

the flower has petals that are bright pinkish purple with white stigma



this white and yellow flower have thin white petals and a round yellow stamen



this small bird has a pink breast and crown, and black primaries and secondaries.



this magnificent fellow is almost all black with a red crest, and white cheek patch.



Generative Adversarial Text to Image Synthesis^[*]

**Scott Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran
Bernt Schiele, Honglak Lee**

REEDSCOT¹, AKATA², XCYAN¹, LLAJAN¹
SCHIELE², HONGLAK¹

¹ University of Michigan, Ann Arbor, MI, USA (UMICH.EDU)

² Max Planck Institute for Informatics, Saarbrücken, Germany (MPI-INF.MPG.DE)

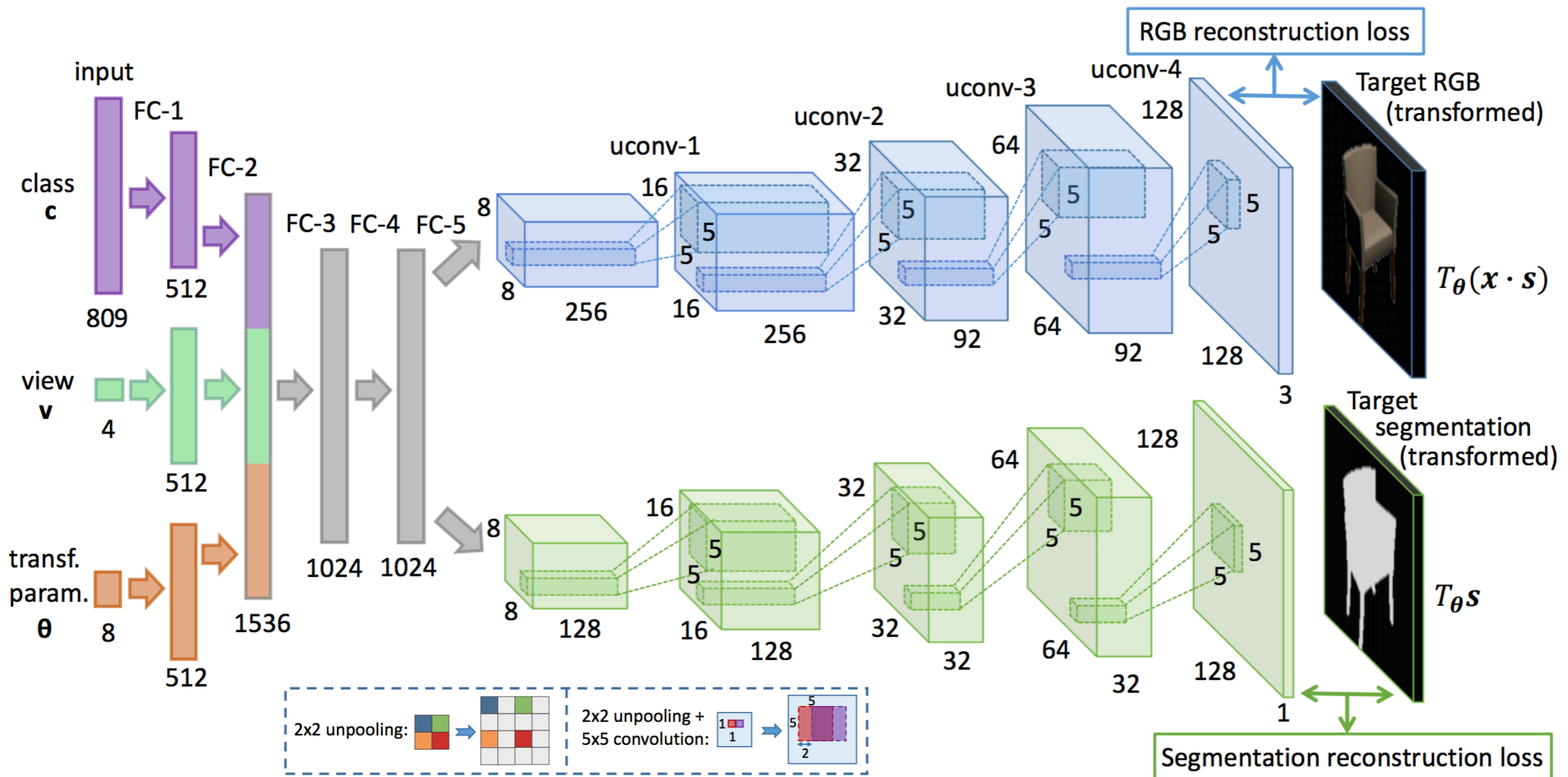
Presented By: Ezgi Ekiz

[*] Reed, Scott, et al. "Generative adversarial text to image synthesis." Proceedings of The 33rd International Conference on Machine Learning. Vol. 3. 2016.

