

StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks

Han Zhang^{*1}, Tao Xu², Hongsheng Li³, Shaoting Zhang⁴,
Xiaolei Huang², Xiaogang Wang³, Dimitris Metaxas¹

¹Department of Computer Science, Rutgers University

²Department of Computer Science and Engineering, Lehigh University

³Department of Electronic Engineering, The Chinese University of Hong Kong

⁴Department of Computer Science, University of North Carolina at Charlotte

Abstract

Synthesizing photo-realistic images from text descriptions is a challenging problem in computer vision and has many practical applications. Samples generated by existing text-to-image approaches can roughly reflect the meaning of the given descriptions, but they fail to contain necessary details and vivid object parts. In this paper, we propose stacked Generative Adversarial Networks (StackGAN) to generate photo-realistic images conditioned on text descriptions. The Stage-I GAN sketches the primitive shape and basic colors of the object based on the given text description, yielding Stage-I low resolution images. The Stage-II GAN takes Stage-I results and text descriptions as inputs, and generates high resolution images with photo-realistic details. The Stage-II GAN is able to rectify defects and add compelling details with the refinement process. Samples generated by StackGAN are more plausible than those generated by existing approaches. Importantly, our StackGAN for the first time generates realistic 256×256 images conditioned on only text descriptions, while state-of-the-art methods can generate at most 128×128 images. To demonstrate the effectiveness of the proposed StackGAN, extensive experiments are conducted on CUB and Oxford-102 datasets, which contain enough object appearance variations and are widely-used for text-to-image generation analysis.

1. Introduction

Generating photo-realistic images from text descriptions is a challenging problem. It has tremendous applications including photo editing and computer-aided design when

^{*}han.zhang@cs.rutgers.edu

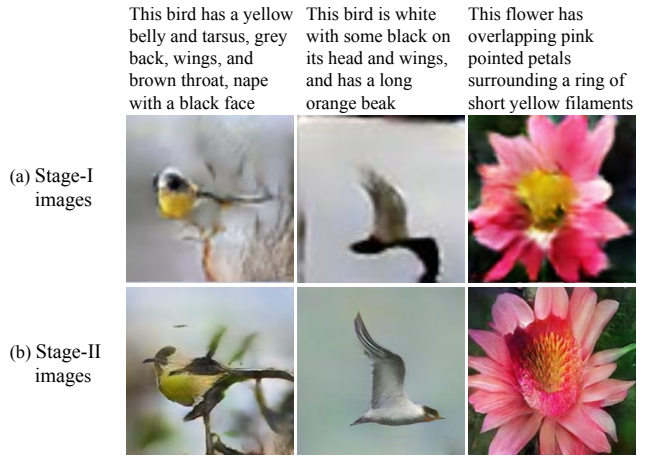


Figure 1. Photo-realistic images generated by our StackGAN from unseen text descriptions. Descriptions for birds and flowers are from CUB [32] and Oxford-102 [18] datasets, respectively. (a) Given text descriptions, Stage-I of StackGAN sketches rough shapes and basic colors of objects, yielding low resolution images. (b) Stage-II of StackGAN takes Stage-I results and text descriptions as inputs, and generates high resolution images with photo-realistic details.

fully automatic synthesis systems are available. However, even the most advanced methods failed in generating high resolution images with photo-realistic details using text descriptions. The main challenge of this problem is that the space of plausible images given text descriptions is multimodal. There are a large number of images that correctly fit the given text description.

Recently, Generative Adversarial Networks (GAN) [7, 3, 19] have shown promising results in modeling complex multimodal data and synthesizing real-world images. Reed *et al.* [22, 20] demonstrated that GAN can effectively gen-

erate images conditioned on text descriptions. They succeeded in generating plausible 64×64 images conditioned on text descriptions [22]. However, their synthesized images in many cases lack details and vivid object parts, *e.g.*, beaks and eyes of birds. Moreover, they were unable to synthesize higher resolution images (*e.g.*, 128×128) without providing additional spatial annotations of objects.

To tackle the challenges, we decompose the problem of text to photo-realistic image synthesis into two more manageable sub-problems with stacked Generative Adversarial Networks (StackGAN). A low resolution image is generated using our Stage-I GAN (Figure 1(a)). This GAN learns to draw rough shape and basic colors of the generated object conditioned on the given text description, and generate background regions from a random noise vector sampled from a prior distribution. The generated low resolution image is generally coarse and has a lot of defects, *e.g.*, object shape distortions and absence of object parts. It might not look real because some convincing details might be missing. On top of our Stage-I GAN, we stack Stage-II GAN to generate realistic high resolution images conditioned on the low resolution image and the corresponding text description (Figure 1(b)). Since Stage-I GAN generates rough shape and layout for both the object and background, Stage-II GAN only needs to focus on drawing details and rectifying defects in low resolution images. This task is much easier than directly drawing a high resolution image from scratch. By conditioning on the text again, Stage-II GAN learns to capture text information that are omitted by Stage-I GAN and draws more details for the object.

The main contribution of our paper is the design of the Stacked Generative Adversarial Networks (StackGAN), which can synthesize *photo-realistic* images from text descriptions. Compared to existing text-to-image generative models, our StackGAN generates images with more realistic details and achieves 28.47% and 20.30% improvements in terms of inception scores on CUB [32] and Oxford-102 [18] datasets, respectively. High-resolution 256×256 images can be generated by our proposed StackGAN, while state-of-the-art methods have difficulty generating images larger than 64×64 without additional spatial annotations.

