
There is cosine similarities between the generated image vectors and text vectors.

Evaluation Metric	Inception Score		R-precision (top-1)	
	CUB	COCO	CUB	COCO
MirrorGAN w/o GA, $\lambda=0$	$3.91 \pm .09$	$19.01 \pm .42$	39.09	50.69
MirrorGAN w/o GA, $\lambda=20$	$4.47 \pm .07$	$25.99 \pm .41$	55.67	73.28
MirrorGAN, $\lambda=5$	$4.01 \pm .06$	$21.85 \pm .43$	32.07	52.55
MirrorGAN, $\lambda=10$	$4.30 \pm .07$	$24.11 \pm .31$	43.21	63.40
MirrorGAN, $\lambda=20$	$4.54 \pm .17$	$26.47 \pm .41$	57.67	74.52

GA = Global Attention

λ expresses the importance of STREAM

Result



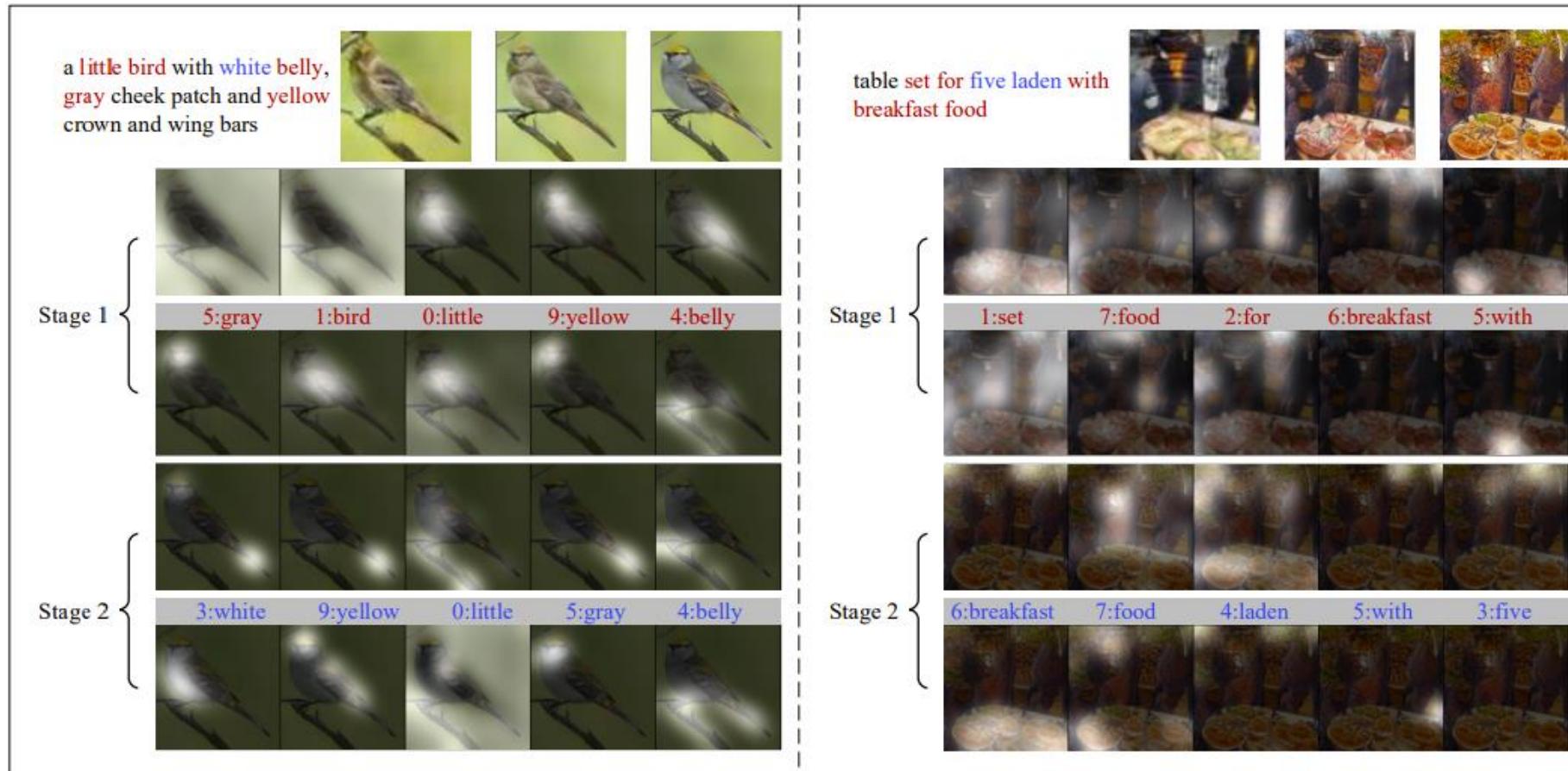
Results of Human perceptual test

A higher value of the Authenticity Test means more convincing images.

A higher value of the Semantic Consistency Test means a closer semantics between input text and generated images.

5

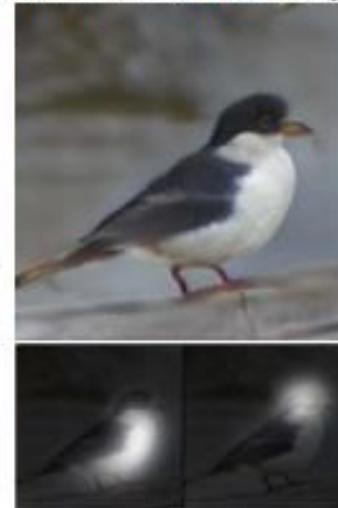
Result



this bird has a **yellow** crown and a **white** belly



this bird has a **black** crown and a **white** belly



this bird has a **black** crown and a **red** belly



this bird has **blue** wings and a **red** belly



Images generated by MirrorGAN by modifying the text descriptions by a single word and the corresponding top-2 attention maps in the last stage

Reference

- <https://arxiv.org/pdf/1903.05854v1.pdf>
- <https://arxiv.org/pdf/1612.03242.pdf>
- <https://arxiv.org/pdf/1711.10485.pdf>
- <https://arxiv.org/pdf/1703.10593.pdf>
- <https://cs.stanford.edu/people/karpathy/cvpr2015.pdf>
- <https://drive.google.com/file/d/1luXzEwW1NYX0s2SOuTLMetnIv0I5RR7M/view>