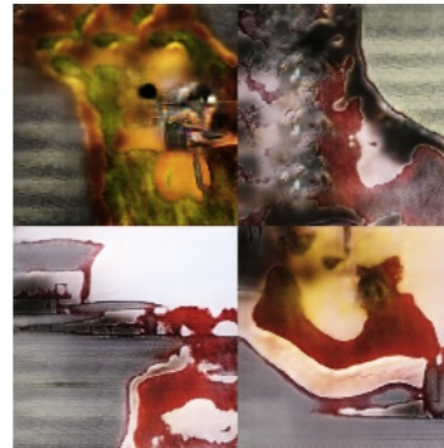
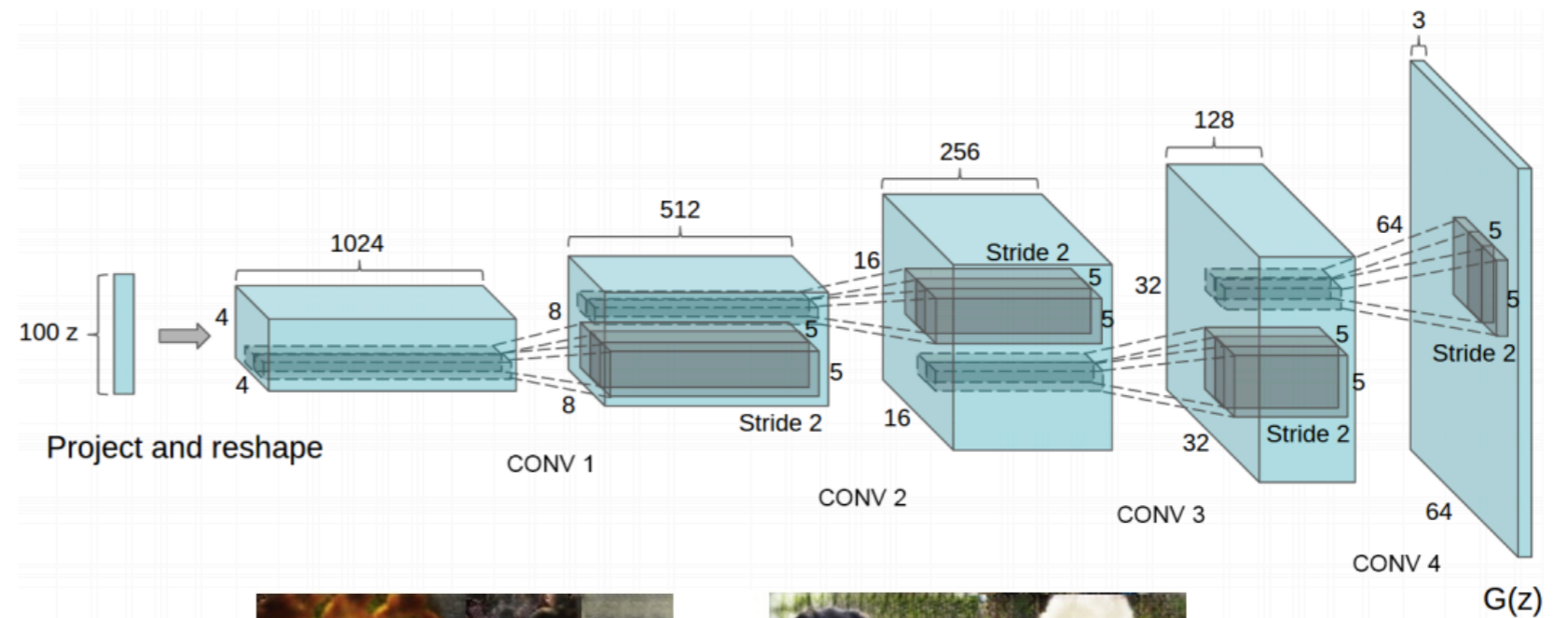
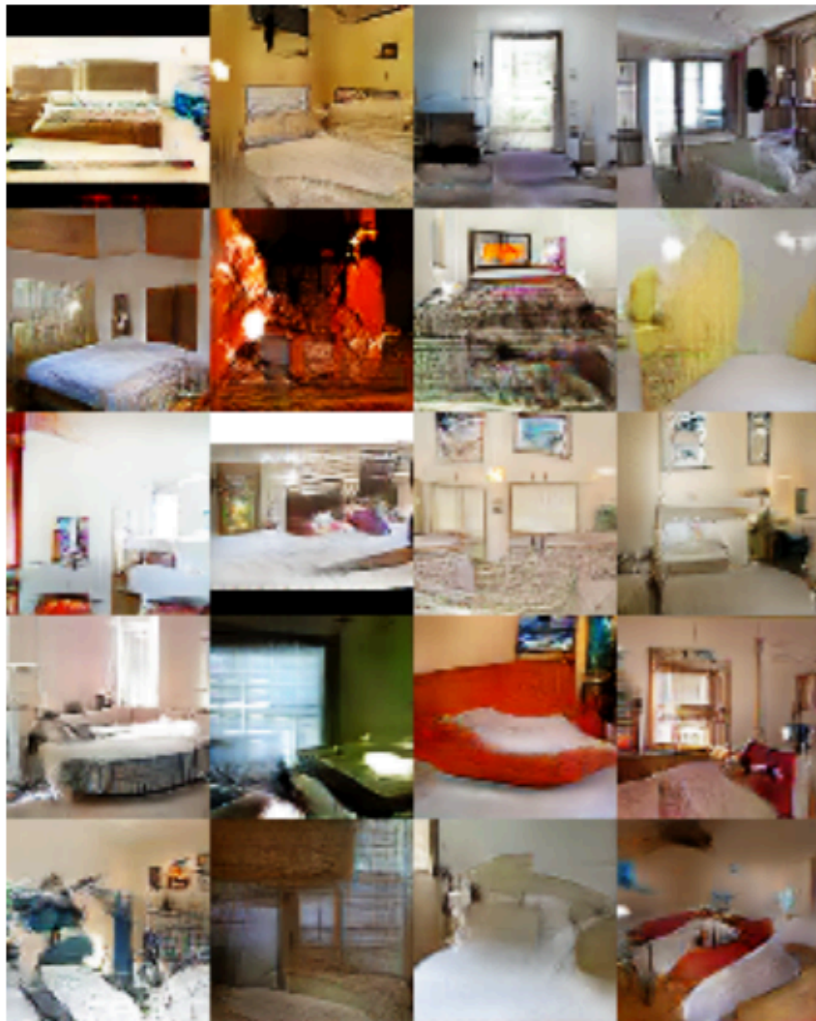
Outline

- The goal is to synthesize images that are mistakable for real from textual description. The method is built upon:
 - Text encoding that captures important visual details
 - Generative Adversarial Networks (GAN) and GAN-CLS
 - Manifold interpolation
 - Style transfer