2. Related Work

Generative image modeling is a fundamental problem in computer vision. Recently, there has been remarkable progress in this direction with the emergence of deep learning techniques. Dosovitskiy *et al.* [5] trained deconvolutional neural networks to generate 3D chairs, tables and cars. Others groups [36, 24] proposed deterministic neural networks as function approximators for image synthesis. Compared to these deterministic approaches, Variational Autoencoders (VAE) [12, 25] formulated the problem with probabilistic graphical models whose goal was to

maximize the lower bound of data likelihood. Gregor *et al.* [8] proposed the DRAW model, which applied the recurrent variational autoencoder together with the attention mechanism to generate realistic images for house numbers. Autoregressive models (*e.g.*, PixelRNN) [30] that utilized neural networks to model the conditional distribution of the pixel space have also generated appealing synthetic images. Recently, Generative Adversarial Networks (GAN) [7] have shown promising performance for generating sharper samples. Given the unstable training dynamics of GAN, several techniques [19, 26] have been proposed to stabilize the training process and generate compelling results on synthesizing faces, room interiors, and CIFAR-10 images. An energy-based GAN [37] has also been proposed for more stable training behavior.

Built upon these generative models, conditional image generation has also been studied. Most methods utilized simple conditioning variables such as attributes or class labels [35, 31, 2]. There is also work conditioned on images to generate images, such as photo editing [1, 38], domain transfer [29] and super-resolution [27, 14]. However, super-resolution methods [27, 14] can only add limited details to the low resolution images and can not correct large defects as our proposed StackGAN does. Recently, several methods have been developed to generate images from unstructured text. Mansimov *et al.* [16] built an AlignDRAW model by learning to estimate alignment between text and the generating canvas. Reed *et al.* [23] used conditional PixelCNN to generate images using the text descriptions and object location constraints. Using conditional GAN, Reed *et al.* [22] successfully generated plausible 64×64 images for birds and flowers from their text descriptions. Their follow-up work [20] was able to generate higher resolution images by controlling the object locations with additional annotations, such as bounding box or part keypoints of the object.

Besides using a single GAN for generating images, there is also work [33, 3] that utilized a series of GAN for image generation. Wang *et al.* [33] factorized the indoor scene generation process into structure generation and style generation with the proposed S^2 -GAN. The input images of the Style-GAN are samples drawn from the Structure-GAN. The two GANs focused on generating images of two different modalities. In contrast, the second stage of our StackGAN aims to complete object details and correct defects of Stage-I results based on text descriptions. Denton *et al.* [3] built a series of GAN within a Laplacian pyramid framework. At each level of the pyramid, a residual image was generated conditioned on the image of the previous stage and then added back to the input image to produce the input for the next stage. However, they only succeeded in generating 96×96 images, while our method utilizes a simpler architecture to generate 256×256 images with photo-realistic details and six times more pixels.

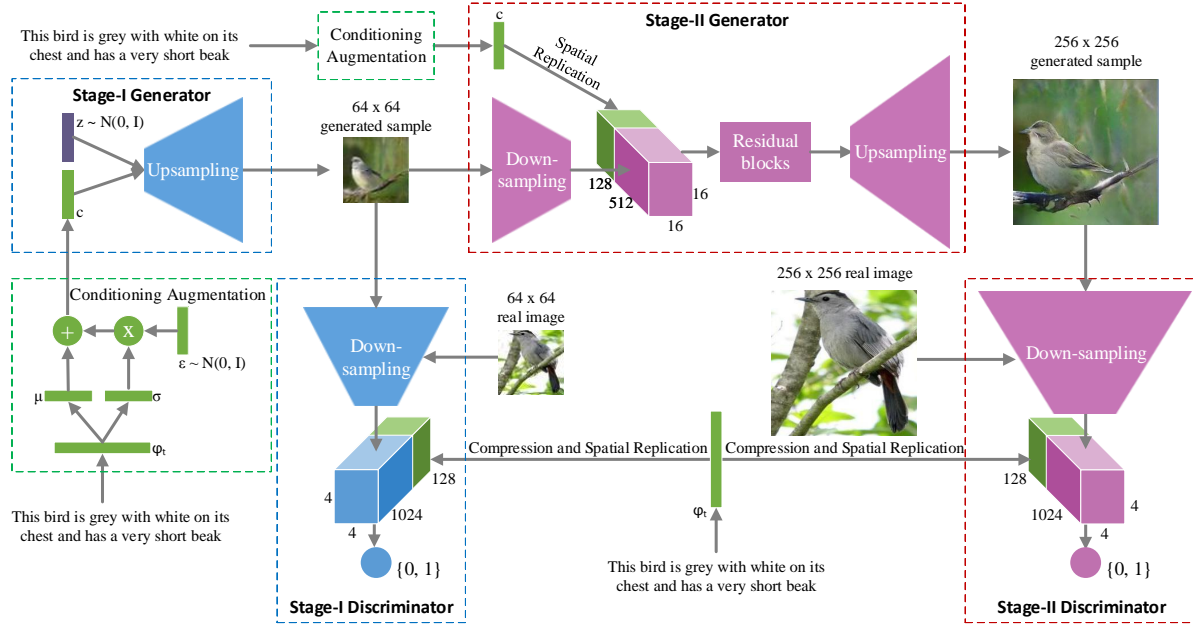


Figure 2. The architecture of the proposed StackGAN. The Stage-I generator draws a low resolution image by sketching rough shape and basic colors of the object from the given text and painting the background from a random noise vector. The Stage-II generator generates a high resolution image with photo-realistic details by conditioning on both the Stage-I result and the text again.

3. Stacked Generative Adversarial Networks

To generate high resolution images with photo-realistic details, we propose a simple yet effective stacked Generative Adversarial Networks (StackGAN). It decomposes the text-to-image generative process into two stages.

- **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 and completes details of the object by reading the text description again, producing a high resolution photo-realistic image.

