# Generative Adversarial Text to Image Synthesis

Reed S, Akata Z, Yan X, et al. ICML 2016

# Text2Image

**DCGAN + CGAN**，文字描述由一个现成的词嵌入方法生成向量*

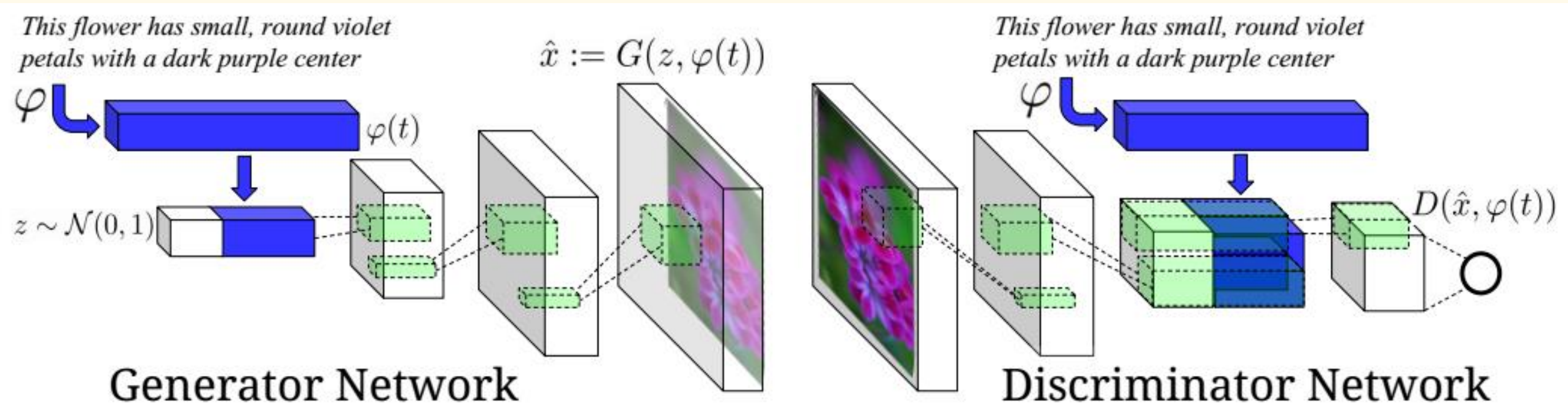$$\min_{G} \max_{D} V(D,G) = \mathbb{E}_{x\sim p_{data}(x)}[\log D(x)]+$$
$$\mathbb{E}_{x\sim p_z(z)}[\log(1 - D(G(z)))]$$



* Reed, S., Akata, Z., Lee, H., and Schiele, B. Learning deep representations for fine-grained visual descriptions. In CVPR, 2016.

# Text2Image

**DCGAN + CGAN**，文字描述**t**由一个现成的词嵌入方法生成向量$\varphi(t)$

| | | | |
|---|---|---|---|
| GT | an all black bird with a distinct thick, rounded bill. | this small bird has a yellow breast, brown crown, and black superciliary | a tiny bird, with a tiny beak, tarsus and feet, a blue crown, blue coverts, and black cheek patch |
| GAN | | | |

# GAN-CLS 通过加入反例来训练D(与G无关)

**Intuition: the discriminator has no explicit notion of whether real training images match the text embedding context.**

