

Generative Models: Recent Progresses and Applications

CEDL 2016

Topics

- Inference?
 - Inverse mapping to the latent space for GAN
- Extensions and applications

Adversarially Learned Inference

Vincent Dumoulin*

vincent.dumoulin@umontreal.ca

Ishmael Belghazi*

ishmael.belghazi@gmail.com

Ben Poole [†]

poole@cs.stanford.edu

Alex Lamb*

alex6200@gmail.com

Martin Arjovsky*

martinarjovsky@gmail.com

Olivier Mastropietro*

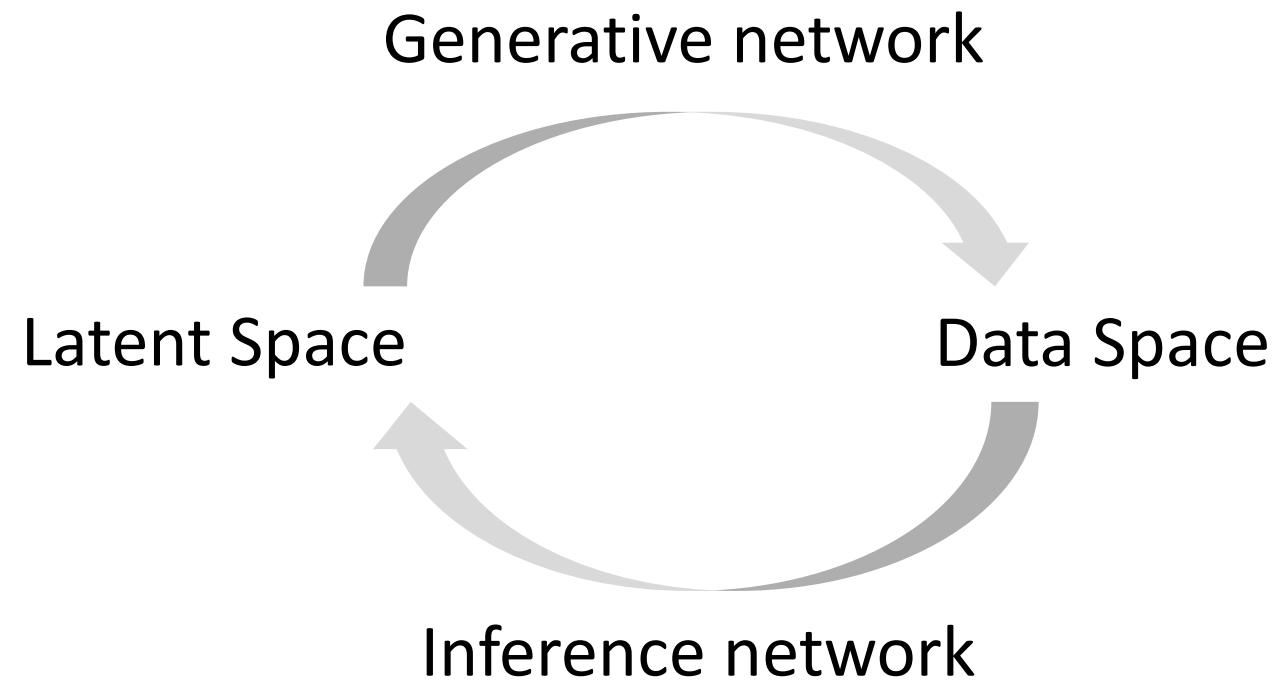
oli.mastro@gmail.com

Aaron Courville,* CIFAR Fellow

aaron.courville@gmail.com

Abstract

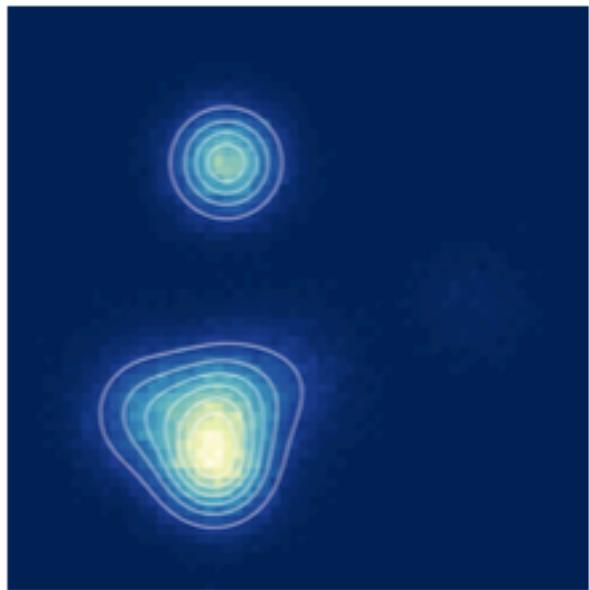
We introduce the adversarially learned inference (ALI) model, which jointly learns a generation network and an inference network using an adversarial process. The generation network maps samples from stochastic latent variables to the data space while the inference network maps training examples in data space to the space of latent variables. An adversarial game is cast between these two networks and a discriminative network that is trained to distinguish between joint latent/data-space samples from the generative network and joint samples from the inference network. We illustrate the ability of the model to learn mutually coherent inference and generation networks through the inspections of model samples and reconstructions and confirm the usefulness of the learned representations by obtaining a performance competitive with other recent approaches on the semi-supervised SVHN task.



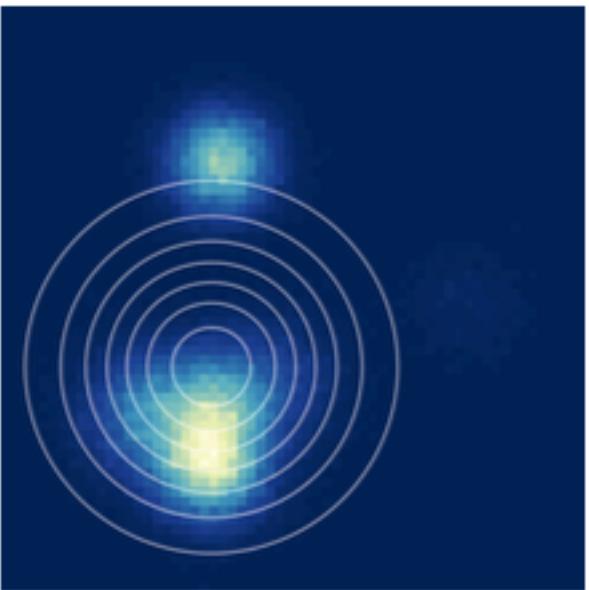
VAE vs. GAN

- VAE
 - Maximizing log-likelihood $p(x)$
 - Approximate posterior $q(z|x)$
 - Diffusion problem: All probability mass lying outside the relatively restrictive subset of pixel space occupied by natural images
 - Blurry images
- GAN
 - Hard to extend to semi-supervised learning
 - No efficient inference mechanism