3.1. Preliminaries

Generative Adversarial Networks (GAN) [7] are composed of two models that are alternatively trained to compete with each other. The generator G is optimized to reproduce the true data distribution p_{data} by generating images that are difficult for the discriminator D to differentiate from real images. Meanwhile, D is optimized to distinguish real images and synthetic images generated by G . Overall, the training procedure is similar to a two-player min-max game with the following objective function,

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

where x is a real image from the true data distribution p_{data} , and z is a noise vector sampled from distribution p_z (e.g., uniform or Gaussian distribution).

Conditional GAN [6, 17] is an extension of GAN where both the generator and discriminator receive additional conditioning variables c , yielding $G(z, c)$ and $D(x, c)$. This formulation allows G to generate images conditioned on variables c .

3.2. Stage-I GAN

Instead of directly generating a high resolution image conditioned on the text description, we simplify the task to first generate a low resolution image. Stage-I GAN is designed to generate such low resolution images and focus on drawing rough shape and correct colors for the object.

As shown in Figure 2, the conditioning text description t is first encoded by an encoder, yielding a text embedding φ_t . In previous works [22, 20], the text embedding is nonlinearly transformed to generate conditioning latent variables for the generator. However, latent space conditioned on text is usually high dimensional (> 100 dimensions). With limited amount of data, it usually causes discontinuity in the latent data manifold, which is not desirable for learning the generator.

To mitigate this problem, we introduce a conditioning augmentation technique to produce more conditioning variables for the generator. We randomly sample latent variables from an independent Gaussian distribution $\mathcal{N}(\mu(\varphi_t), \Sigma(\varphi_t))$, where the mean $\mu(\varphi_t)$ and diagonal co-

variance matrix $\Sigma(\varphi_t)$ are functions of the text embedding φ_t . The proposed formulation encourages robustness to small perturbations along the conditioning manifold, and thus yields more training pairs given a small number of image-text pairs. To further enforce the smoothness over the conditioning manifold and avoid overfitting [4, 13], we add the following regularization term to the objective of the generator during training,

$$D_{KL}(\mathcal{N}(\mu(\varphi_t), \Sigma(\varphi_t)) \parallel \mathcal{N}(0, I)), \quad (2)$$

which is the Kullback-Leibler divergence (KL divergence) between the standard Gaussian distribution and the conditioning Gaussian distribution.

Conditioned on Gaussian latent variables c_0 , Stage-I GAN trains the discriminator D_0 and the generator G_0 by alternatively maximizing \mathcal{L}_{D_0} in Eq. (3) and minimizing \mathcal{L}_{G_0} in Eq. (4),

$$\mathcal{L}_{D_0} = \mathbb{E}_{(I_0, t) \sim p_{data}} [\log D_0(I_0, \varphi_t)] + \mathbb{E}_{z \sim p_z, t \sim p_{data}} [\log(1 - D_0(G_0(z, c_0), \varphi_t))], \quad (3)$$

$$\mathcal{L}_{G_0} = \mathbb{E}_{z \sim p_z, t \sim p_{data}} [\log(1 - D_0(G_0(z, c_0), \varphi_t))] + \lambda D_{KL}(\mathcal{N}(\mu_0(\varphi_t), \Sigma_0(\varphi_t)) \parallel \mathcal{N}(0, I)), \quad (4)$$

where the real image I_0 and the text description t are from the true data distribution p_{data} . z is a noise vector randomly sampled from a given distribution p_z (e.g., Gaussian distribution used in this paper). λ is a regularization parameter that controls the balance between the two terms in Eq. (4). We use $\lambda = 1$ for all our experiments. φ_t is the text embedding, which is generated by a pre-trained encoder [21] in this paper. Gaussian conditioning variables c_0 are sampled from $\mathcal{N}(\mu_0(\varphi_t), \Sigma_0(\varphi_t))$ to reflect the text description. Using the reparameterization trick introduced in [12], both $\mu_0(\varphi_t)$ and $\Sigma_0(\varphi_t)$ are learned jointly with the rest of the network.

Model Architecture. For the generator, the text embedding φ_t is fed into a fully connected layer to generate μ_0 and σ_0 (σ_0 are the values in the diagonal of Σ_0) for Gaussian distribution $\mathcal{N}(\mu_0(\varphi_t), \Sigma_0(\varphi_t))$. Our N_g dimensional conditioning vector c_0 is computed by $c_0 = \mu_0 + \sigma_0 \odot \epsilon$ (where \odot is the element-wise multiplication, $\epsilon \sim \mathcal{N}(0, I)$). Then, c_0 is concatenated with a N_z dimensional noise vector to generate a $W_0 \times H_0$ image by a series of up-sampling blocks.

For the discriminator, the text embedding φ_t is first compressed to N_d dimensions using a fully-connected layer and then spatially replicated to form a $M_d \times M_d \times N_d$ tensor. Meanwhile, the image is fed through a series of down-sampling blocks until it has $M_d \times M_d$ spatial dimension. Then, the image filter map is concatenated along the channel dimension with the text tensor. The resulting tensor is further fed to a 1×1 convolutional layer to jointly learn

features across the image and the text. Finally, a fully-connected layer with one node is used to produce the decision score.

3.3. Stage-II GAN

Low resolution images generated by Stage-I GAN lack vivid object parts and might also contain shape distortions. In addition, some details in the text might be omitted in the first stage. This is important information needed to generate a photo-realistic image. Stage-II GAN is built upon Stage-I GAN to generate photo-realistic high resolution images. It conditions on low resolution images generated by the previous stage, and also the text embedding again to correct defects in Stage-I results and encourage the model to extract previously ignored information in the text to generate more photo-realistic details.