Similar Work



Similar Work



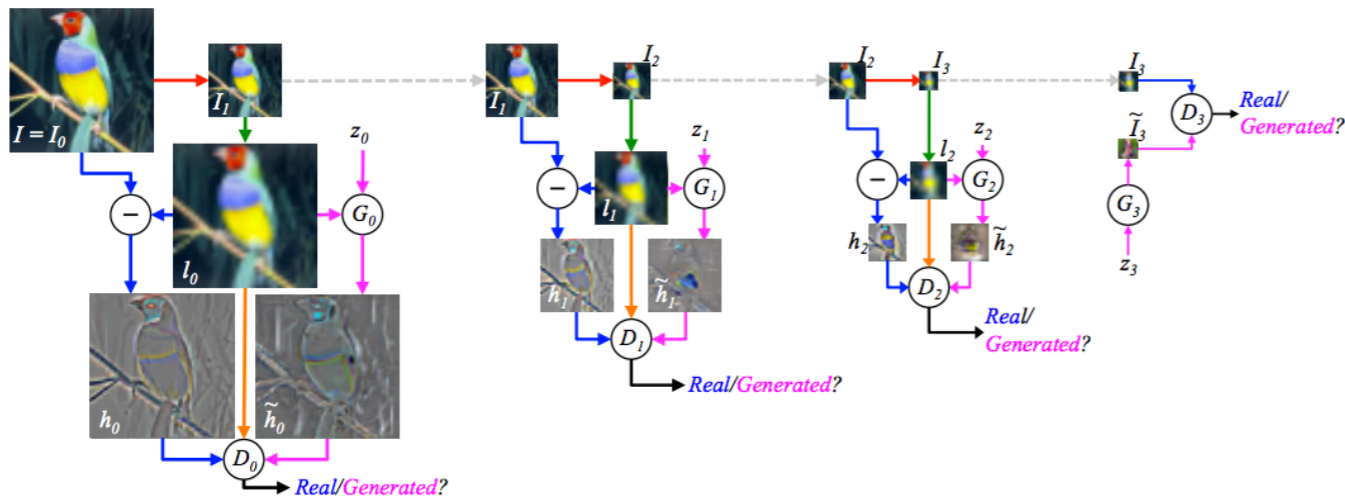
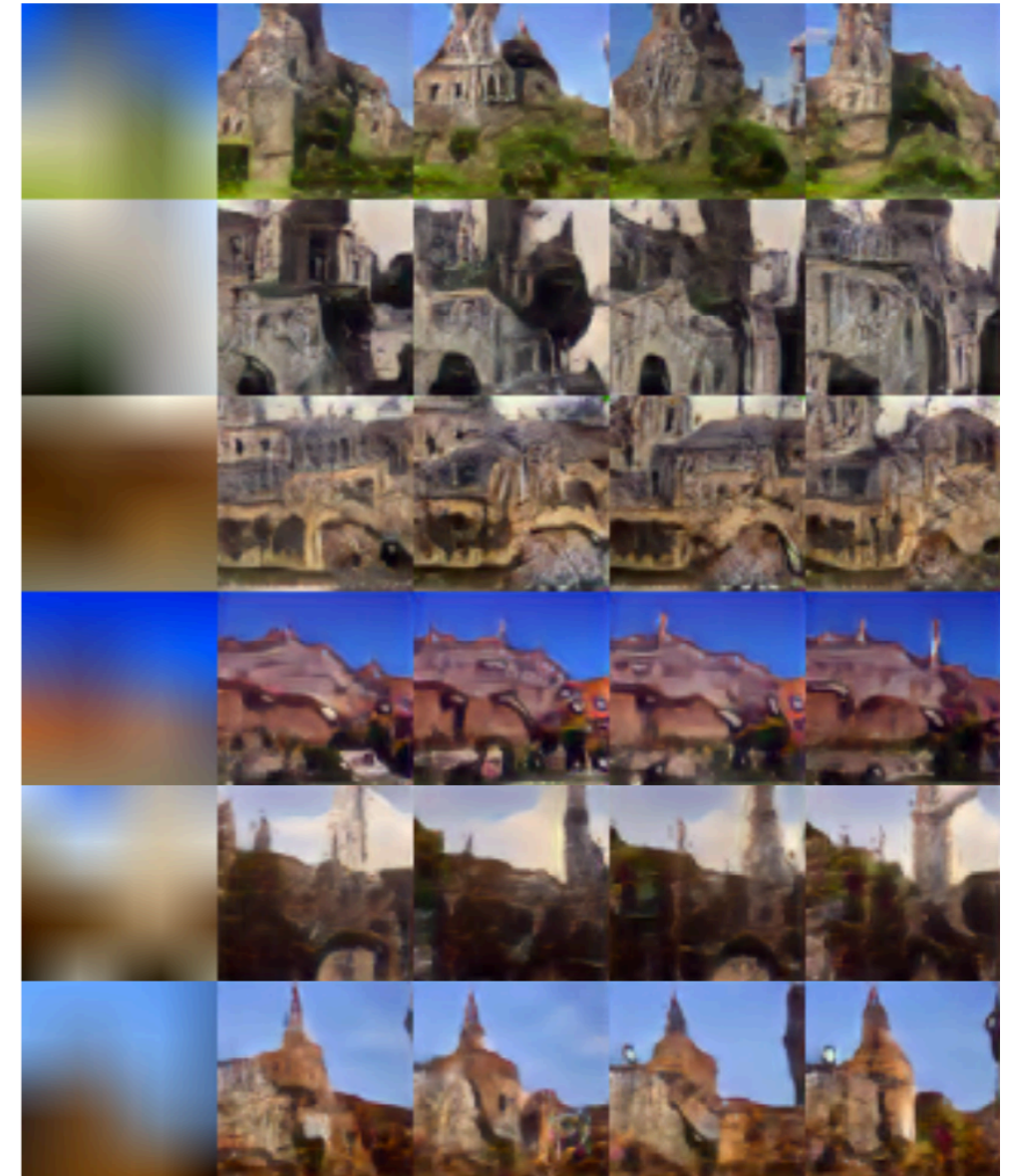
DCGAN



Improved DCGAN

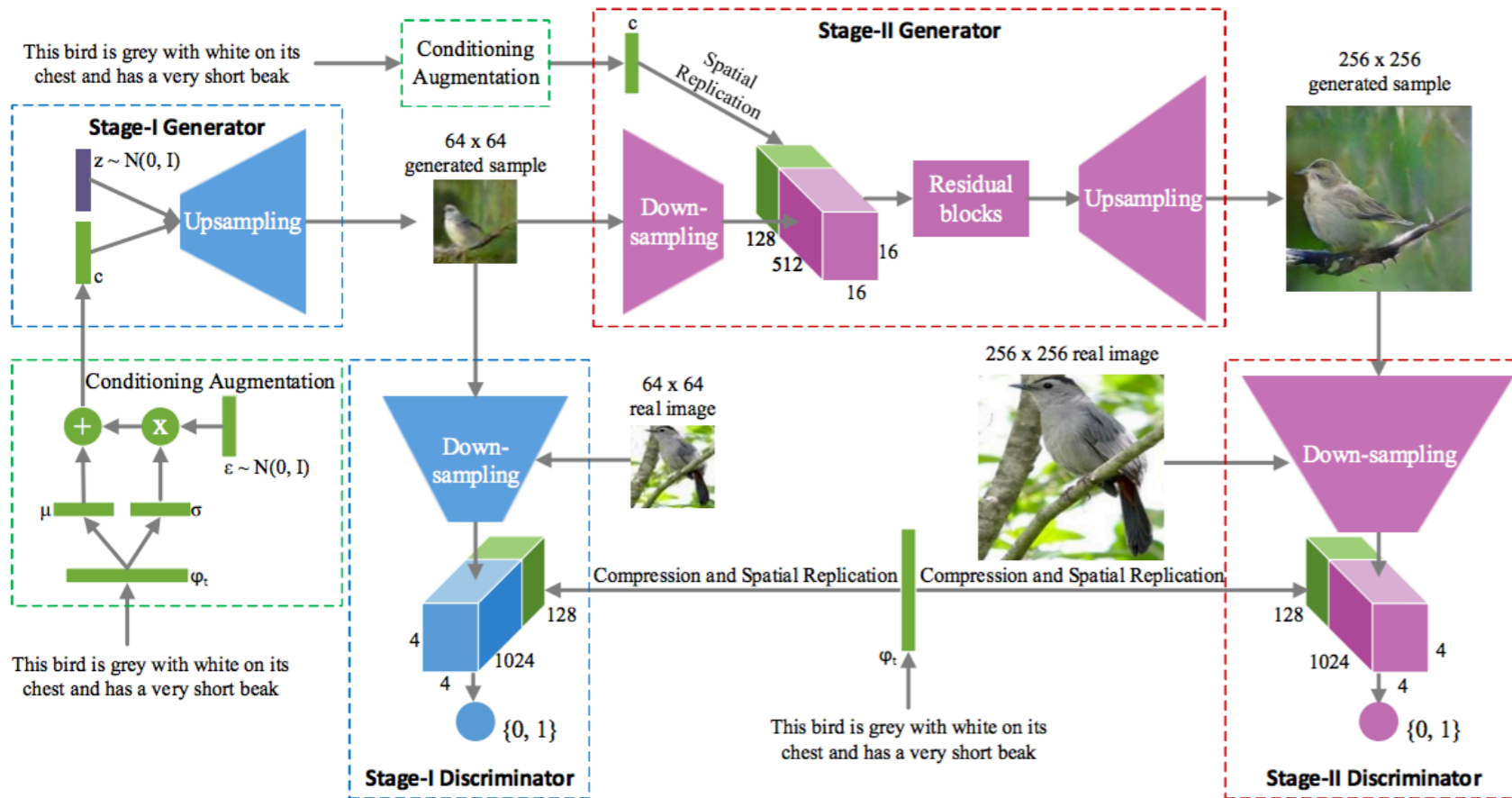
- Radford, Alec, Luke Metz, and Soumith Chintala. "Unsupervised representation learning with deep convolutional generative adversarial networks." arXiv preprint arXiv:1511.06434 (2015).
- Salimans, Tim, et al. "Improved techniques for training gans." Advances in Neural Information Processing Systems. 2016.