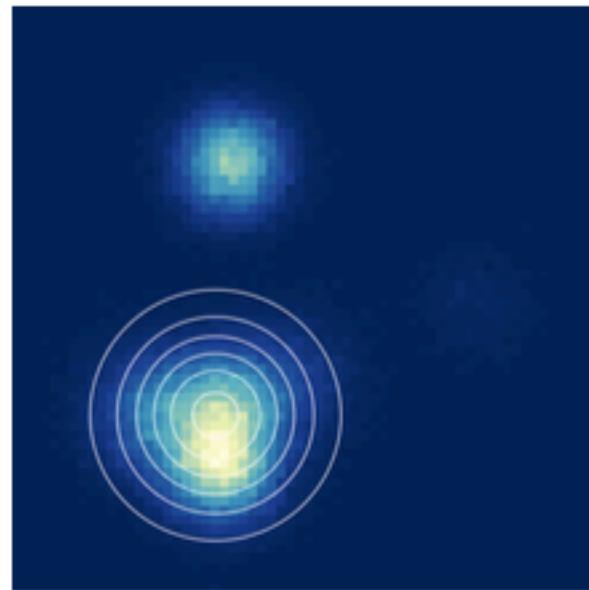
Data



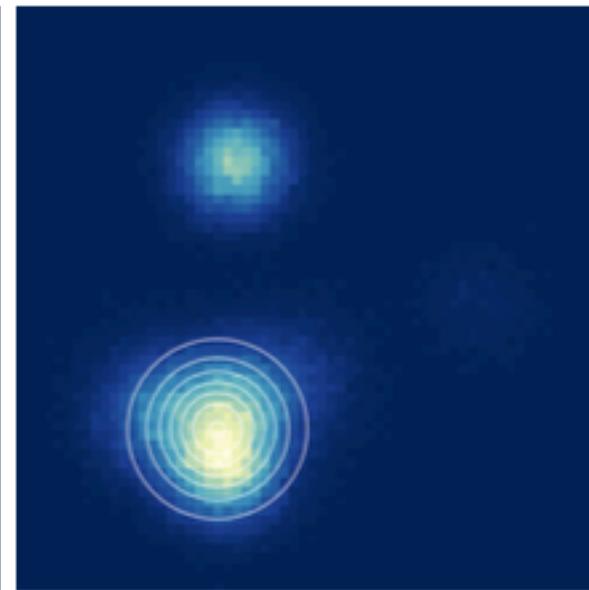
KLD

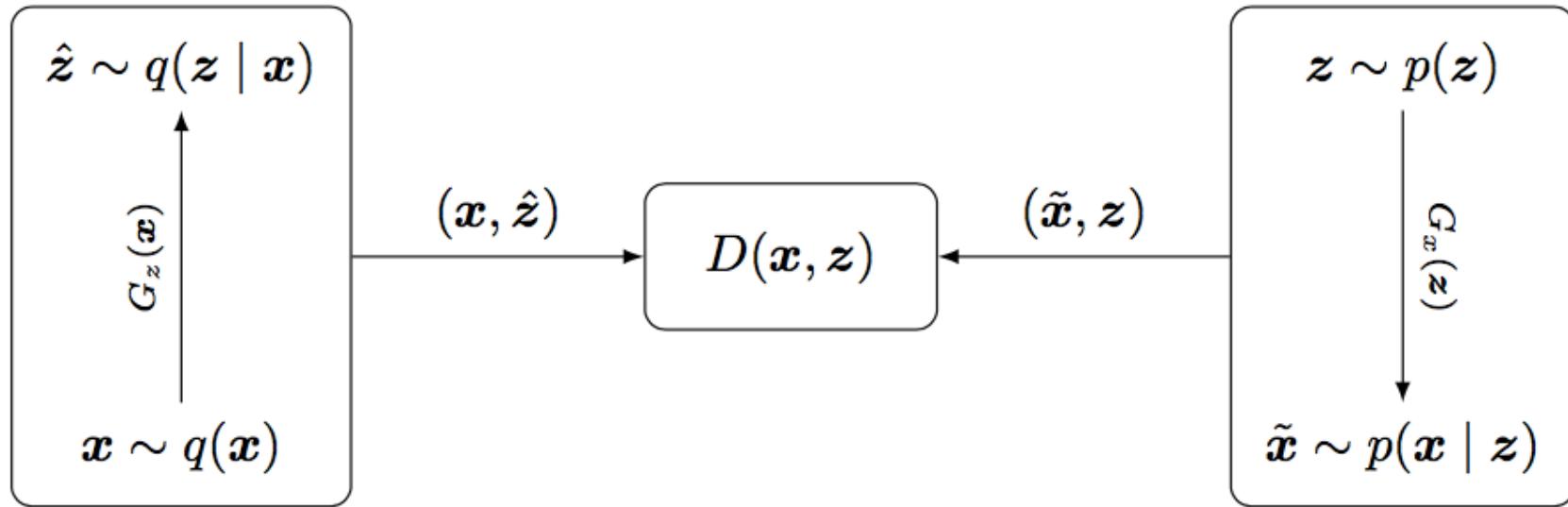


MMD



JSD





- the *encoder* joint distribution $q(\mathbf{x}, \mathbf{z}) = q(\mathbf{x})q(\mathbf{z} | \mathbf{x})$,
- the *decoder* joint distribution $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x} | \mathbf{z})$.

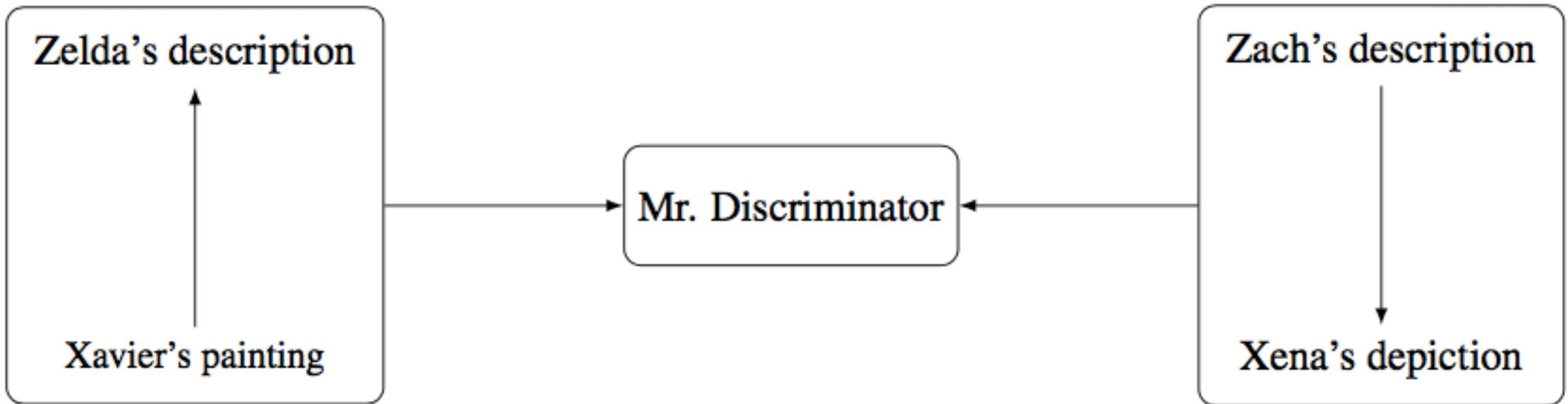


Figure 8: A Circle of Infinite Painters' view of the ALI game.

These two distributions have marginals that are known to us: the encoder marginal $q(\mathbf{x})$ is the empirical data distribution and the decoder marginal $p(\mathbf{z})$ is usually defined to be a simple, factorized distribution, such as the standard Normal distribution $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$. As such, the generative process between $q(\mathbf{x}, \mathbf{z})$ and $p(\mathbf{x}, \mathbf{z})$ is reversed.

Game Value Function

$$\begin{aligned}\min_G \max_D V(D, G) &= \mathbb{E}_{q(\mathbf{x})}[\log(D(\mathbf{x}, G_z(\mathbf{x})))] + \mathbb{E}_{p(\mathbf{z})}[\log(1 - D(G_x(\mathbf{z}), \mathbf{z}))] \\ &= \iint q(\mathbf{x})q(\mathbf{z} \mid \mathbf{x}) \log(D(\mathbf{x}, \mathbf{z})) d\mathbf{x} d\mathbf{z} \\ &\quad + \iint p(\mathbf{z})p(\mathbf{x} \mid \mathbf{z}) \log(1 - D(\mathbf{x}, \mathbf{z})) d\mathbf{x} d\mathbf{z}.\end{aligned}$$

Sampling from $q(\mathbf{z} \mid \mathbf{x})$ amounts to drawing an \mathbf{x} from the data distribution and propagating the sample through the encoder network $(\mu_z, \sigma_z) = G_z(\mathbf{x})$, then generating $\hat{\mathbf{z}} \sim \mathcal{N}(\mu_z, \sigma_z)$. Sampling from $p(\mathbf{x} \mid \mathbf{z})$ amounts to sampling from $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ and propagating the sample through the decoder network to recover $\tilde{\mathbf{x}}$.