Conditioning on the low resolution sample s_0 and Gaussian latent variables c , the discriminator D and generator G in Stage-II GAN is trained by alternatively maximizing \mathcal{L}_D in Eq. (5) and minimizing \mathcal{L}_G in Eq. (6),

$$\mathcal{L}_D = \mathbb{E}_{(I, t) \sim p_{data}} [\log D(I, \varphi_t)] + \mathbb{E}_{s_0 \sim p_{G_0}, t \sim p_{data}} [\log(1 - D(G(s_0, c), \varphi_t))], \quad (5)$$

$$\mathcal{L}_G = \mathbb{E}_{s_0 \sim p_{G_0}, t \sim p_{data}} [\log(1 - D(G(s_0, c), \varphi_t))] + \lambda D_{KL}(\mathcal{N}(\mu(\varphi_t), \Sigma(\varphi_t)) \parallel \mathcal{N}(0, I)), \quad (6)$$

where $s_0 = G_0(z, c_0)$ is generated by Stage-I GAN. Different from the original GAN formulation, the random noise z is not used in this stage with the assumption that the randomness has already been preserved in s_0 . Gaussian conditioning variables c used in this stage and c_0 used in Stage-I GAN share the same pre-trained text encoder, generating the same text embedding φ_t . But, they utilize different fully connected layers for generating different means and standard deviations. In this way, Stage-II GAN learns to capture useful information in the text embedding that is omitted by Stage-I GAN.

Model Architecture. For the generator, similar to the previous stage, φ_t is used to generate our N_g dimensional Gaussian conditioning vector c , which is spatially replicated to form a $M_g \times M_g \times N_g$ tensor. Meanwhile, the sample s_0 generated by Stage-I GAN is fed into several down-sampling blocks until it has a spatial size of $M_g \times M_g$. Then, the image filter map and the text tensor are concatenated along the channel dimension. The resulting tensor is fed into several residual blocks [11, 9] to jointly encode the image and text features, and finally a series of up-sampling blocks are used to generate a $W \times H$ image.

For the discriminator, its structure is similar to that of Stage-I discriminator with only extra down-sampling blocks since the image size is larger in this stage. To explicitly enforce GAN to learn better alignment between the image and

the conditioning text, rather than using the naive discriminator, we adopt the matching-aware discriminator proposed by Reed *et al.* [22] for both stages. During training, the discriminator takes real images and their corresponding text descriptions as positive sample pairs, whereas negative sample pairs consist of two groups. The first is real images with mismatched text embeddings, while the second is synthetic images with conditioning text embeddings.

3.4. Implementation details

The up-sampling blocks consist of the nearest-neighbor upsampling followed by a 3×3 stride 1 convolution. Batch normalization [10] and ReLU activation are applied after every convolution except the last one. The residual blocks consist of 3×3 stride 1 convolutions, Batch normalization and ReLU. Two residual blocks are used in 128×128 StackGAN models while four are used in 256×256 models. The down-sampling blocks consist of 4×4 stride 2 convolutions, Batch normalization and LeakyReLU [15, 34], except that the first one does not have Batch normalization.

By default, $N_g = 128$, $N_z = 100$, $M_g = 16$, $M_d = 4$, $N_d = 128$, $W_0 = H_0 = 64$ and $W = H = 256$. For training, we first iteratively train D_0 and G_0 of Stage-I GAN for 600 epochs by fixing Stage-II GAN. Then we iteratively train D and G of Stage-II GAN for another 600 epochs by fixing Stage-I GAN. All networks are trained using ADAM solver with batch size 64 and an initial learning rate of 0.0002. The learning rate is decayed to 1/2 of its previous value every 100 epochs.

4. Experiments

To validate our method, we conduct extensive quantitative and qualitative evaluations on Caltech-UCSD Bird (CUB) [32] and Oxford-102 flower [18] datasets, which contain enough object appearance variations and are widely-used for text-to-image generation analysis. Two state-of-the-art methods on text-to-image synthesis, GAN-INT-CLS [22] and GAWWN [20], are compared. Results by the two compared methods are generated using the code and models released by their authors. In addition, we design several baseline models to investigate the design and important components of our proposed StackGAN. For the first baseline, we utilize only Stage-I GAN of our StackGAN for generating 64×64 images to investigate whether the stack structure is beneficial. Then we modify our StackGAN to generate 128×128 images to investigate whether larger images by our method results in clearer images. We also investigate whether inputting text at both stages and adding the proposed conditioning augmentation are helpful for generating images of better quality.

4.1. Datasets and evaluation metrics

CUB [32] contains 200 bird species with 11,788 images. Since 80% of birds in this dataset have object-image size ratios of less than 0.5 [32], as a pre-processing step, we crop all images to ensure that bounding boxes of birds have greater-than-0.75 object-image size ratios. Oxford-102 [18] contains 8,189 images of flowers from 102 different categories. For this dataset, center cropping is not needed because most flowers are relatively large in the images. For every image in CUB and Oxford-102 datasets, 10 captions are provided by [21]. Following the experimental setup in [22], we split each dataset into class-disjoint training and test sets. CUB is split to 150 train classes and 50 test classes; Oxford-102 is split to 82 train and 20 test classes. During training, we randomly crop and flip images, and generate the corresponding text embedding as the average of four randomly selected descriptions corresponding to the image.

Evaluation metrics. It is difficult to evaluate the performance of generative models (*e.g.*, GAN). Asking human annotators to determine the visual quality of samples is most intuitive and reliable. We also choose a recently proposed numerical assessment approach “inception score” [26] for quantitative evaluation,

$$I = \exp(\mathbb{E}_{\mathbf{x}} D_{KL}(p(y|\mathbf{x}) || p(y))), \quad (7)$$

where \mathbf{x} denotes one generated sample, and y is the label predicted by the Inception model [28]. The intuition behind this metric is that good models should generate diverse but meaningful images. Therefore, the KL divergence between the marginal distribution $p(y)$ and the conditional distribution $p(y|\mathbf{x})$ should be large. In our experiments, we fine-tune two Inception models on CUB and Oxford-102 datasets, respectively. As suggested in [26], we evaluate this metric on a large number of samples (*i.e.*, 30k) for each model.