**Algorithm 1** GAN-CLS training algorithm with step size $\alpha$, using minibatch SGD for simplicity.

1: **Input:** minibatch images $x$, matching text $t$, mismatching $\hat{t}$, number of training batch steps $S$
2: **for** $n = 1$ **to** $S$ **do**
3:      $h \leftarrow \varphi(t)$ {Encode matching text description}
4:      $\hat{h} \leftarrow \varphi(\hat{t})$ {Encode mis-matching text description}
5:      $z \sim \mathcal{N}(0,1)^Z$ {Draw sample of random noise}
6:      $\hat{x} \leftarrow G(z, h)$ {Forward through generator}
7:      $s_r \leftarrow D(x, h)$ {real image, right text}
8:      $s_w \leftarrow D(x, \hat{h})$ {real image, wrong text}
9:      $s_f \leftarrow D(\hat{x}, h)$ {fake image, right text}
10:     $\mathcal{L}_D \leftarrow \log(s_r) + (\log(1 - s_w) + \log(1 - s_f))/2$
11:     $D \leftarrow D - \alpha \partial \mathcal{L}_D / \partial D$ {Update discriminator}
12:     $\mathcal{L}_G \leftarrow \log(s_f)$
13:     $G \leftarrow G - \alpha \partial \mathcal{L}_G / \partial G$ {Update generator}
14: **end for**

一张图片和他的一个文字描述是一个样本对

**GT**

an all black bird with a distinct thick, rounded bill.

this small bird has a yellow breast, brown crown, and black superciliary

a tiny bird, with a tiny beak, tarsus and feet, a blue crown, blue coverts, and black cheek patch

**GAN**

**GAN - CLS**

**Deep networks have been shown to learn representations in which interpolations between embedding pairs tend to be near the data manifold**

$$\mathbb{E}_{t_1, t_2 \sim p_{data}}[\log(1 - D(G(z, \beta t_1 + (1 - \beta)t_2)))]$$

**Because the interpolated embeddings are synthetic, the discriminator D does not have "real" corresponding image and text pairs to train on. However, D learns to predict whether image and text pairs match or not. Thus, if D does a good job at this, then by satisfying D on interpolated text embeddings G can learn to fill in gaps on the data manifold in between training points.**

原来的样本数量不足以使得判别器**D**判断图文的关系，通过插值法增加样本数使得**D**能学得他们之间的关系。

**Algorithm 1** GAN-CLS training algorithm with step size $\alpha$, using minibatch SGD for simplicity.

1: **Input:** minibatch images $x$, matching text $t$, mis-matching $\hat{t}$, number of training batch steps $S$
2: **for** $n = 1$ **to** $S$ **do**
3:     $h \leftarrow \varphi(t)$ {Encode matching text description}
4:     $\hat{h} \leftarrow \varphi(\hat{t})$ {Encode mis-matching text description}
5:     $z \sim \mathcal{N}(0, 1)^Z$ {Draw sample of random noise}
6:     $\hat{x} \leftarrow G(z, h)$ {Forward through generator}
7:     $s_r \leftarrow D(x, h)$ {real image, right text}
8:     $s_w \leftarrow D(x, \hat{h})$ {real image, wrong text}
9:     $s_f \leftarrow D(\hat{x}, h)$ {fake image, right text}
10:    $\mathcal{L}_D \leftarrow \log(s_r) + (\log(1 - s_w) + \log(1 - s_f))/2$
11:    $D \leftarrow D - \alpha \partial \mathcal{L}_D / \partial D$ {Update discriminator}
12:    $\mathcal{L}_G \leftarrow \log(s_f)$
13:    $G \leftarrow G - \alpha \partial \mathcal{L}_G / \partial G$ {Update generator}
14: **end for**

差值得到的文字描述当作真实的对应文本

GT

an all black bird with a distinct thick, rounded bill.

this small bird has a yellow breast, brown crown, and black superciliary

a tiny bird, with a tiny beak, tarsus and feet, a blue crown, blue coverts, and black cheek patch

GAN

GAN - CLS

GAN - INT

GAN - INT - CLS

| | | | |
|---|---|---|---|
| **GT** | this flower is white and pink in color, with petals that have veins. | these flowers have petals that start off white in color and end in a dark purple towards the tips. | bright droopy yellow petals with burgundy streaks, and a yellow stigma. | a flower with long pink petals and raised orange stamen. |