Generator: minimizing Jenson-Shannon Divergence between $q(\mathbf{x}, \mathbf{z})$ and $p(\mathbf{x}, \mathbf{z})$

$$\begin{aligned}\text{JSD}(P \parallel Q) &= \frac{1}{2}D(P \parallel M) + \frac{1}{2}D(Q \parallel M) \\ M &= \frac{1}{2}(P + Q)\end{aligned}$$

$$\text{JSD}_{\pi_1, \dots, \pi_n}(P_1, P_2, \dots, P_n) = H\left(\sum_{i=1}^n \pi_i P_i\right) - \sum_{i=1}^n \pi_i H(P_i)$$

Compared with GAN

- The generator has two components: the encoder, $G_z(\mathbf{x})$, which maps data samples \mathbf{x} to z -space, and the decoder $G_x(z)$, which maps samples from the prior $p(z)$ (a source of

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

- The discriminator is trained to distinguish between joint pairs $(\mathbf{x}, \hat{\mathbf{z}} = G_x(\mathbf{x}))$ and $(\tilde{\mathbf{x}} = G_x(\mathbf{z}), \mathbf{z})$, as opposed to marginal samples $\mathbf{x} \sim q(\mathbf{x})$ and $\tilde{\mathbf{x}} \sim p(\mathbf{x})$.

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

Reparametrization Trick (VAE)

$$q(z | x) = \mathcal{N}(\mu(x), \sigma^2(x)I)$$

$$z = \mu(x) + \sigma(x) \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

Maximizing, preventing zero gradient

$$V'(D, G) = \mathbb{E}_{q(\mathbf{x})}[\log(1 - D(\mathbf{x}, G_{\mathbf{z}}(\mathbf{x})))] + \mathbb{E}_{p(\mathbf{z})}[\log(D(G_{\mathbf{x}}(\mathbf{z}), \mathbf{z}))]$$

$$D^*(x, z) = \frac{q(x, z)}{q(x, z) + p(x, z)}.$$

Algorithm 1 The ALI training procedure.

 $\theta_g, \theta_d \leftarrow$ initialize network parameters**repeat**

- | | |
|---|---|
| $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)} \sim q(\mathbf{x})$ | ▷ Draw M samples from the dataset and the prior |
| $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(M)} \sim p(\mathbf{z})$ | |
| $\hat{\mathbf{z}}^{(i)} \sim q(\mathbf{z} \mid \mathbf{x} = \mathbf{x}^{(i)}), \quad i = 1, \dots, M$ | ▷ Sample from the conditionals |
| $\tilde{\mathbf{x}}^{(j)} \sim p(\mathbf{x} \mid \mathbf{z} = \mathbf{z}^{(j)}), \quad j = 1, \dots, M$ | |
| $\rho_q^{(i)} \leftarrow D(\mathbf{x}^{(i)}, \hat{\mathbf{z}}^{(i)}), \quad i = 1, \dots, M$ | ▷ Compute discriminator predictions |
| $\rho_p^{(j)} \leftarrow D(\tilde{\mathbf{x}}^{(j)}, \mathbf{z}^{(j)}), \quad j = 1, \dots, M$ | |
| $\mathcal{L}_d \leftarrow -\frac{1}{M} \sum_{i=1}^M \log(\rho_q^{(i)}) - \frac{1}{M} \sum_{j=1}^M \log(1 - \rho_p^{(j)})$ | ▷ Compute discriminator loss |
| $\mathcal{L}_g \leftarrow -\frac{1}{M} \sum_{i=1}^M \log(1 - \rho_q^{(i)}) - \frac{1}{M} \sum_{j=1}^M \log(\rho_p^{(j)})$ | ▷ Compute generator loss |
| $\theta_d \leftarrow \theta_d - \nabla_{\theta_d} \mathcal{L}_d$ | ▷ Gradient update on discriminator network |
| $\theta_g \leftarrow \theta_g - \nabla_{\theta_g} \mathcal{L}_g$ | ▷ Gradient update on generator networks |

until convergence

Latent Space Interpolation



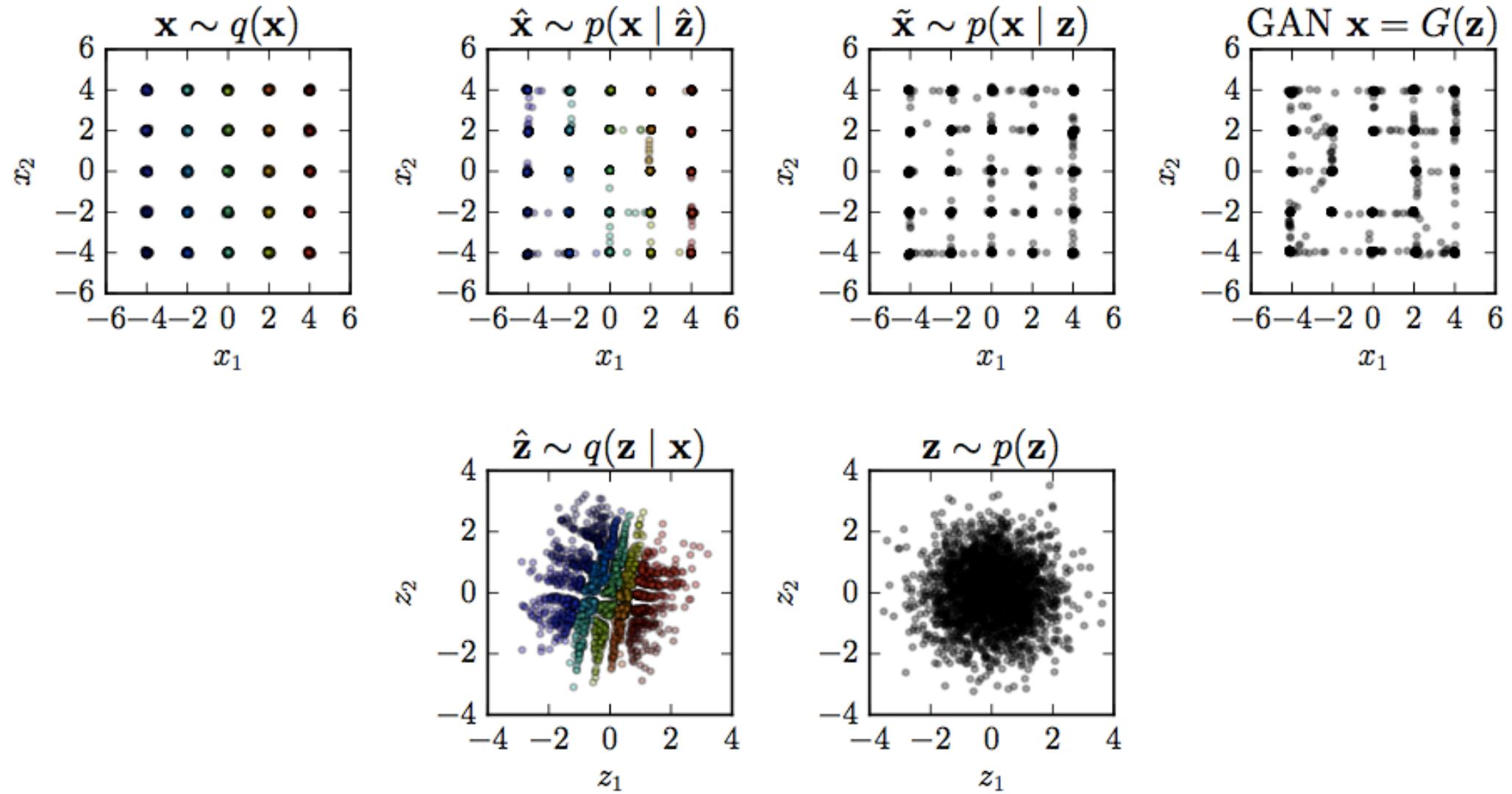


Figure 7: Comparison of ALI and GAN on a synthetic gaussian mixture dataset. In the top row, from left to right: data samples, ALI reconstructions, ALI samples, GAN samples. In the bottom row, from left to right: ALI encodings, prior samples. The nine mixture modes, their ALI encodings and their ALI reconstructions have been color-coded to facilitate visualization.

Semi-supervised

We follow the procedure outlined by Radford *et al.* (2015). We train an L2-SVM on the learned representations of a model trained on SVHN. The last three hidden layers of the encoder as well as its output are concatenated to form a 8960-dimensional feature vector. A 10,000 example held-out validation set is taken from the training set and is used for model selection. The SVM is trained on 1000 examples taken at random from the remainder of the training set. The test error rate is measured for 100 different SVMs trained on different random 1000-example training sets, and the average error rate is measured along with its standard deviation.