Similar Work



Denton, Emily L., Soumith Chintala, and Rob Fergus. "Deep Generative Image Models using a Laplacian Pyramid of Adversarial Networks." Advances in neural information processing systems. 2015.

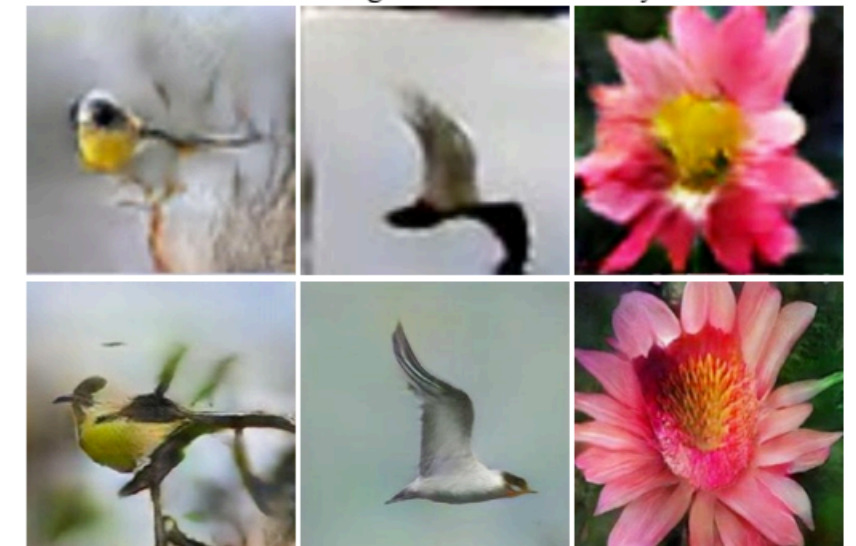
Similar Work



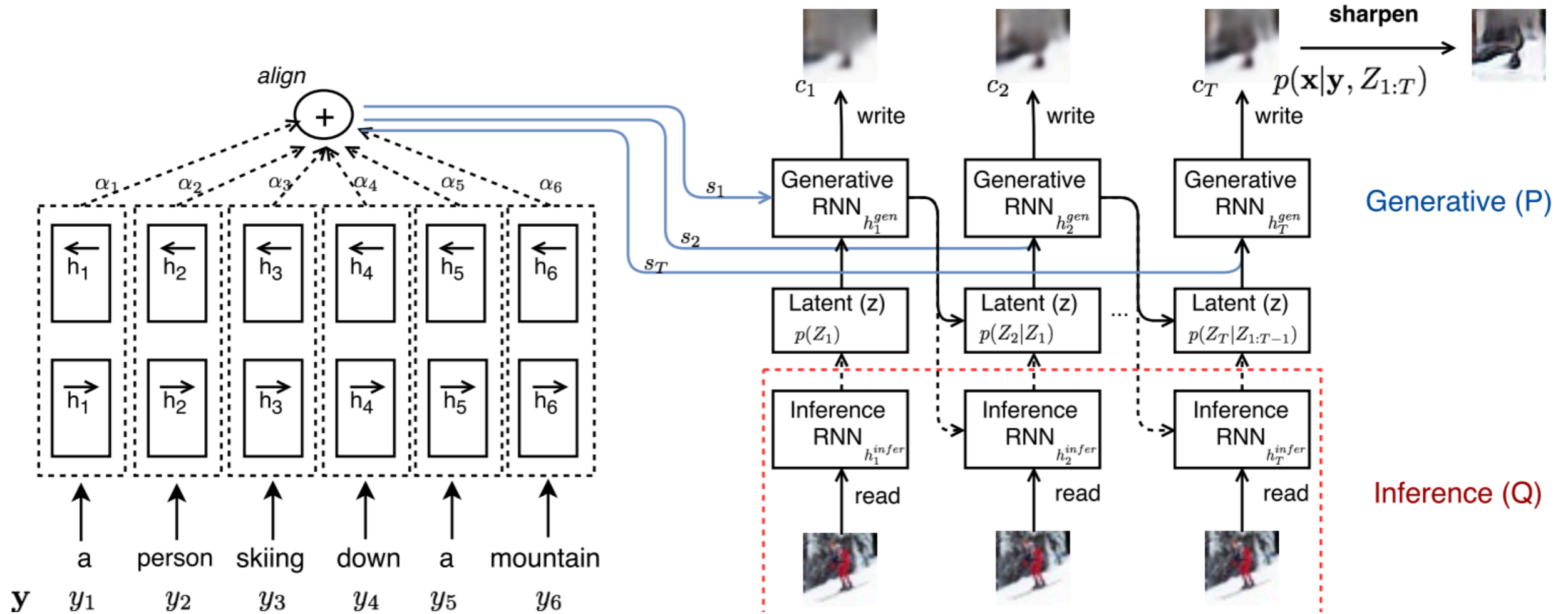
This bird has a yellow belly and tarsus, grey back, wings, and brown throat, nape with a black face

This bird is white with some black on its head and wings, and has a long orange beak

This flower has overlapping pink pointed petals surrounding a ring of short yellow filaments



Similar Work



A stop sign is flying in blue skies.



