

Probabilistic Graphical Models

Deep Generative Models - II

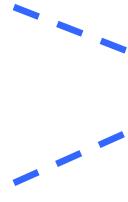
Zhiteng Hu

Lecture 18, March 25, 2019

Reading: see class homepage



Overview: Deep Learning & Generative Models

- 3/6 Lecture 15 Statistical and Algorithmic Foundations of Deep Learning
- 3/18 Lecture 16 Building blocks of DL  *Model Architectures*
- 3/20 Lecture 17 Deep generative models (part 1)  *Inference & Learning*
- 3/25 Lecture 18 Deep generative models (part 2)
- 3/27 Lecture 19 Case Study: Text Generation  *Advanced Topics*





Outline

- Generative Adversarial Networks (GANs)
 - GANs Progress
 - Vanilla GAN, Wasserstein GAN, Progressive GAN, BigGAN
- Normalizing Flow (NF)
 - Basic Concepts
 - GLOW
- Integrating Domain Knowledge into Deep Learning





Outline

- Generative Adversarial Networks (GANs)
 - GANs Progress
 - Vanilla GAN, Wasserstein GAN, Progressive GAN, BigGAN
- Normalizing Flow (NF)
 - Basic Concepts
 - GLOW
- Integrating Domain Knowledge into Deep Learning





GAN Progress on Face Generation



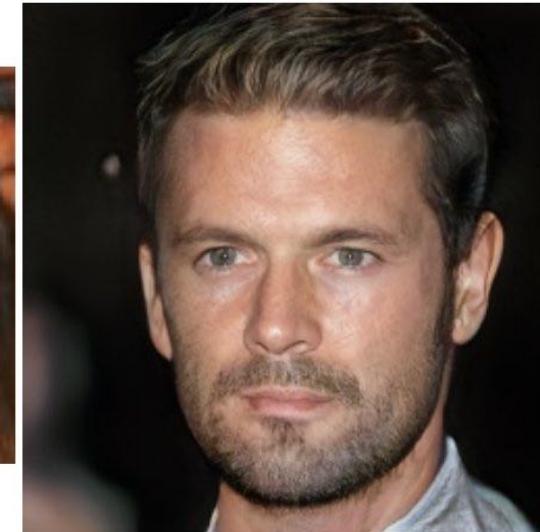
2014



2015



2016



2017



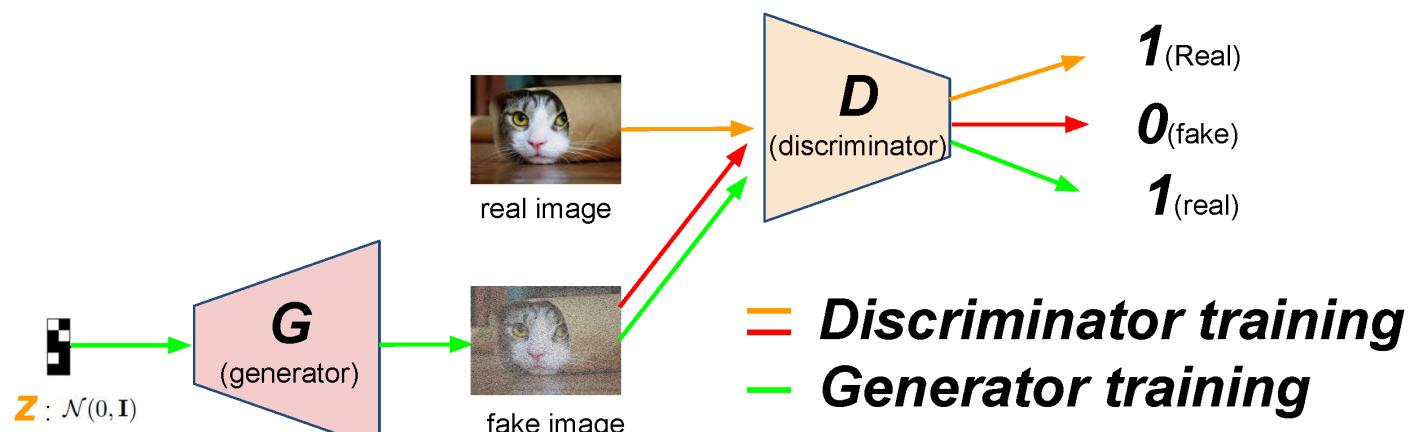
2018





Recap: Generative Adversarial Nets (GANs)