Although the inception score has shown to well correlate with human perception on visual quality of samples [26], it cannot reflect whether the generated images are well conditioned on the given text descriptions. Therefore, we also conduct human evaluation. We randomly select 50 text descriptions for each class and generate 5 images for each sentence. Given the same text descriptions, 10 users (not including any of the authors) are asked to rank the results by different methods. The average human ranks are calculated to evaluate all compared methods.

4.2. Quantitative and qualitative results

We compare our results with the state-of-the-art text-to-image methods [20, 22] on CUB and Oxford-102 datasets. The inception scores and average human ranks for our proposed StackGAN and compared methods are reported in Ta-



Figure 3. Example results by our proposed StackGAN, GAWWN [20], and GAN-INT-CLS [22] conditioned on text descriptions from CUB test set. GAWWN and GAN-INT-CLS generate 16 images for each text description, respectively. We select the best one for each of them to compare with our StackGAN.

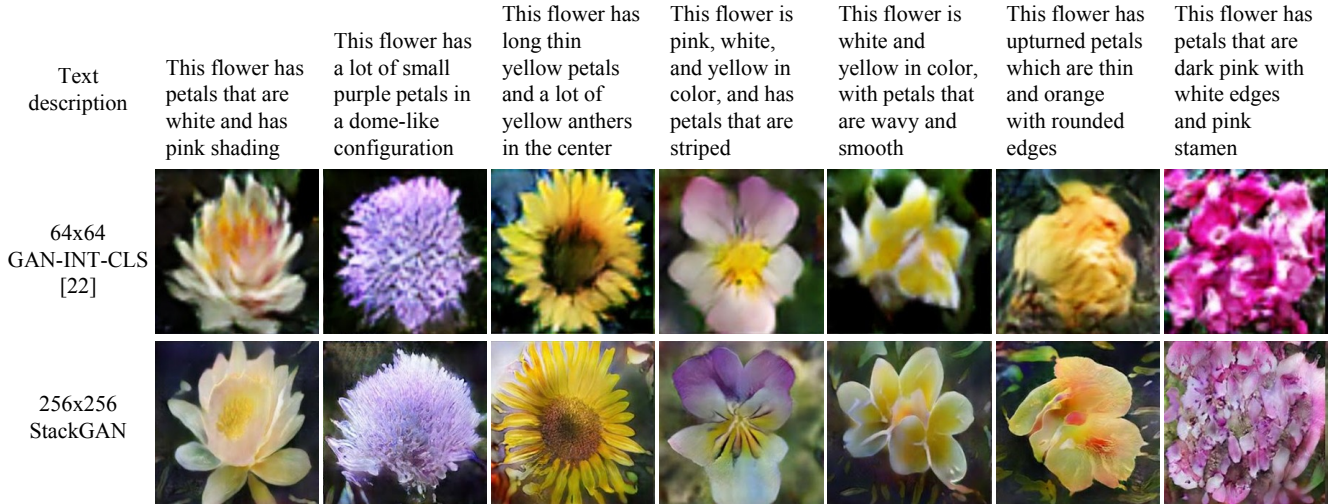


Figure 4. Example results by our proposed StackGAN and GAN-INT-CLS [22] conditioned on text descriptions from Oxford-102 test set.

ble 1. Representative examples generated by text descriptions by different methods are shown in Figures 3 and 4.

Our StackGAN achieves the best inception score and average human rank on both datasets. Compared with GAN-INT-CLS [22], StackGAN achieves 28.47% improvement in terms of inception score on CUB dataset (from 2.88 to 3.70), and 20.30% improvement on Oxford-102 (from 2.66 to 3.20). The better average human rank of our StackGAN also indicates our proposed method is able to generate more realistic samples conditioned on text descriptions.

As shown in Figure 3, the 64×64 samples generated by GAN-INT-CLS can only reflect the general shape and color of the birds. Their results lack vivid parts (*e.g.*, beak and legs) and convincing details in most cases, which make them neither realistic enough nor have sufficiently high resolution. By using additional conditioning variables on location constraints, GAWWN [20] obtains a better inception score on CUB dataset, which is still slightly lower than ours. It generates higher resolution images with more details than GAN-INT-CLS, as shown in Figure 3. However, as men-



Figure 5. Samples generated by our StackGAN from unseen texts in CUB test set. Each column lists the text description, images generated from the text by Stage-I and Stage-II of StackGAN. To illustrate the capability of StackGAN on generating diverse images, those examples are selected from different classes. State-II GAN is able to correct defects and completing details of the Stage-I results for both foreground objects and background.

Method	Inception scores		Human rank	
	CUB	Oxford-102	CUB	Oxford-102
GAN-INT-CLS [22]	$2.88 \pm .04$	$2.66 \pm .03$	$2.81 \pm .03$	$1.87 \pm .03$
GAWWN [20]	$3.62 \pm .07$	/	$1.99 \pm .04$	/
Our StackGAN	$3.70 \pm .04$	$3.20 \pm .01$	$1.37 \pm .02$	$1.13 \pm .03$

Table 1. Inception scores and average human ranks of our proposed StackGAN, GAWWN [20], and GAN-INT-CLS [22] on CUB and Oxford-102 datasets. The inception score highly correlates with human perception in terms of image quality. Higher inception scores mean better image quality. The smaller average human rank indicates better results from human’s perspective.



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 L_2 distances between features are calculated for nearest-neighbor retrieval.

tioned by its authors, GAWWN fails to generate any plausible images if only conditioned on text descriptions [20]. In comparison, our StackGAN can generate 256×256 photo-realistic images from only text descriptions.