A herd of elephants flying in the blue skies.



A toilet seat sits open in the grass field.



A person skiing on sand clad vast desert.

Text Feature Representation

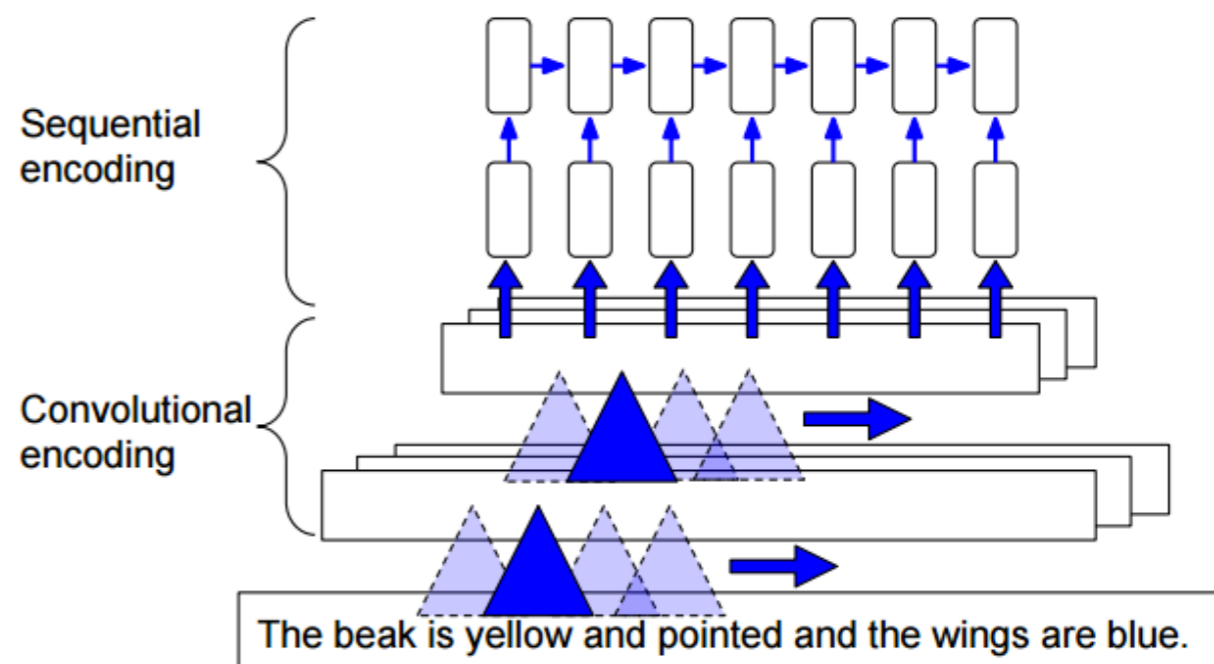
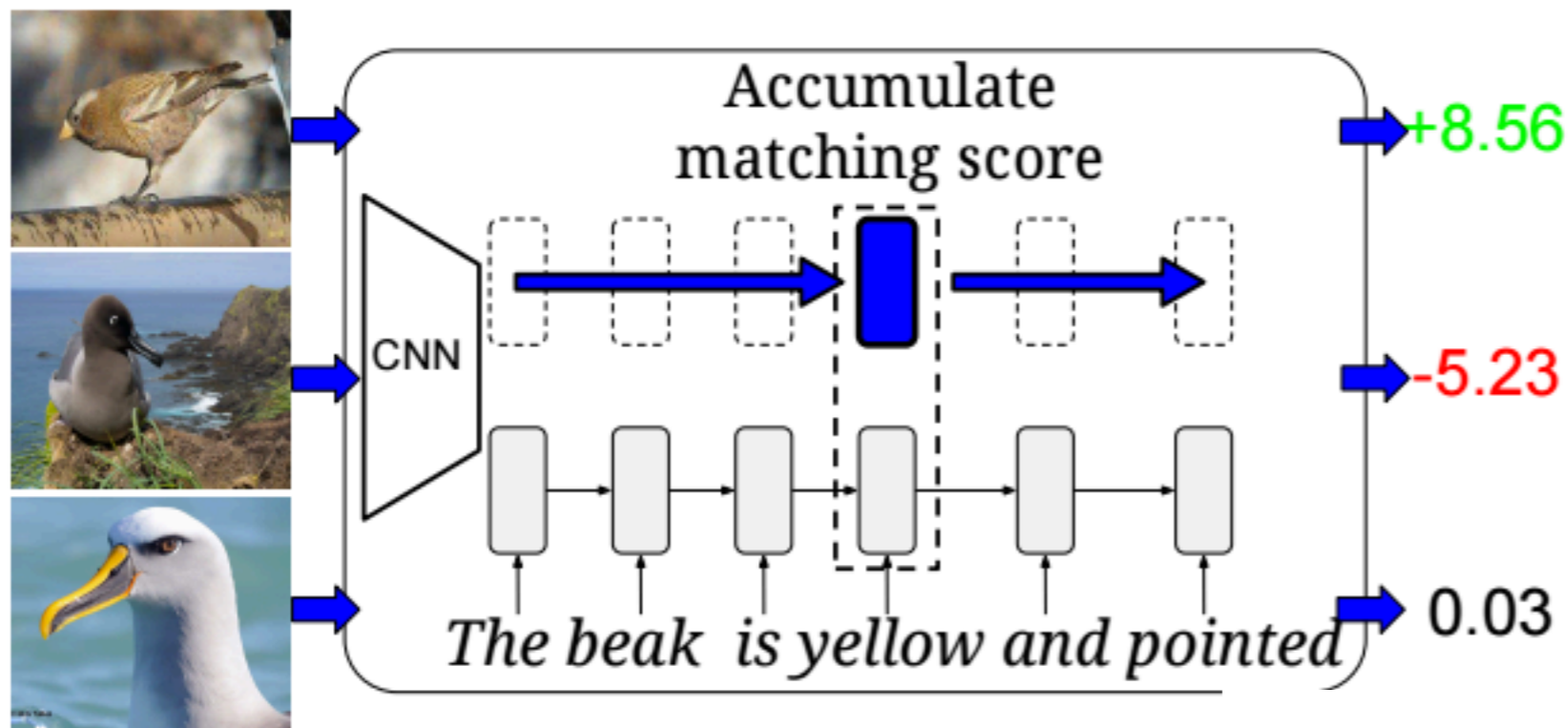
- The representation should capture important visual details
- Word/character based convolutional recurrent network is used

$$\frac{1}{N} \sum_{n=1}^N \Delta(y_n, f_v(v_n)) + \Delta(y_n, f_t(t_n))$$

$$f_v(v) = \arg \max_{y \in \mathcal{Y}} \mathbb{E}_{t \sim \mathcal{T}(y)} [\phi(v)^T \varphi(t)]$$

$$f_t(t) = \arg \max_{y \in \mathcal{Y}} \mathbb{E}_{v \sim \mathcal{V}(y)} [\phi(v)^T \varphi(t)]$$