Method	Error rate
KNN (as reported in Zhao <i>et al.</i> (2015))	77.93%
TSVM (Vapnik, 1998)	66.55%
VAE (M1 + M2) (Kingma <i>et al.</i> , 2014)	36.02%
SWWAE without dropout (Zhao <i>et al.</i> , 2015)	27.83%
SWWAE with dropout (Zhao <i>et al.</i> , 2015)	23.56%
DCGAN + L2-SVM (Radford <i>et al.</i> , 2015)	22.18% ($\pm 1.13\%$)
SDGM (Maaløe <i>et al.</i>, 2016)	16.61% ($\pm 0.24\%$)
ALI (ours)	19.14% ($\pm 0.50\%$)

Table 1: SVHN test set error rates for semi-supervised learning. 1000 training examples were used.

	Operation	Kernel	Strides	Feature maps	BN?	Dropout	Nonlinearity						
$G_z(x) - 3 \times 32 \times 32$ input													
Convolution	5×5	1×1	32		✓	0.0	Leaky ReLU						
Convolution	4×4	2×2	64		✓	0.0	Leaky ReLU						
Convolution	4×4	1×1	128		✓	0.0	Leaky ReLU						
Convolution	4×4	2×2	256		✓	0.0	Leaky ReLU						
Convolution	4×4	1×1	512		✓	0.0	Leaky ReLU						
Convolution	1×1	1×1	512		✓	0.0	Leaky ReLU						
Convolution	1×1	1×1	512		✗	0.0	Linear						
$G_x(z) - 256 \times 1 \times 1$ input													
Transposed convolution	4×4	1×1	256		✓	0.0	Leaky ReLU						
Transposed convolution	4×4	2×2	128		✓	0.0	Leaky ReLU						
Transposed convolution	4×4	1×1	64		✓	0.0	Leaky ReLU						
Transposed convolution	4×4	2×2	32		✓	0.0	Leaky ReLU						
Transposed convolution	5×5	1×1	32		✓	0.0	Leaky ReLU						
Convolution	1×1	1×1	32		✓	0.0	Leaky ReLU						
Convolution	1×1	1×1	3		✗	0.0	Sigmoid						
$D(x) - 3 \times 32 \times 32$ input													
Convolution	5×5	1×1	32		✗	0.2	Leaky ReLU						
Convolution	4×4	2×2	64		✓	0.2	Leaky ReLU						
Convolution	4×4	1×1	128		✓	0.2	Leaky ReLU						
Convolution	4×4	2×2	256		✓	0.2	Leaky ReLU						
Convolution	4×4	1×1	512		✓	0.2	Leaky ReLU						
$D(z) - 256 \times 1 \times 1$ input													
Convolution	1×1	1×1	512		✗	0.2	Leaky ReLU						
Convolution	1×1	1×1	512		✗	0.2	Leaky ReLU						
$D(x, z) - 1024 \times 1 \times 1$ input													
Concatenate $D(x)$ and $D(z)$ along the channel axis													
Convolution	1×1	1×1	1024		✗	0.2	Leaky ReLU						
Convolution	1×1	1×1	1024		✗	0.2	Leaky ReLU						
Convolution	1×1	1×1	1		✗	0.2	Sigmoid						
Optimizer	Adam ($\alpha = 10^{-4}$, $\beta_1 = 0.5$, $\beta_2 = 10^{-3}$)												
Batch size	100												
Epochs	100												
Leaky ReLU slope 0.01													
Weight, bias initialization Isotropic gaussian ($\mu = 0$, $\sigma = 0.01$), Constant(0)													

Table 3: SVHN model hyperparameters.

Adversarial Feature Learning

Jeff Donahue, Philipp Krähenbühl, Trevor Darrell

Computer Science Division

University of California, Berkeley

{jdonahue,philk,trevor}@cs.berkeley.edu

Abstract

The ability of the Generative Adversarial Networks (GANs) framework to learn generative models mapping from simple latent distributions to arbitrarily complex data distributions has been demonstrated empirically, with compelling results showing generators learn to “linearize semantics” in the latent space of such models. Intuitively, such latent spaces may serve as useful feature representations for auxiliary problems where semantics are relevant. However, in their existing form, GANs have no means of learning the inverse mapping – projecting data back into the latent space. We propose Bidirectional Generative Adversarial Networks (BiGANs) as a means of learning this inverse mapping, and demonstrate that the resulting learned feature representation is useful for auxiliary supervised discrimination tasks, competitive with contemporary approaches to unsupervised and self-supervised feature learning.

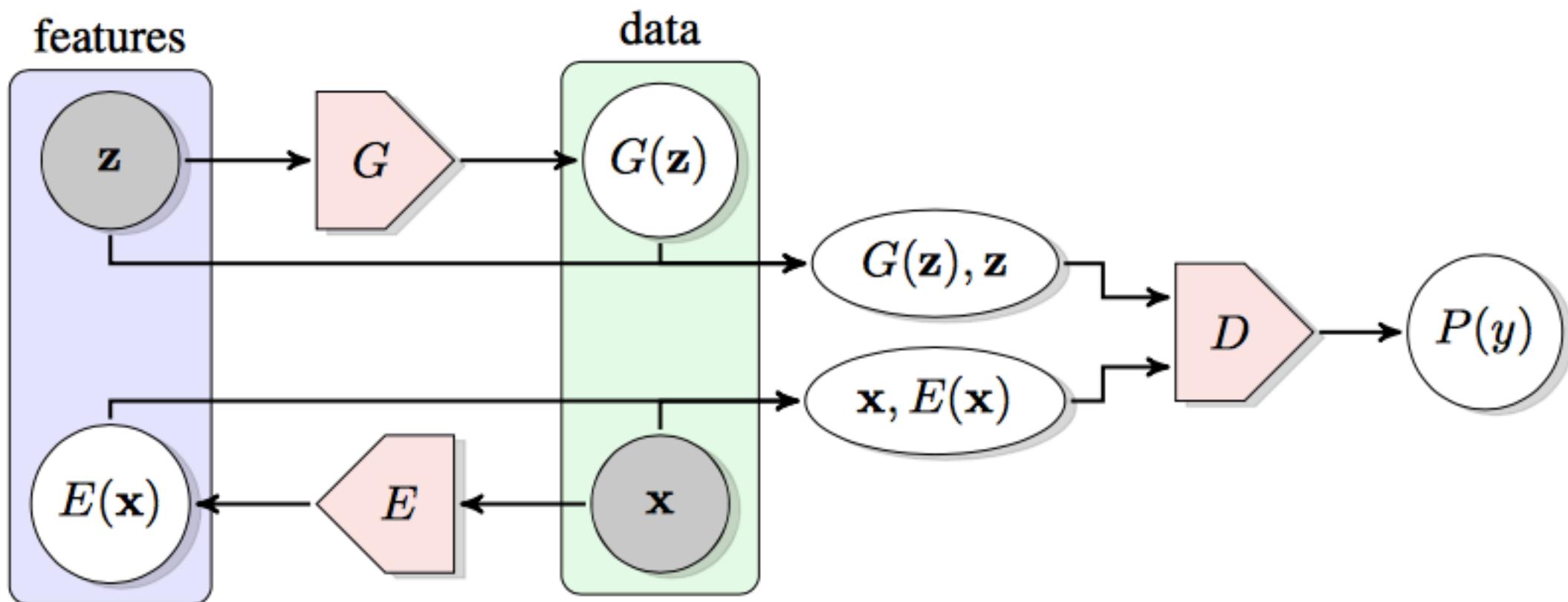


Figure 1: The structure of a Bidirectional Generative Adversarial Network (BiGAN).

$$V(D, E, G) = \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{X}}} \left[\underbrace{\mathbb{E}_{\mathbf{z} \sim p_E(\cdot | \mathbf{x})} [\log D(\mathbf{x}, \mathbf{z})]}_{\log D(\mathbf{x}, E(\mathbf{x}))} \right] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{Z}}} \left[\underbrace{\mathbb{E}_{\mathbf{x} \sim p_G(\cdot | \mathbf{z})} [\log (1 - D(\mathbf{x}, \mathbf{z}))]}_{\log(1 - D(G(\mathbf{z}), \mathbf{z}))} \right]$$