Figure 5 illustrates some examples of the Stage-I and

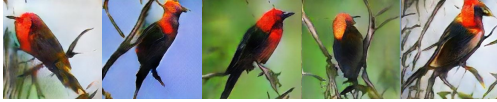
Stage-II images generated by our StackGAN. As shown in the first row of Figure 5, in most cases, Stage-I GAN is able to draw rough shapes and colors of objects given text descriptions. However, Stage-I images are usually blurry with various defects and missing details, especially for foreground objects. As shown in the second row, Stage-II GAN generates $4 \times$ higher resolution images with more convincing details to better reflect corresponding text descriptions. For cases where Stage-I GAN has generated plausible shapes and colors, Stage-II GAN completes the details. For instance, in the 1st column of Figure 5, with a satisfactory Stage-I result, Stage-II GAN focuses on drawing the short beak and white color described in the text as well as details for the tail and legs. In all other examples, different degrees of details are added to Stage-II images. In many other cases, Stage-II GAN is able to correct the defects of Stage-I results by processing the text description again. For example, while the Stage-I image in the 5th column has a blue crown rather than the reddish brown crown described in the text, the defect is corrected by Stage-II GAN. In some extreme cases (e.g., the 7th column of Figure 5), even when Stage-I GAN fails to draw a plausible shape, Stage-II GAN is able to generate reasonable objects. We also observe that StackGAN has the ability to transfer background from Stage-I images and fine-tune them to be more realistic with higher resolution at Stage-II.

Importantly, the StackGAN does not achieve good results by simply memorizing training samples but by capturing the complex underlying language-image relations. We extract visual features from our generated images and all training images by the Stage-II discriminator D of our StackGAN. For each generated image, its nearest neighbors from the training set can be retrieved. By visually inspect-

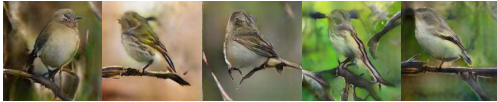
This small blue bird has a short pointy beak and brown on its wings



This bird is completely red with black wings and pointy beak



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



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



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.

ing the retrieved training images, we can conclude that the generated images have some similar characteristics with the retrieved training images but are essentially different. Three examples of retrieving nearest training images with the generated images on CUB and Oxford-102 datasets are shown in Figure 6.

4.3. Component analysis

In this subsection, we analyze different components of our proposed StackGAN on CUB dataset with our baseline models, which further validate the design of StackGAN.

The design of StackGAN. The inception scores for different baselines of our StackGAN are reported in Table 2. As shown in the first row of Table 2, if only Stage-I GAN is utilized, the inception score decreases significantly from 3.70 to 2.66. Such performance drop can also be well illustrated by Stage-I results in Figure 5, which demonstrates the necessity of our two-stage structure. By decreasing the output resolution from 256×256 to 128×128 , the inception score decreases from 3.70 to 3.5. Note that all images are scaled to 299×299 by default before calculating the inception score. Therefore, if our StackGAN just increases the image size without adding more information, the inception score will stay the same for different resolution samples. Therefore, the decrease of inception score by 128×128 StackGAN demonstrates that our 256×256 StackGAN does add more details into the larger images. For the 128×128 StackGAN, if the text is only input at the Stage-I GAN (denoted as “no Text twice”), the inception score decreases from 3.35 to 3.13, which shows that Stage-II GAN needs to process text descriptions again in order to better refine Stage-I results.

Conditioning augmentation. We also investigate the effectiveness of our proposed conditioning augmentation. By removing it from StackGAN 128×128 (denoted as “no CA” in Table 2), the inception score decreases from 3.35

Method	CA	Text twice	Inception score
64×64 Stage-I GAN	no	/	$2.66 \pm .03$
128×128 StackGAN	yes	no	$3.13 \pm .03$
	no	yes	$3.20 \pm .03$
256×256 StackGAN	yes	yes	$3.35 \pm .02$
	yes	yes	$3.70 \pm .04$

Table 2. Inception scores calculated with 30,000 samples generated by different baseline models of our StackGAN. Higher inception scores mean better image quality.

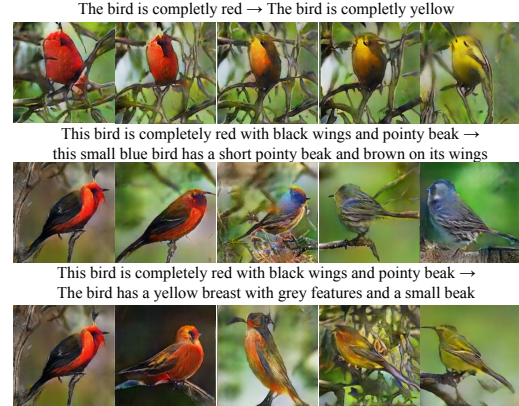


Figure 8. (Left to right) Images generated by interpolating two sentence embeddings. Gradual appearance changes from the first sentence’s meaning to that of the second sentence can be observed. The noise vector z is fixed to be zeros for each row.

to 3.20. Moreover, Figure 7 shows that we can generate birds with different poses and viewpoints from the same text embedding with the proposed conditioning augmentation. While all inputs of StackGAN are fixed, conditioning augmentation samples different conditioning variables from the Gaussian distribution controlled by the text embedding. As a result, the conditioning augmentation encourages robustness to small perturbations along the latent manifold, yielding more diverse samples.

Sentence embedding interpolation. To demonstrate that our StackGAN learns a smooth latent data manifold, we use it to generate images from linearly interpolated sentence embeddings, as shown in Figure 8. We fix the noise vector z so the generated image is inferred from the given text description only. Images in the first row are generated by simple sentences made up by us. Those sentences contain only simple color descriptions. The results show that the generated images from interpolated embeddings can accurately reflect color changes and generate plausible bird shapes. We do observe that StackGAN prefers to generate a simple image without parts and details if the text descriptions are too simple. The last two rows illustrate samples generated from more complex sentences. Those sentences contain more details on bird appearances. The generated images from interpolated embeddings succeed in reflecting changes of the sentence meanings. For example, images in the second row change their primary color from red to blue, and change the wing color from black to brown.