**GAN**

**GAN - CLS**

**GAN - INT**

**GAN - INT - CLS**

**CUB contains 200 bird species with 11,788 images.**
**Oxford-102 contains 8,189 images of flowers from 102 different categories.**

**Solution: train a convolutional network S to invert G to regress from samples x ˆ back onto z.**

$$\mathcal{L}_{style} = \mathbb{E}_{t,z\sim\mathcal{N}(0,1)}\big|\big|z - S(G(z,\varphi(t)))\big|\big|_2^2$$

where $S$ is the style encoder network. With a trained generator and style encoder, style transfer from a query image $x$ onto text $t$ proceeds as follows:

$$s \leftarrow S(x), \hat{x} \leftarrow G(s,\varphi(t))$$

where $\hat{x}$ is the result image and $s$ is the predicted style.

**Text descriptions (content)** — **Images (style)**

$$s \leftarrow S(x)$$

$$\hat{x} \leftarrow G(s, \varphi(t))$$

The bird has a **yellow breast** with **grey** features and a small beak.

This is a large **white** bird with **black wings** and a **red head**.

A small bird with a **black head and wings** and features grey wings.

This bird has a **white breast**, brown and white coloring on its head and wings, and a thin pointy beak.

A small bird with **white base** and **black stripes** throughout its belly, head, and feathers.

A small sized bird that has a cream belly and a short pointed bill.

This bird is **completely red**.

This bird is **completely white**.

This is a **yellow** bird. The **wings are bright blue**.

'Blue bird with black beak' → 'Red bird with black beak'
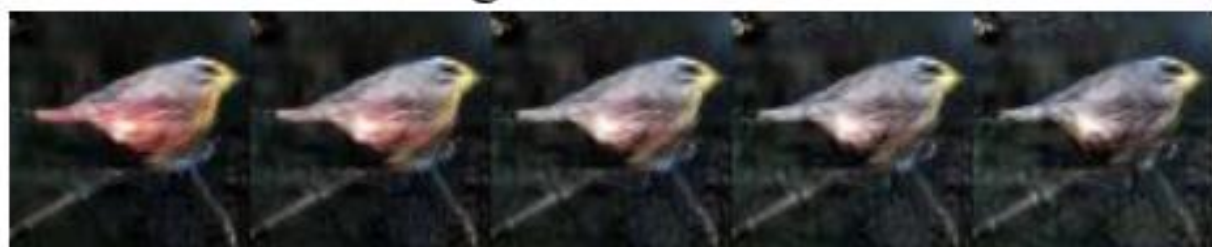
'This bird is completely red with black wings'

'Small blue bird with black wings' → 'Small yellow bird with black wings'

'this bird is all blue, the top part of the bill is blue, but the bottom half is white'

'This bird is bright.' → 'This bird is dark.'

'This is a yellow bird. The wings are bright blue'

# Problem

**the generated scenes are not usually coherent; for example the human-like blobs in the baseball scenes lack clearly articulated parts.**