Text Feature Representation



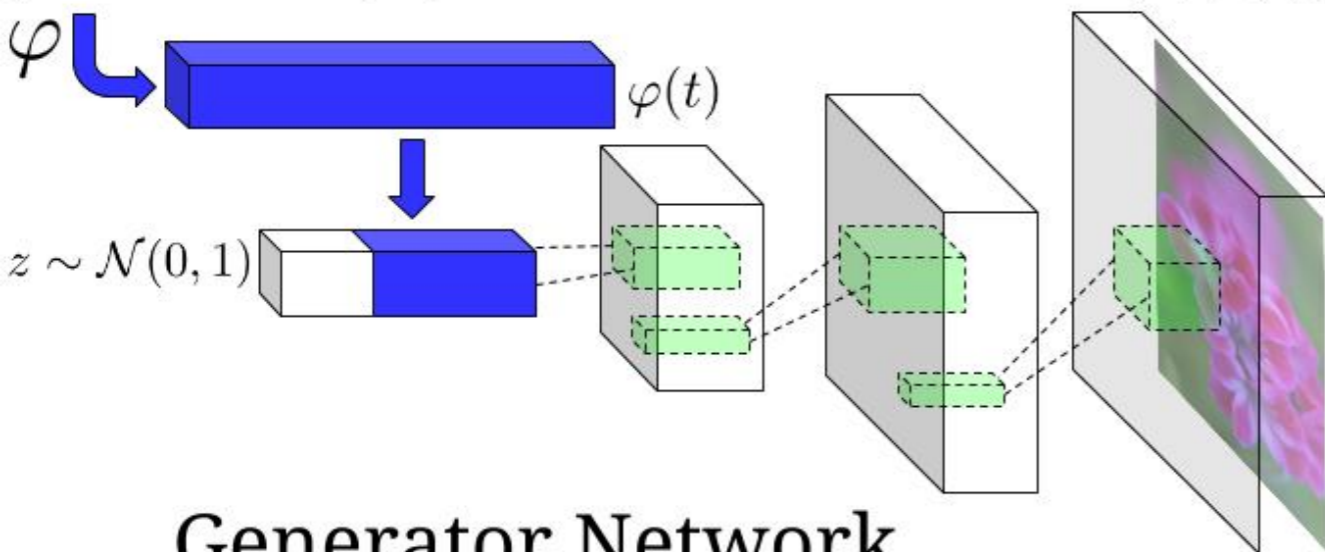
Multimodality

- A mapping between text and pixels should be learned: GAN is used
- In GAN, the generator network tries to fool adversarially trained discriminator network
 - both are conditioned on text
 - Discriminator acts as a smart adaptive loss function

GAN

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

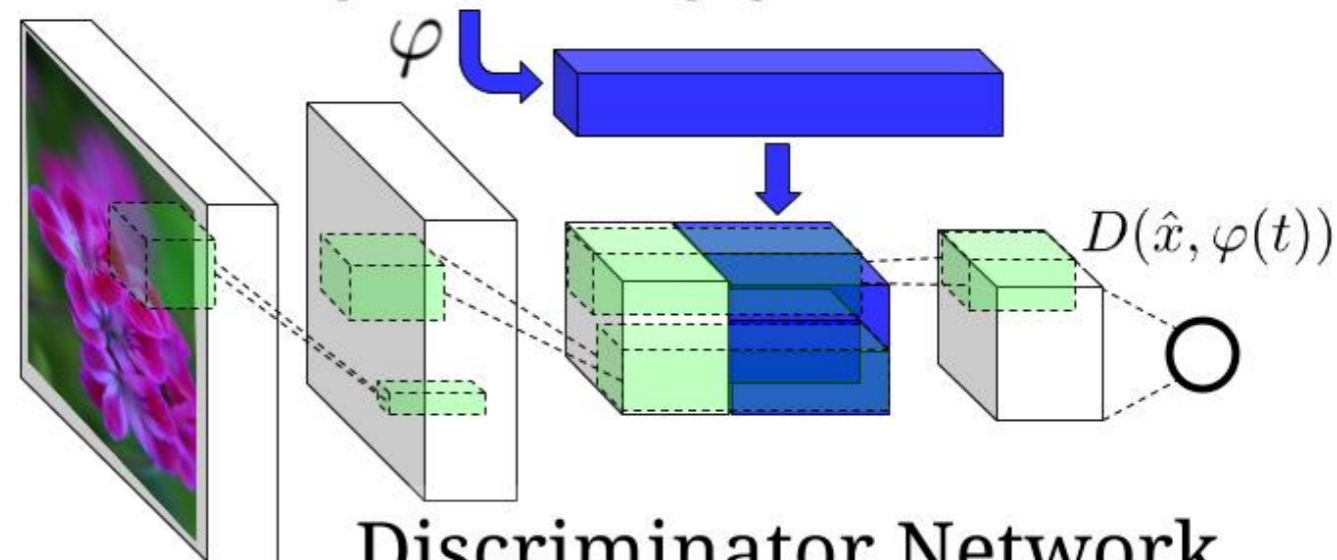
This flower has small, round violet petals with a dark purple center



Generator Network

- Fully connected layer (dim-reduc.)
- Leaky ReLU
- Concatenation
- Deconvolution

This flower has small, round violet petals with a dark purple center



Discriminator Network

- Several layers of stride-2 conv. (with spatial batch normalization)
- Leaky ReLU
- Fully connected layer (dim-reduc.) + rectification (text)
- Depth Concatenation
- conv, rectification, conv.