5. Conclusions

In this paper, we propose stacked Generative Adversarial Networks (StackGAN) for synthesizing *photo-realistic* images. The proposed method decomposes the synthesis process into two more manageable stages. Stage-I GAN draws the object following basic colors and shape constraints from given text descriptions, and Stage-II GAN corrects the defects in Stage-I results and adds more photo-realistic details. Extensive quantitative and qualitative results demonstrate the effectiveness of our proposed method. Compared to existing text-to-image generative models, our method generates higher resolution images (*e.g.*, 256×256) with more photo-realistic details.

References

- [1] A. Brock, T. Lim, J. M. Ritchie, and N. Weston. Neural photo editing with introspective adversarial networks. *arXiv:1609.07093*, 2016. 2
- [2] X. Chen, Y. Duan, R. Houthoofd, J. Schulman, I. Sutskever, and P. Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *NIPS*, 2016. 2
- [3] E. L. Denton, S. Chintala, A. Szlam, and R. Fergus. Deep generative image models using a laplacian pyramid of adversarial networks. In *NIPS*, 2015. 1, 2
- [4] C. Doersch. Tutorial on variational autoencoders. *arXiv:1606.05908*, 2016. 4
- [5] A. Dosovitskiy, J. T. Springenberg, and T. Brox. Learning to generate chairs with convolutional neural networks. In *CVPR*, 2015. 2
- [6] J. Gauthier. Conditional generative adversarial networks for convolutional face generation. *Technical report*, 2015. 3
- [7] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014. 1, 2, 3
- [8] K. Gregor, I. Danihelka, A. Graves, D. J. Rezende, and D. Wierstra. DRAW: A recurrent neural network for image generation. In *ICML*, 2015. 2
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 4
- [10] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 5
- [11] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 4
- [12] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 2, 4
- [13] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther. Autoencoding beyond pixels using a learned similarity metric. In *ICML*, 2016. 4
- [14] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi. Photo-realistic single image super-resolution using a generative adversarial network. *arXiv:1609.04802*, 2016. 2
- [15] A. L. Maas, A. Y. Hannun, and A. Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *ICML*, 2013. 5
- [16] E. Mansimov, E. Parisotto, L. J. Ba, and R. Salakhutdinov. Generating images from captions with attention. In *ICLR*, 2016. 2
- [17] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv:1411.1784*, 2014. 3
- [18] M.-E. Nilsback and A. Zisserman. Automated flower classification over a large number of classes. In *ICCVGIP*, 2008. 1, 2, 5
- [19] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016. 1, 2
- [20] S. Reed, Z. Akata, S. Mohan, S. Tenka, B. Schiele, and H. Lee. Learning what and where to draw. In *NIPS*, 2016. 1, 2, 3, 5, 6, 7
- [21] S. Reed, Z. Akata, B. Schiele, and H. Lee. Learning deep representations of fine-grained visual descriptions. In *CVPR*, 2016. 4, 5
- [22] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee. Generative adversarial text-to-image synthesis. In *ICML*, 2016. 1, 2, 3, 5, 6, 7
- [23] S. Reed, A. van den Oord, N. Kalchbrenner, V. Bapst, M. Botvinick, and N. de Freitas. Generating interpretable images with controllable structure. *Technical report*, 2016. 2
- [24] S. E. Reed, Y. Zhang, Y. Zhang, and H. Lee. Deep visual analogy-making. In *NIPS*, 2015. 2
- [25] D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *ICML*, 2014. 2
- [26] T. Salimans, I. J. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. In *NIPS*, 2016. 2, 5
- [27] C. K. Snderby, J. Caballero, L. Theis, W. Shi, and F. Huszar. Amortised map inference for image super-resolution. *arXiv:1610.04490*, 2016. 2
- [28] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016. 5
- [29] Y. Taigman, A. Polyak, and L. Wolf. Unsupervised cross-domain image generation. *arXiv:1611.02200*, 2016. 2
- [30] A. van den Oord, N. Kalchbrenner, and K. Kavukcuoglu. Pixel recurrent neural networks. In *ICML*, 2016. 2
- [31] A. van den Oord, N. Kalchbrenner, O. Vinyals, L. Espeholt, A. Graves, and K. Kavukcuoglu. Conditional image generation with pixelcnn decoders. In *NIPS*, 2016. 2
- [32] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011. 1, 2, 5
- [33] X. Wang and A. Gupta. Generative image modeling using style and structure adversarial networks. In *ECCV*, 2016. 2
- [34] B. Xu, N. Wang, T. Chen, and M. Li. Empirical evaluation of rectified activations in convolutional network. In *ICML Workshop*, 2015. 5

- [35] X. Yan, J. Yang, K. Sohn, and H. Lee. Attribute2image: Conditional image generation from visual attributes. In *ECCV*, 2016. [2](#)
- [36] J. Yang, S. E. Reed, M. Yang, and H. Lee. Weakly-supervised disentangling with recurrent transformations for 3d view synthesis. In *NIPS*, 2015. [2](#)
- [37] J. Zhao, M. Mathieu, and Y. LeCun. Energy-based generative adversarial network. *arXiv:1609.03126*, 2016. [2](#)
- [38] J. Zhu, P. Krähenbühl, E. Shechtman, and A. A. Efros. Generative visual manipulation on the natural image manifold. In *ECCV*, 2016. [2](#)

Supplementary Materials

Additional Results on CUB Dataset

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



This smaller brown bird has white stripes on the coverts, wingbars and secondaries



A cardinal looking bird, but fatter with gray wings, an orange head, and black eyerings