**Low resolution (64x64像素)**

# StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks
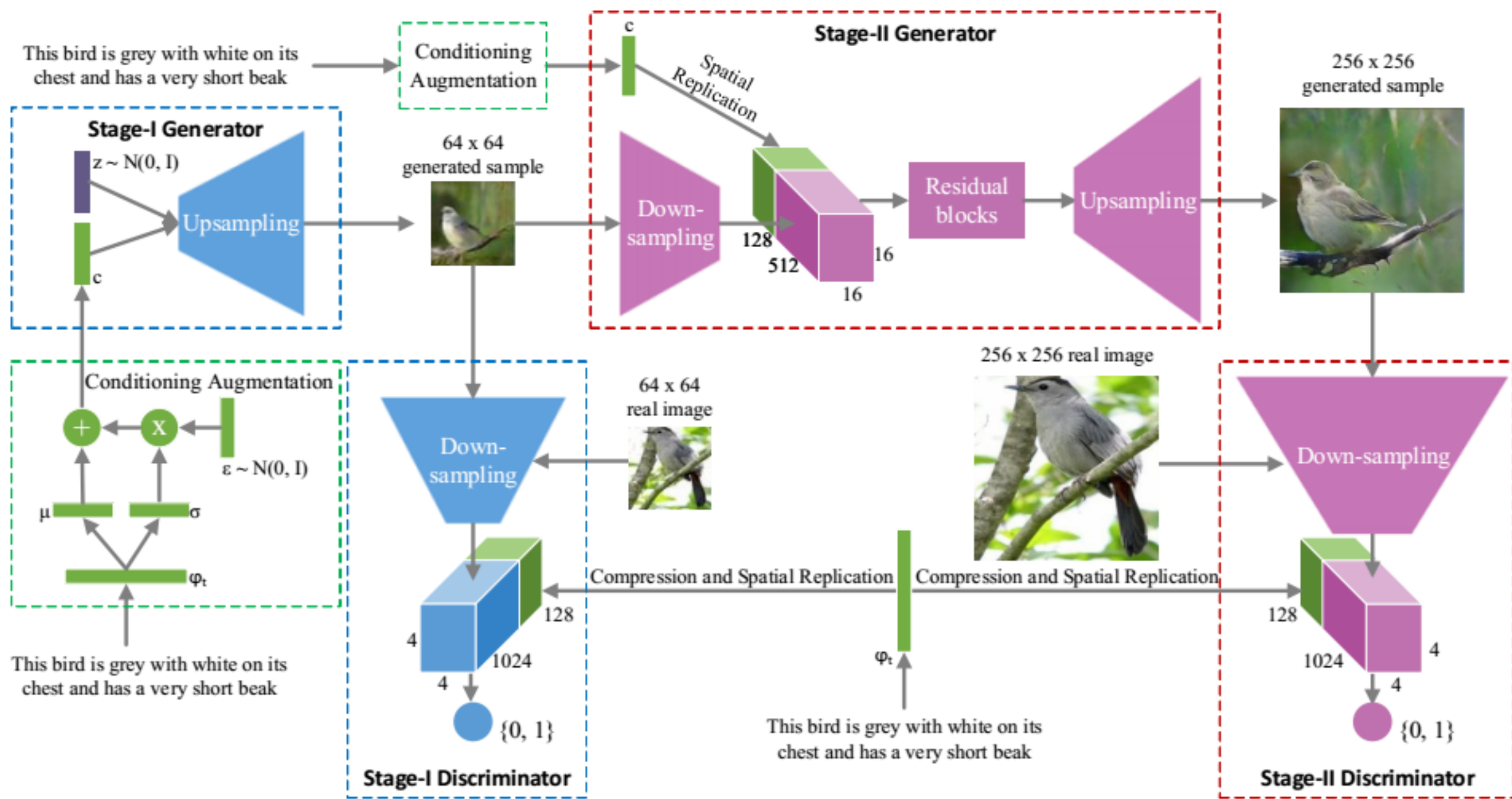
Zhang H, Xu T, Li H, et al.
arXiv preprint. 2016.

**the space of plausible images given text descriptions is multimodal. There are a large number of images that correctly fit the given text description**

给定文字描述作为条件，**输入是噪声z拼上文字向量**。一阶段只生成大概的轮廓和颜色的低分辨率图像

给定文字描述和一阶段的图像作为输入，二阶段生成细节和对一阶段改进

we adopt the GAN-CLS for both stages.