- Same value function
- Swapped labels
- Jenson-Shannon

Learning What and Where to Draw

Scott Reed^{1,*}

reedscot@google.com

Zeynep Akata²

akata@mpi-inf.mpg.de

Santosh Mohan¹

santoshm@umich.edu

Samuel Tenka¹

samtenka@umich.edu

Bernt Schiele²

schiele@mpi-inf.mpg.de

Honglak Lee¹

honglak@umich.edu

¹University of Michigan, Ann Arbor, USA

²Max Planck Institute for Informatics, Saarbrücken, Germany

Abstract

Generative Adversarial Networks (GANs) have recently demonstrated the capability to synthesize compelling real-world images, such as room interiors, album covers, manga, faces, birds, and flowers. While existing models can synthesize images based on global constraints such as a class label or caption, they do not provide control over pose or object location. We propose a new model, the Generative Adversarial What-Where Network (GAWWN), that synthesizes images given instructions describing what content to draw in which location. We show high-quality 128×128 image synthesis on the Caltech-UCSD Birds dataset, conditioned on both informal text descriptions and also object location. Our system exposes control over both the bounding box around the bird and its constituent parts. By modeling the conditional distributions over part locations, our system also enables conditioning on arbitrary subsets of parts (e.g. only the beak and tail), yielding an efficient interface for picking part locations. We also show preliminary results on the more challenging domain of text- and location-controllable synthesis of images of human actions on the MPII Human Pose dataset.



This bird is completely black.



This bird is bright blue.



a man in an orange jacket, black pants and a black cap wearing sunglasses skiing

Figure 1: Text-to-image examples. Locations can be specified by keypoint or bounding box.

Text- and location-controllable image synthesis

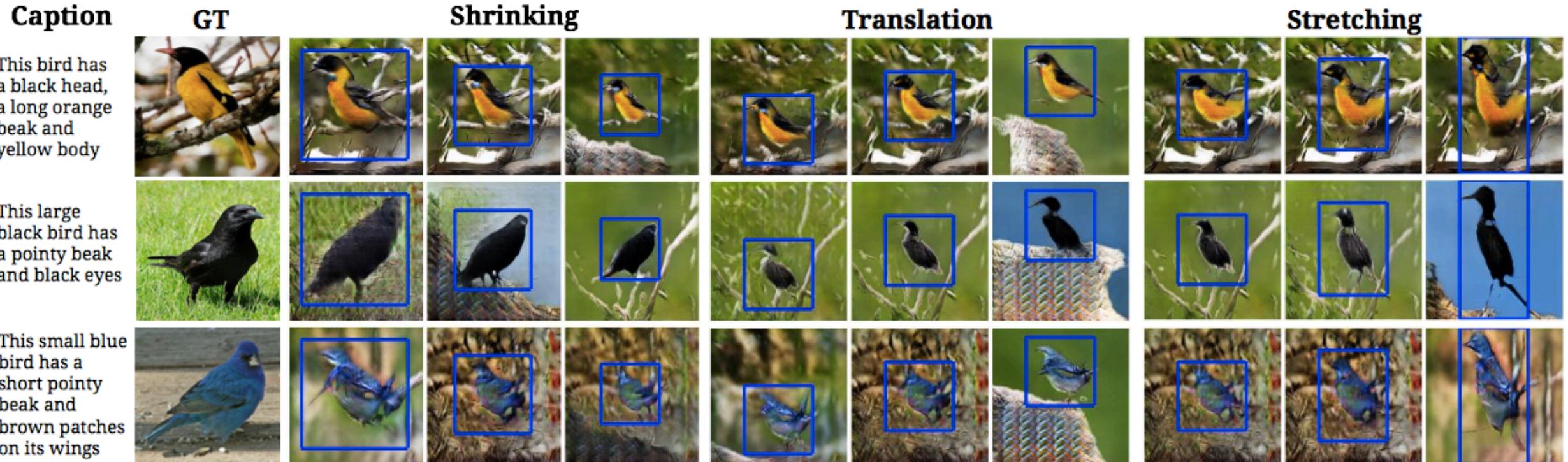


Figure 4: Controlling the bird's position using bounding box coordinates. and previously-unseen text.

Structured Joint Embedding

- Text & image

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

image encoder text encoder

training

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

Generative Adversarial What-Where Networks (GAWWN)

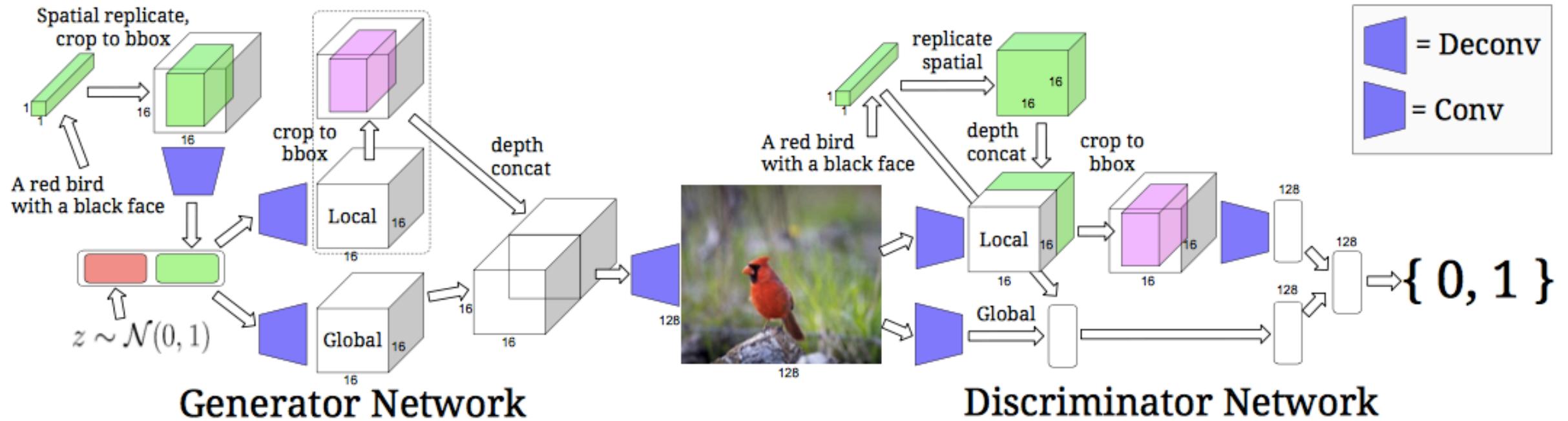


Figure 2: GAWWN with bounding box location control.

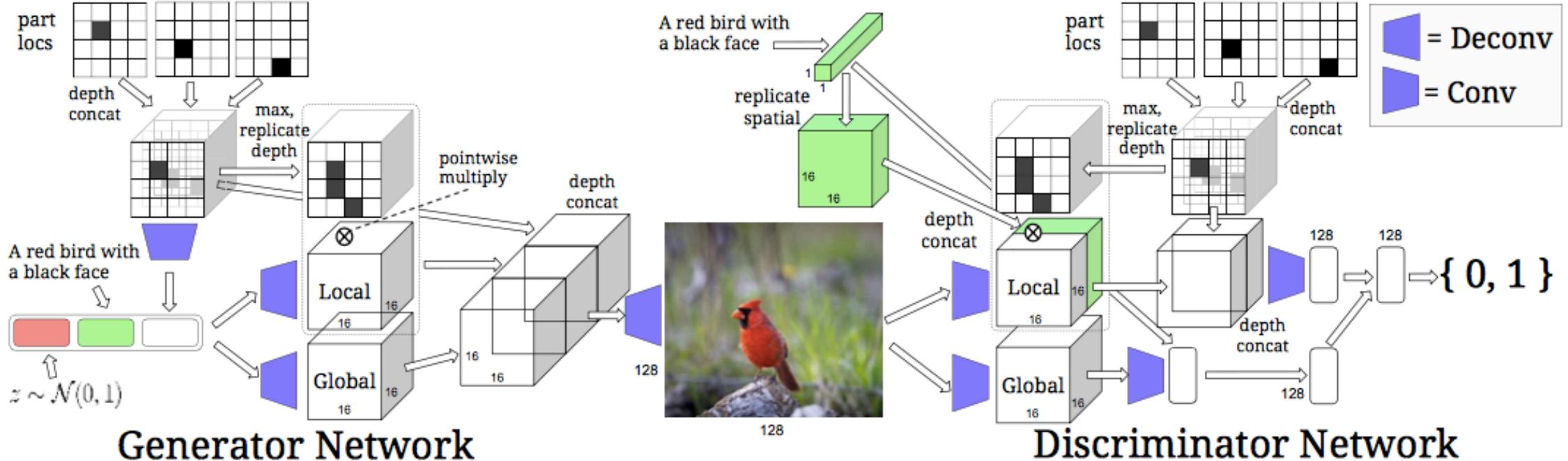


Figure 3: Text and keypoint-conditional GAWWN.. Keypoint grids are shown as 4×4 for clarity of presentation, but in our experiments we used 16×16 .