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



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

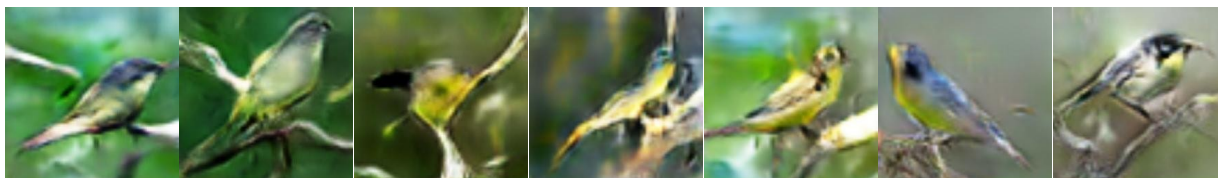


A small bird with orange crown and pointy bill and the bird has mixed color breast and side



This bird is light brown, gray, and yellow in color, with a light colored beak

Stage-I
images

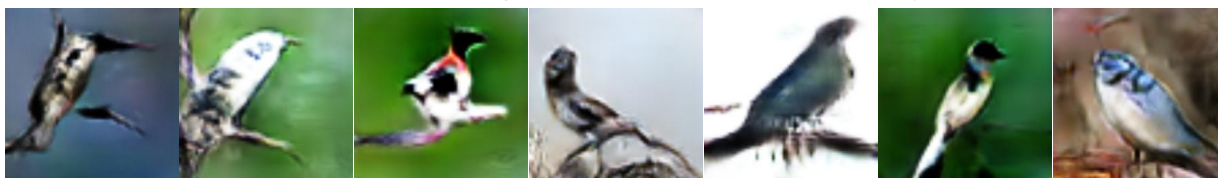


Stage-II
images



This bird has wings that are black and has a white belly

Stage-I
images

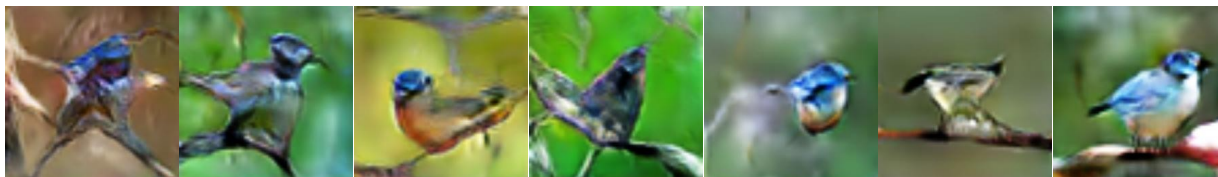


Stage-II
images



This bird has wings that are blue and has a white belly

Stage-I
images



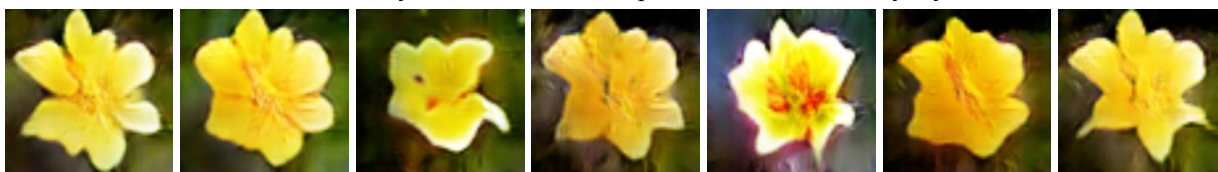
Stage-II
images



Additional Results on Oxford-102 Dataset

This flower is yellow in color, with petals that are vertically layered

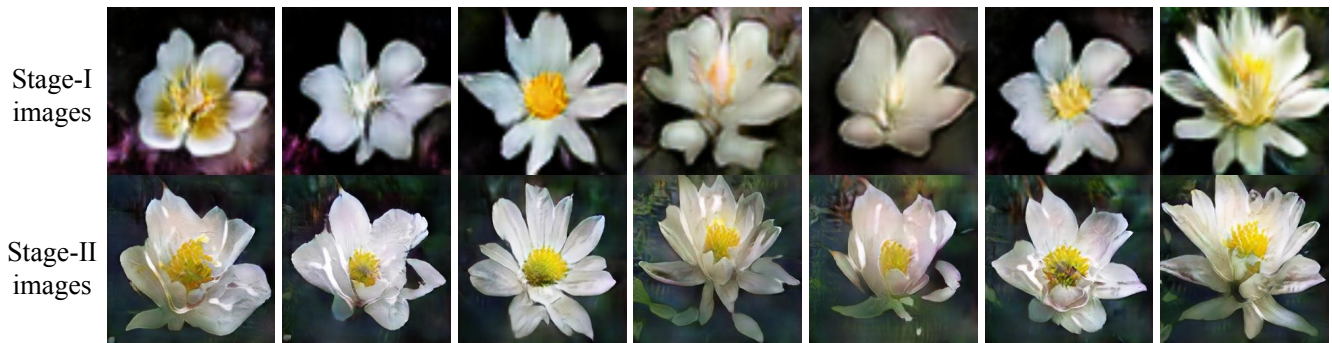
Stage-I
images



Stage-II
images



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



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



This flower is white, pink, and yellow in color, and has petals that are multi colored



This flower has petals that are yellow with shades of orange





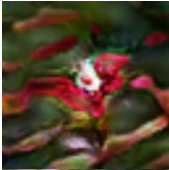




Failure Cases

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

CUB failure cases:

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
Stage-I images							
Stage-II images							

Oxford-102 failure cases:

Text description	The petals of this flower are white with a large stigma	A unique yellow flower with no visible pistils protruding from the center	This flower is pink and yellow in color, with petals that are oddly shaped	This is a light colored flower with many different petals on a green stem	This flower is yellow and green in color, with petals that are ruffled	The flower have large petals that are pink with yellow on some of the petals	A flower that has white petals with some tones of yellow and green filaments
Stage-I images							
Stage-II images	