# Experiment

| Text description | This bird is red and brown in color, with a stubby beak | The bird is short and stubby with yellow on its body | A bird with a medium orange bill white body gray wings and webbed feet | This small black bird has a short, slightly curved bill and long legs | A small bird with varying shades of brown with white under the eyes | A small yellow bird with a black crown and a short black pointed beak | This small bird has a white breast, light grey head, and black wings and tail |

64x64 GAN-INT-CLS [22]

256x256 StackGAN

**CUB contains 200 bird species with 11,788 images.**
**Oxford-102 contains 8,189 images of flowers from 102 different categories.**

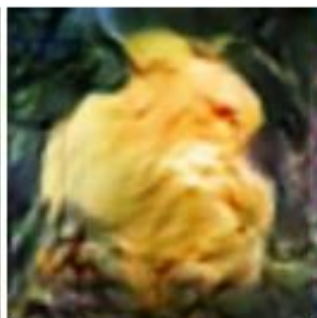| Text description | This flower has petals that are white and has pink shading | This flower has a lot of small purple petals in a dome-like configuration | This flower has long thin yellow petals and a lot of yellow anthers in the center | This flower is pink, white, and yellow in color, and has petals that are striped | This flower is white and yellow in color, with petals that are wavy and smooth | This flower has upturned petals which are thin and orange with rounded edges | This flower has petals that are dark pink with white edges and pink stamen |
|---|---|---|---|---|---|---|---|
| 64x64 GAN-INT-CLS [22] | | | | | | | |
| 256x256 StackGAN | | | | | | | |

一阶段和二阶段生成的图像有显著进步，说明分阶段生成图片的可行性

通过将生成图片从训练
集中找出5个最近邻，
发现生成的图片和训练
集中的图片差别较大，

说明了**stackGAN** 不是
记住了训练样本，而是
找到了语言和图像之间
的联系（更强的泛化能
力）



Figure 6. For generated images (column 1), retrieving their nearest training images (columns 2-6) by utilizing Stage-II discriminator $D$ to extract visual features. The $L2$ distances between features are calculated for nearest-neighbor retrieval.

由于类似**VAE**组建的存在，对输入的文字描述加入了一定的扰动，使得即使固定噪声和文字输入，也能输出多样化的图片



Figure 7. Birds with different poses and viewpoints generated with the same input text embedding by our StackGAN. The noise vector $z$ and text embedding are fixed for each row.