Caption	GT	Samples	Caption	GT	Samples
a woman in a yellow tank top is doing yoga.			a woman in grey shirt is doing yoga.		
the man wearing the red shirt and white pants play golf on the green grass			a man in green shirt and white pants is swinging his golf club.		
a man in a red sweater and grey pants swings a golf club with one hand.			a man in an orange jacket, black pants and a black cap wearing sunglasses skiing.		
a woman wearing goggles swimming through very murky water			a man is skiing and competing for the olympics on the slopes.		

Figure 9: Generating humans. Both the keypoints and the image are generated from unseen text.

Sampling Generative Networks

Notes on a Few Effective Techniques

Tom White
School of Design
Victoria University of Wellington
Wellington, New Zealand
Tom.White@vuw.ac.nz

- Spherical linear Interpolation
$$Slerp(q_1, q_2; \mu) = \frac{\sin (1 - \mu)\theta}{\sin \theta} q_1 + \frac{\sin \mu\theta}{\sin \theta} q_2$$
- Analogy:
 - A:B :: C:?
 - ? = C+B-A
 - A:C :: B:?

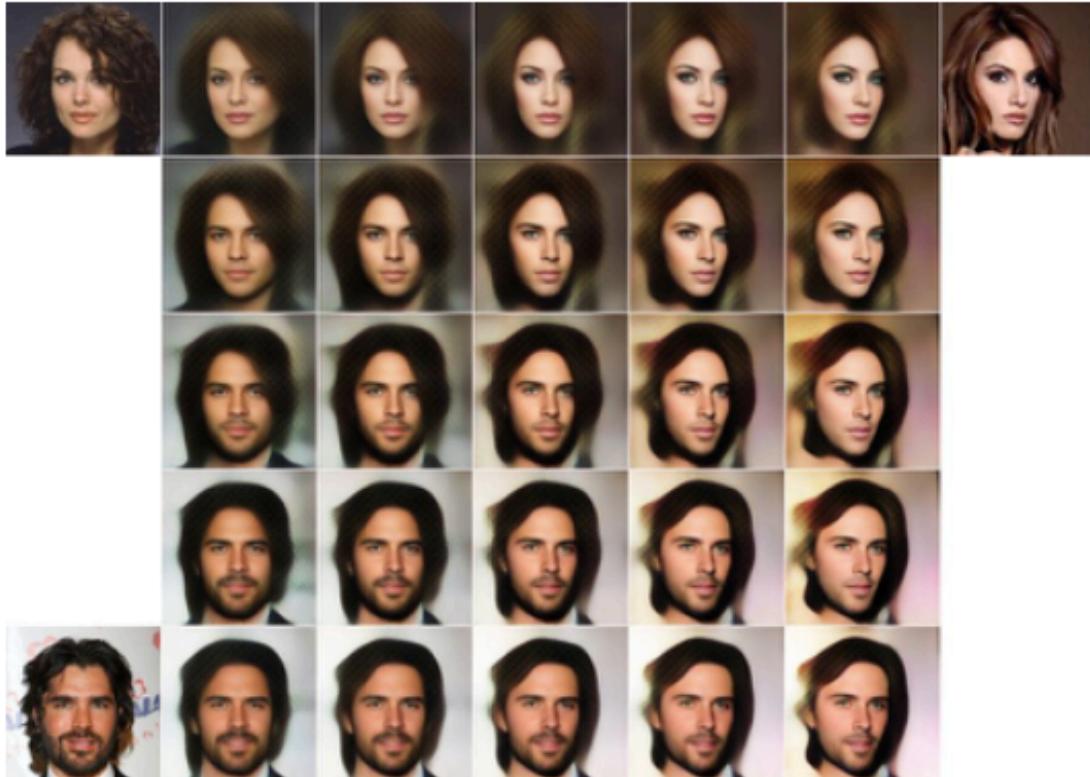


Figure 2: J-Diagram. The three corner images are inputs to the system, with the top left being the “source” (A) and the other two being “analogy targets” (B and C). Adjacent to each is the reconstruction resulting from running the image through both the encoder and decoder of the model. The bottom right image shows the result of applying the analogy operation $(B + C) - A$. All other images are interpolations using the slerp operator. (model: VAE from Lamb 16 on CelebA)

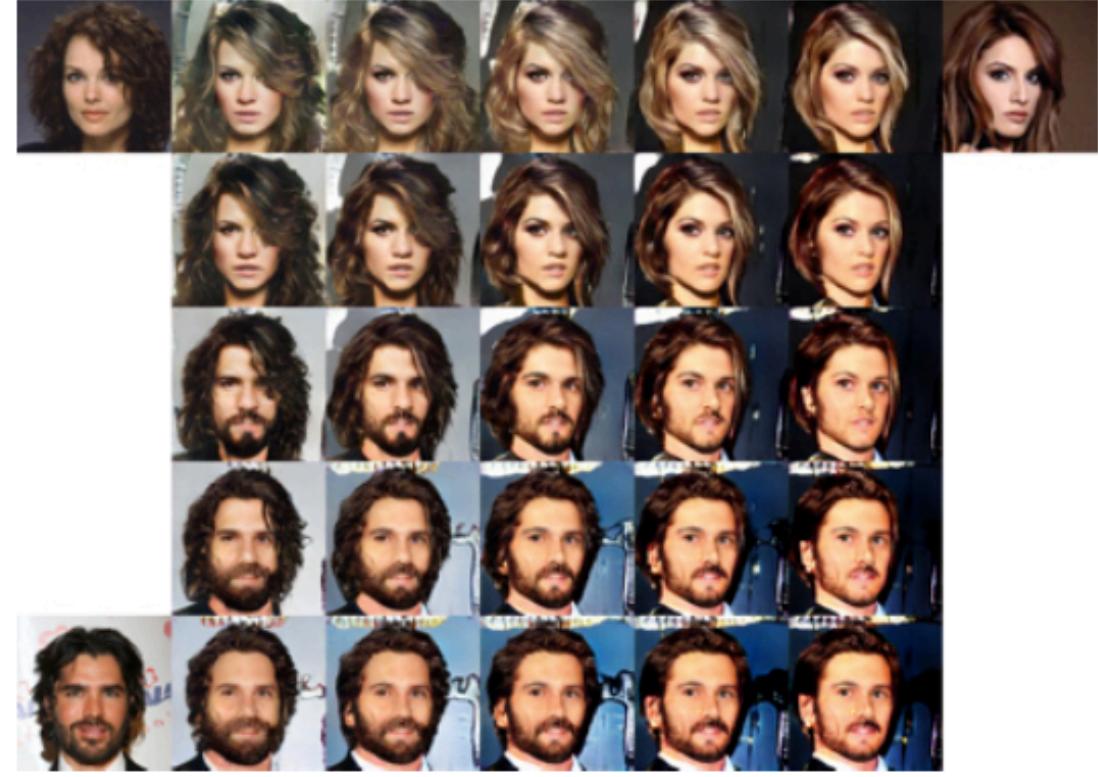
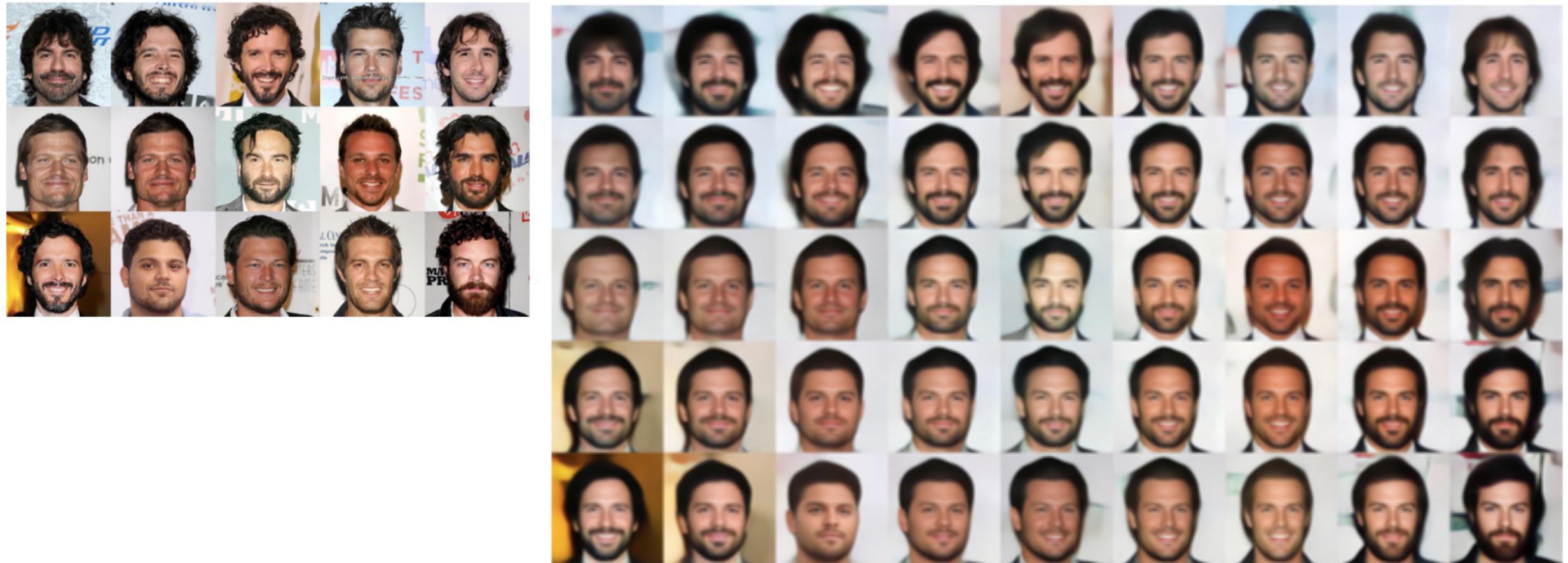