GAN - CLS

- Naive GAN : **<real img, matching text>** : unrealistic images contribute learning, **<synthetic img, arbitrary text>**: wrong class contributes learning
- GAN CLS: **GAN + <real image, mismatched text>** : should be scored as fake, an additional signal provided by discriminator

GAN - CLS

Algorithm 1 GAN-CLS training algorithm with step size α , 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**
-

GAN - INT

- Based on the observation that interpolations between embeddings tend to be near the data manifold, extra amount of text embeddings can be generated (although they don't have a matching text/images, they are useful for D)
- A term added to generator objective:

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

Style Transfer

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

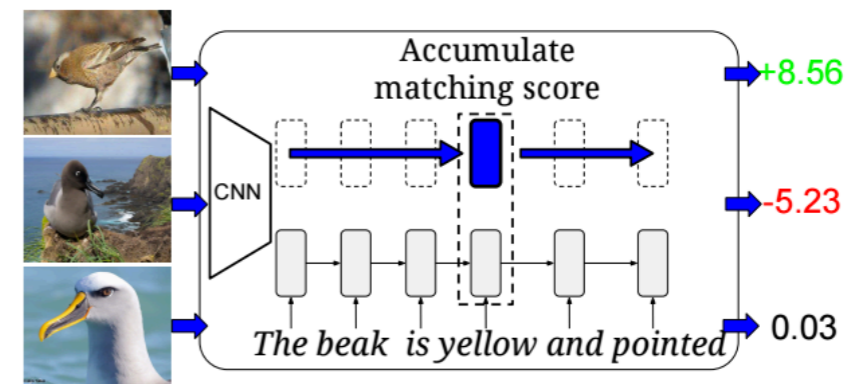
$$\mathcal{L}_{style} = \mathbb{E}_{t, z \sim \mathcal{N}(0,1)} \|z - S(G(z, \varphi(t)))\|_2^2$$

Experiments

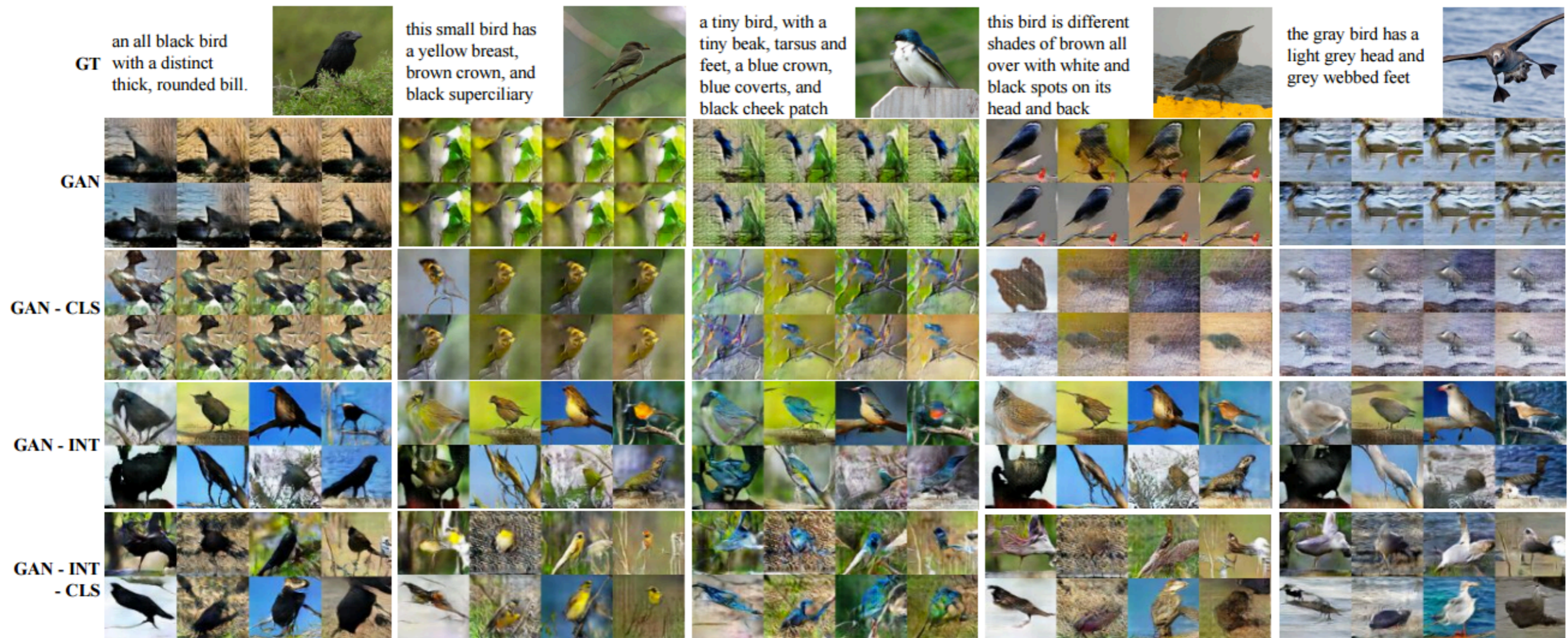
- Datasets:
 - CUB birds (11788 images, 200 classes, 5 captions per image)
 - Split to disjoint classes: 150 train+val, 50 test
 - Oxford-102 flowers (8189 images, 102 categories, 5 captions per image)
 - 82 train+val, 20 test

Experiments

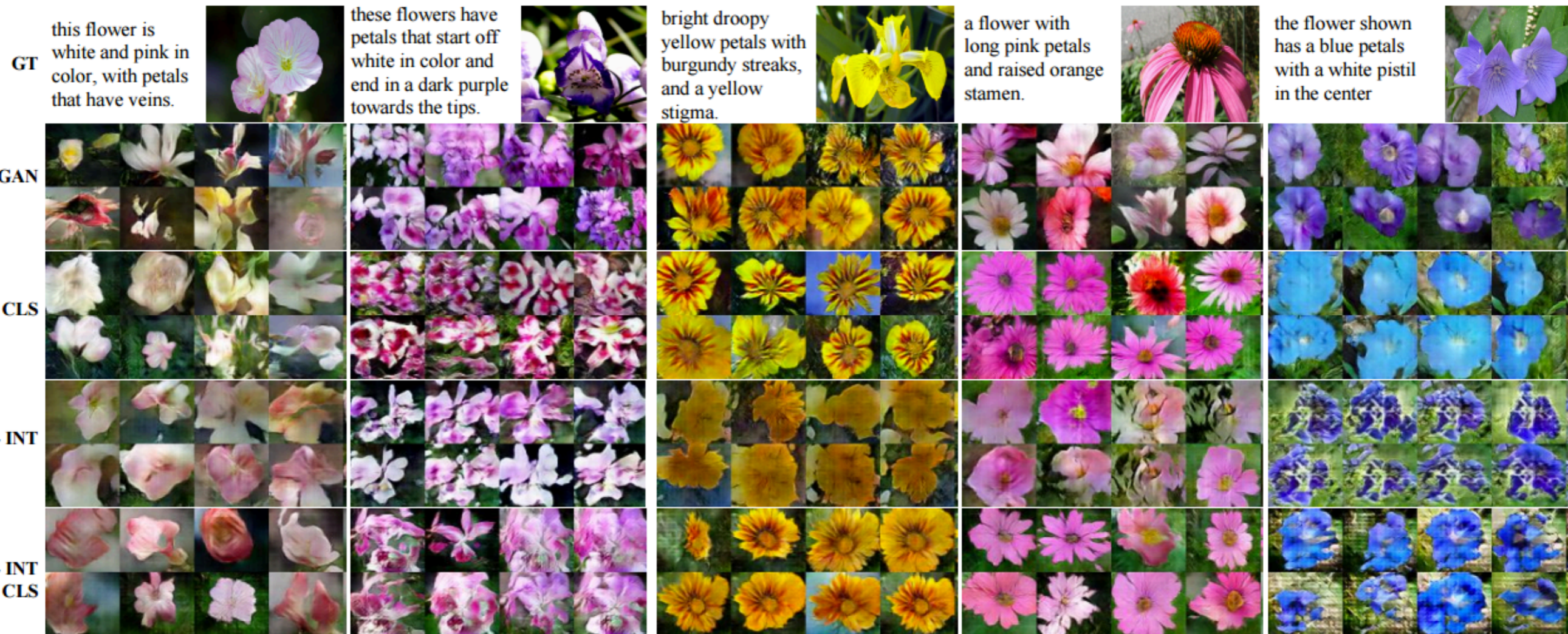
- Text features
 - Pre-training on deep deep convolutional-recurrent text encoder (char level) with Google LeNet image embeddings