- Generative model $\mathbf{x} = G_{\theta}(\mathbf{z})$, $\mathbf{z} \sim p(\mathbf{z})$
 - Map noise variable \mathbf{z} to data space \mathbf{x}
 - Define an **implicit distribution** over \mathbf{x} : $p_{g_{\theta}}(\mathbf{x})$
 - a stochastic process to simulate data \mathbf{x}
 - Intractable to evaluate likelihood
- Discriminator $D_{\phi}(\mathbf{x})$
 - Output the probability that \mathbf{x} came from the data rather than the generator





Recap: Generative Adversarial Nets (GANs)

- Learning
 - A minimax game between the generator and the discriminator
 - Train D to maximize the probability of assigning the correct label to both training examples and generated samples
 - Train G to fool the discriminator

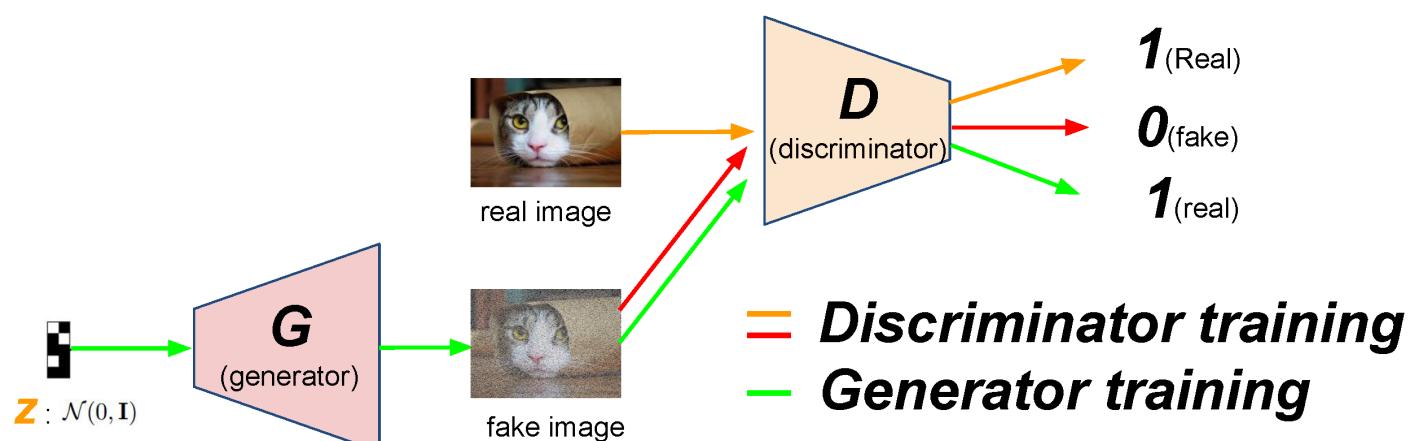
$$\begin{aligned}\max_D \mathcal{L}_D &= \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim G(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(\mathbf{x}))] \\ \min_G \mathcal{L}_G &= \mathbb{E}_{\mathbf{x} \sim G(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(\mathbf{x}))].\end{aligned}$$

- [Goodfellow et al., 2014]

$$\min_{\theta} \text{JSD}(P_{data} \parallel P_{g_{\theta}})$$

- [Hu et al., 2017]

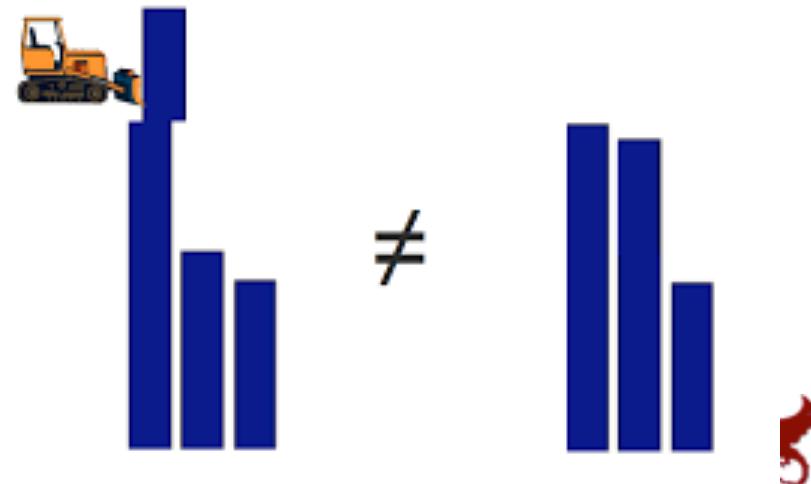
$$\min_{\theta} \text{KL}(P_{\theta} \parallel Q)$$





Wasserstein GAN (WGAN)