To demonstrate that our StackGAN learns a smooth latent data manifold, we use it to generate images from linearly interpolated sentence embeddings
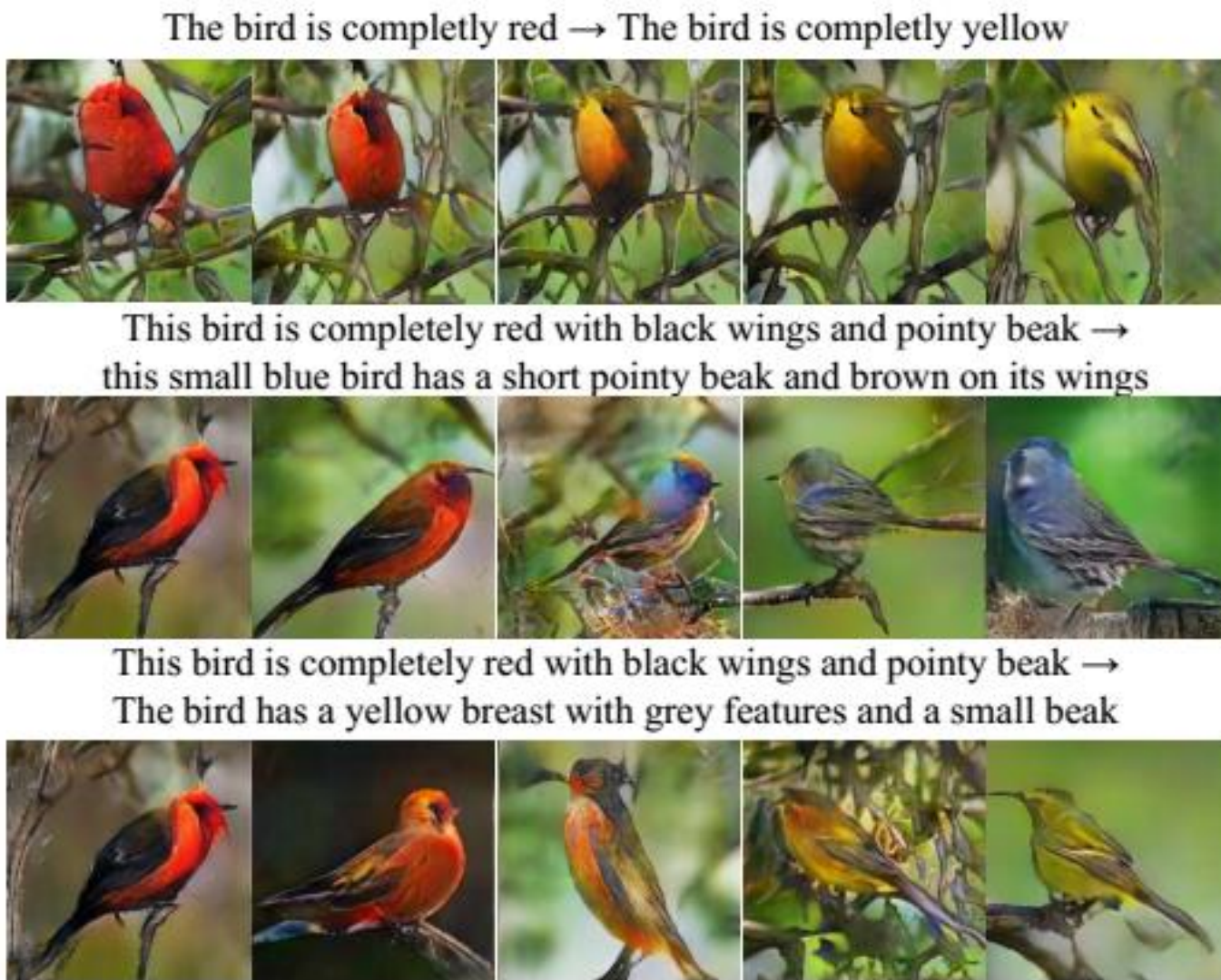


The bird is completely red → The bird is completly yellow

This bird is completely red with black wings and pointy beak → this small blue bird has a short pointy beak and brown on its wings

This bird is completely red with black wings and pointy beak → The bird has a yellow breast with grey features and a small beak

Figure 8. (Left to right) Images generated by interpolating two sentence embeddings. Gradual appearance changes from the first sentence's meaning to that of the second sentence can be observed. The noise vector $z$ is fixed to be zeros for each row.

More example

This bird sits close to the ground with his short yellow tarsus and feet; his bill is long and is also yellow and his color is mostly white with a black crown and primary feathers



A large bird has large thighs and large wings that have white wingbars

The small bird has a red head with feathers that fade from red to gray from head to tail



Stage-I images

Stage-II images

This bird is black with green and has a very short beak



Stage-I images

Stage-II images

## This flower has white petals with a yellow tip and a yellow pistil

Stage-I images

Stage-II images

## A flower with small pink petals and a massive central orange and black stamen cluster

Stage-I images

Stage-II images

# Failure cases

**The main reason for failure cases is that Stage-I GAN fails to generate plausible rough shapes or colors of the objects**



| Text description | This particular bird has a brown body and brown bill | Grey bird with black flat beak with grey and white big wings | Bird has brown body feathers, brown breast feathers, and brown beak | The medium sized bird has a dark grey color, a black downward curved beak, and long wings | Colored bill with a white ring around it on the upper part near the bill | This bird has a dark brown overall body color, with a small white patch around the base of the bill | This medium sized bird is primarily black and has a large wingspan and a long black bill with a strip of white at the beginning of it |

# Failure cases

**The main reason for failure cases is that Stage-I GAN fails to generate plausible rough shapes or colors of the objects**

Thank you