Qualitative Results



Qualitative Results

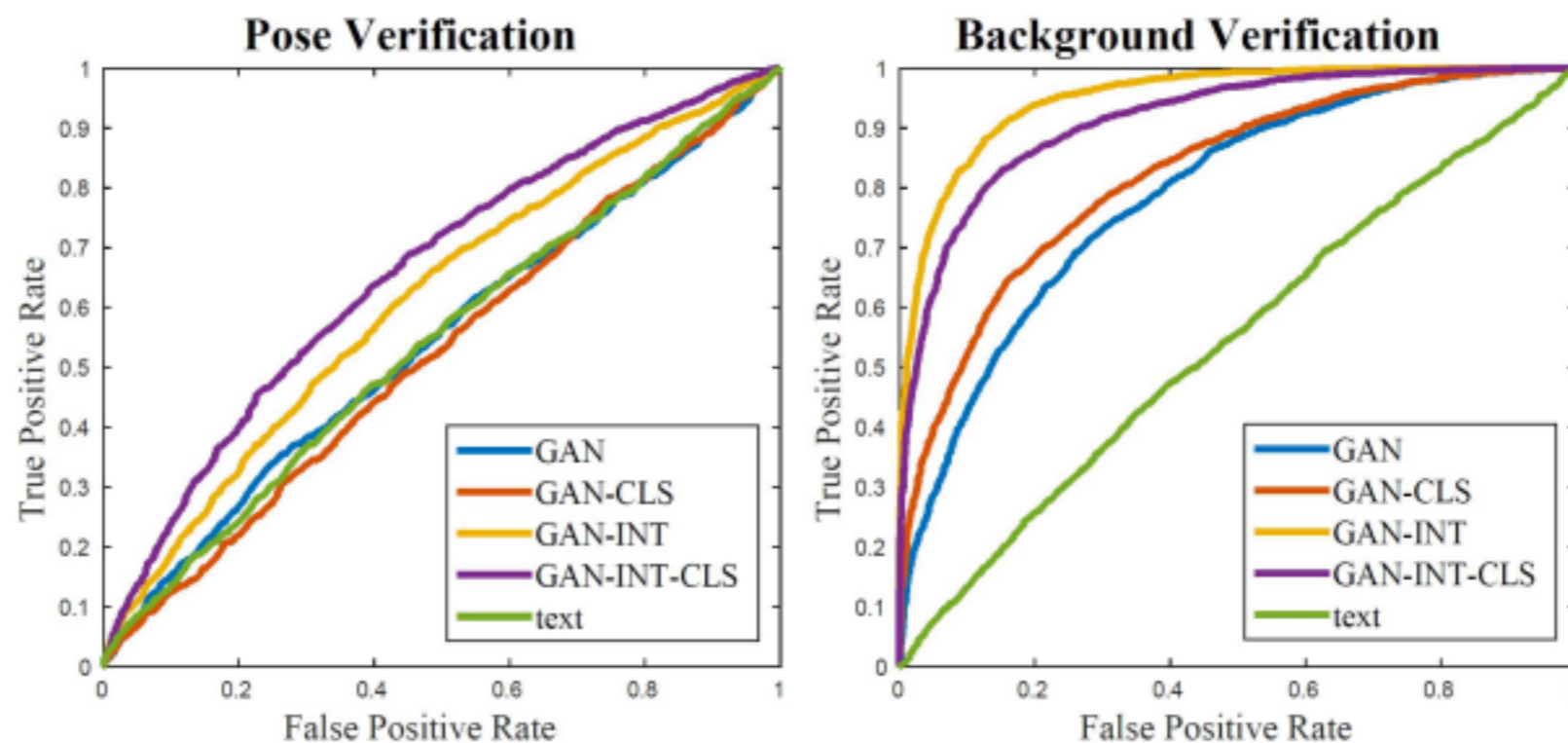


Disentangling Style and Content

- Quantification of success is based on pose verification and background verification
 - Similar pairs of images constructed for each task via K-means:
 - Avg. RGB for background color
 - Keypoint coordinates for pose

Disentangling Style and Content

- Similar and different images are fed into Style network, then cosine similarity is calculated based on the resulting encodings:



Pose and Background Style Transfer

**Text descriptions
(content)**

**Images
(style)**



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**.



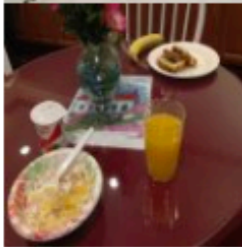

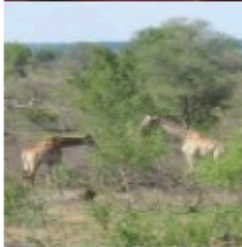













Sentence Interpolation



Figure 8. Left: Generated bird images by interpolating between two sentences (within a row the noise is fixed). Right: Interpolating between two randomly-sampled noise vectors.

GAN CLS on MS COCO

	GT	Ours
a group of people on skis stand on the snow.		
a table with many plates of food and drinks		
two giraffe standing next to each other in a forest.		
a large blue octopus kite flies above the people having fun at the beach.		

	GT	Ours
a man in a wet suit riding a surfboard on a wave.		
two plates of food that include beans, guacamole and rice.		
a green plant that is growing out of the ground.		
there is only one horse in the grassy field.		

	GT	Ours
a pitcher is about to throw the ball to the batter.		
a picture of a very clean living room.		
a sheep standing in a open grass field.		
a toilet in a small room with a window and unfinished walls.		

THANKS FOR YOUR ATTENTION