Figure 3: Same J-Diagram repeated with different model type. To facilitate comparisons (and demonstrate results are not cherry-picked) inputs selected are the first 3 images of the validation set. (model: GAN from Dumoulin 16 on CelebA)

MINE—Manifold Interpolation Neighbor Embedding

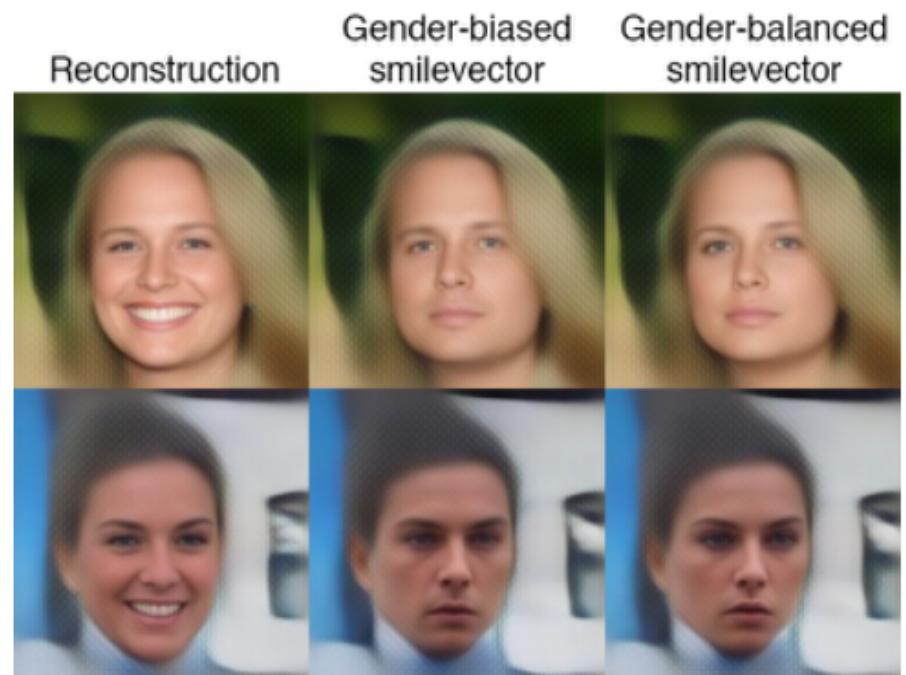


Latent Space Transformation

- Attribute
 - No smile to smile
- Correlated labels
 - Sample bias

	Male	Not Male	Total
Smile	17%	31%	48%
No Smile	25%	27%	52%
Total	42%	58%	

- → replicate data



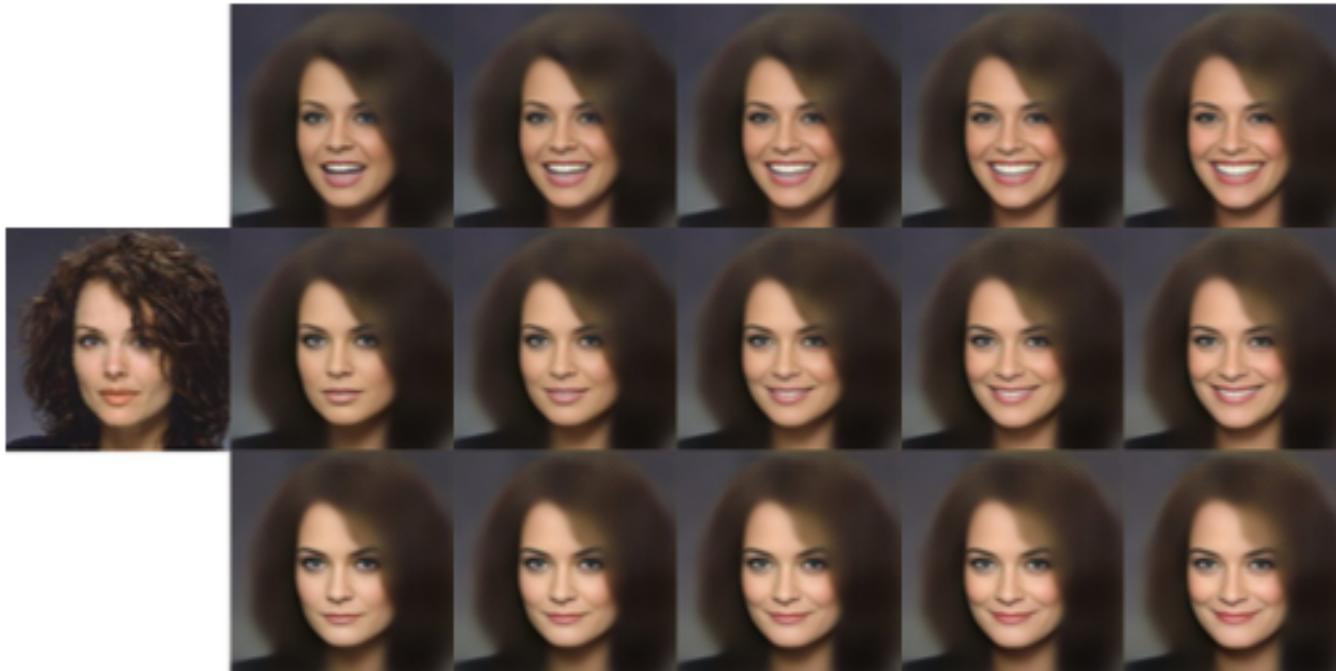


Figure 7: Decoupling attribute vectors for smiling (x-axis) and mouth open (y-axis) allows for more flexible latent space transformations. Input shown at left with reconstruction adjacent. (model: VAE from Lamb 16 on CelebA)

Synthetic Attribute

- Deblurring
 - Extrapolate in the negative direction



Task Specific Adversarial Cost Function

Antonia Creswell, *Student Member, IEEE*, and Anil A. Bharath, *Member, IEEE*

Abstract—The cost function used to train a generative model should fit the purpose of the model. If the model is intended for tasks such as generating perceptually correct samples, it is beneficial to maximise the likelihood of a sample drawn from the model, Q , coming from the same distribution as the training data, P . This is equivalent to minimising the Kullback-Leibler (KL) distance, $KL[Q||P]$. However, if the model is intended for tasks such as retrieval or classification it is beneficial to maximise the likelihood that a sample drawn from the training data is captured by the model, equivalent to minimising $KL[P||Q]$. The cost function used in adversarial training optimises the Jensen-Shannon entropy which can be seen as an even interpolation between $KL[Q||P]$ and $KL[P||Q]$. Here, we propose an alternative adversarial cost function which allows easy tuning of the model for either task. Our task specific cost function is evaluated on a dataset of hand-written characters in the following tasks: Generation, retrieval and one-shot learning.

The original cost function proposed by Goodfellow et al. [9] is:

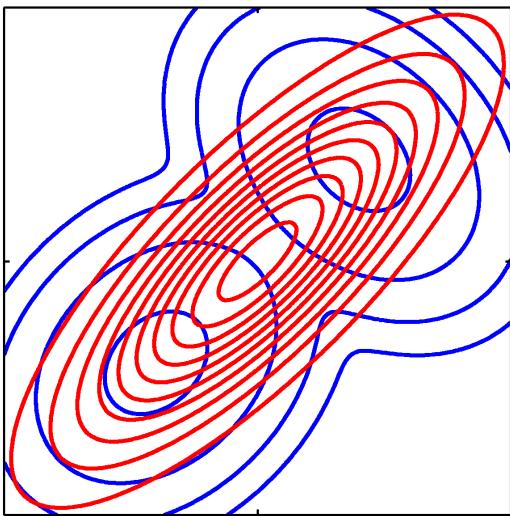
$$\min_G \max_D V(G, D) = \mathbb{E}_{x \sim P} \log D(x) + \mathbb{E}_{z \sim p_z} \log(1 - D(G(z)))$$

We propose the alternative cost function:

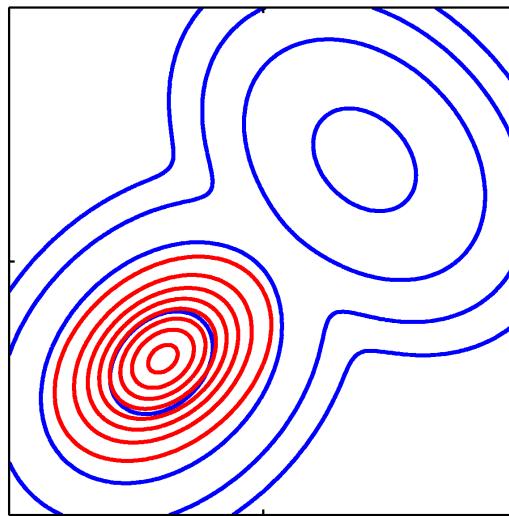
$$\begin{aligned} \min_G \max_D V(G, D) = & \pi \mathbb{E}_{x \sim P} \log D(x) + \\ & (1 - \pi) \mathbb{E}_{z \sim p_z} \log(1 - D(G(z))) \end{aligned}$$

PRML Fig. 10.3 (Bishop)

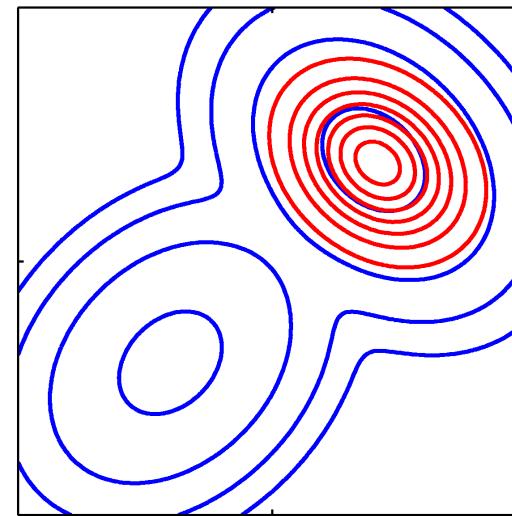
p: blue
q: red



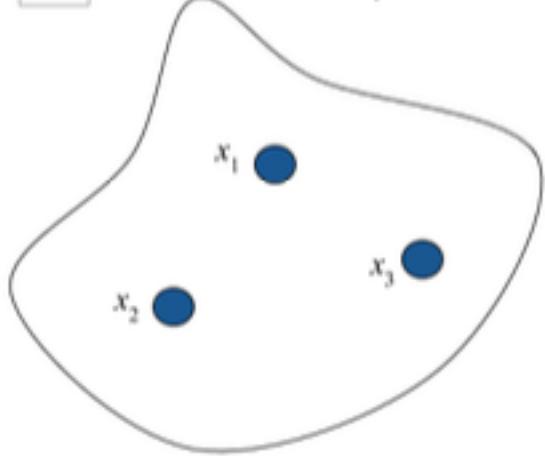
$KL(p \parallel q)$



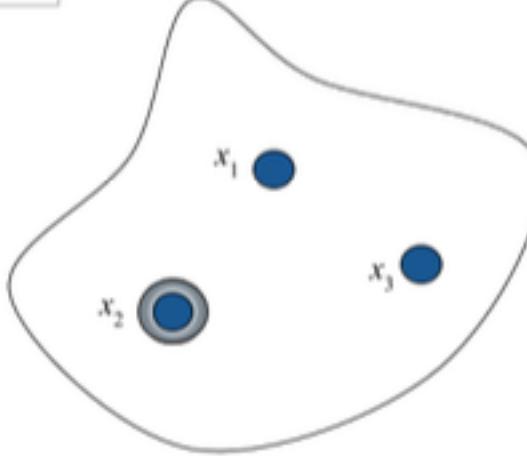
$KL(q \parallel p)$



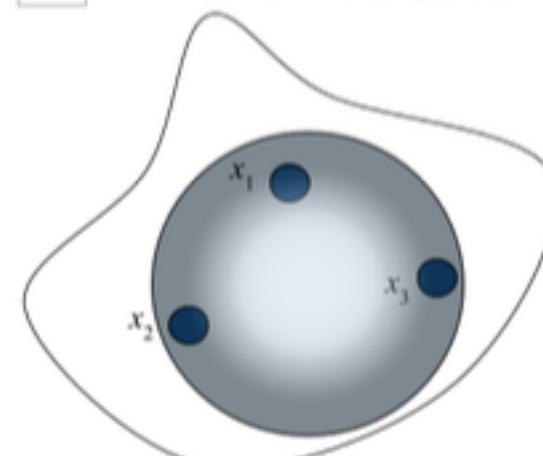
A $P(x) = \delta(x - x_i), i=\{1,2,3\}$



B $\min (KL[Q(x) || P(x)])$



C $\min (KL[P(x) || Q(x)])$



$P(x)=1$



Model, $Q(x) \sim N(\mu, \sigma)$

$$KL[Q(x)||P(x)] \sim 0$$

$$KL[P(x)||Q(x)] \sim 0$$

$$JS_{\pi} = \pi KL[P||\pi P + (1-\pi)Q] + (1-\pi)KL[Q||\pi P + (1-\pi)Q]$$

Huszar et al. [14] showed that JS_{π} is proportional to $KL[P||Q]$ and $KL[Q||P]$, respectively, at the upper and lower limits of π :

$$\frac{JS_{\pi}[P||Q]}{\pi} \rightarrow KL[P||Q] \quad \text{as } \pi \rightarrow 0$$

$$\frac{JS_{\pi}[P||Q]}{(1-\pi)} \rightarrow KL[Q||P] \quad \text{as } \pi \rightarrow 1$$