- If our data are on a **low-dimensional** manifold of a high dimensional space, the model's manifold and the true data manifold can have a **negligible intersection in practice**
- KL divergence is undefined or infinite
- The loss function and gradients may not be continuous and well behaved
- The **Wasserstein Distance** is well defined
 - Earth Mover's Distance
 - Minimum transportation cost for making one pile of dirt in the shape of one probability distribution to the shape of the other distribution





Wasserstein GAN (WGAN)

- Objective

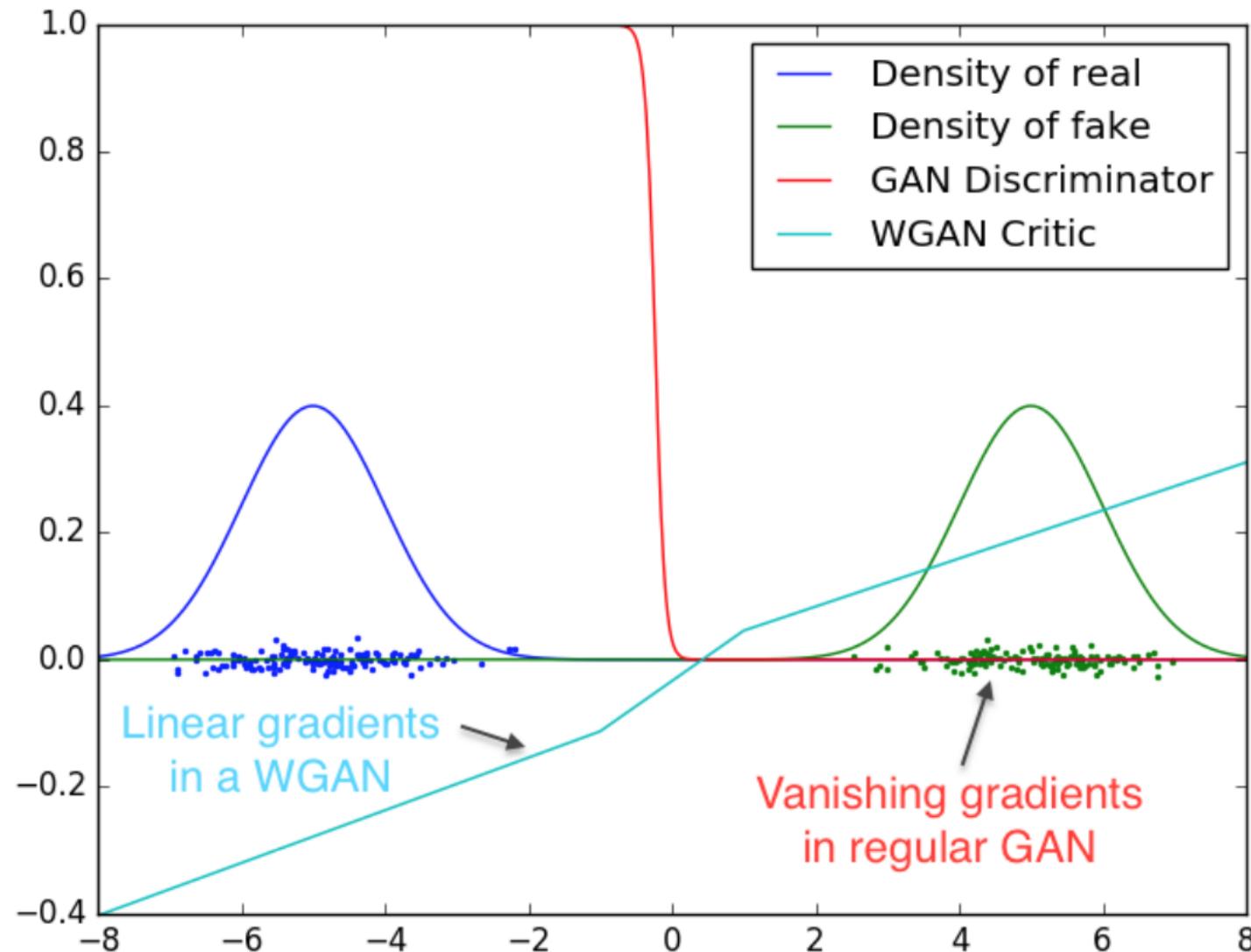
$$W(p_{data}, p_g) = \frac{1}{K} \sup_{\|D\|_L \leq K} \mathbb{E}_{x \sim p_{data}} [D(x)] - \mathbb{E}_{x \sim p_g} [D(x)]$$

- $\|D\|_L \leq K$: K- Lipschitz continuous
- Use gradient-clipping to ensure D has the Lipschitz continuity





Wasserstein GAN (WGAN)





Progressive GAN

Low resolution images

add in
additional
layers

High resolution images

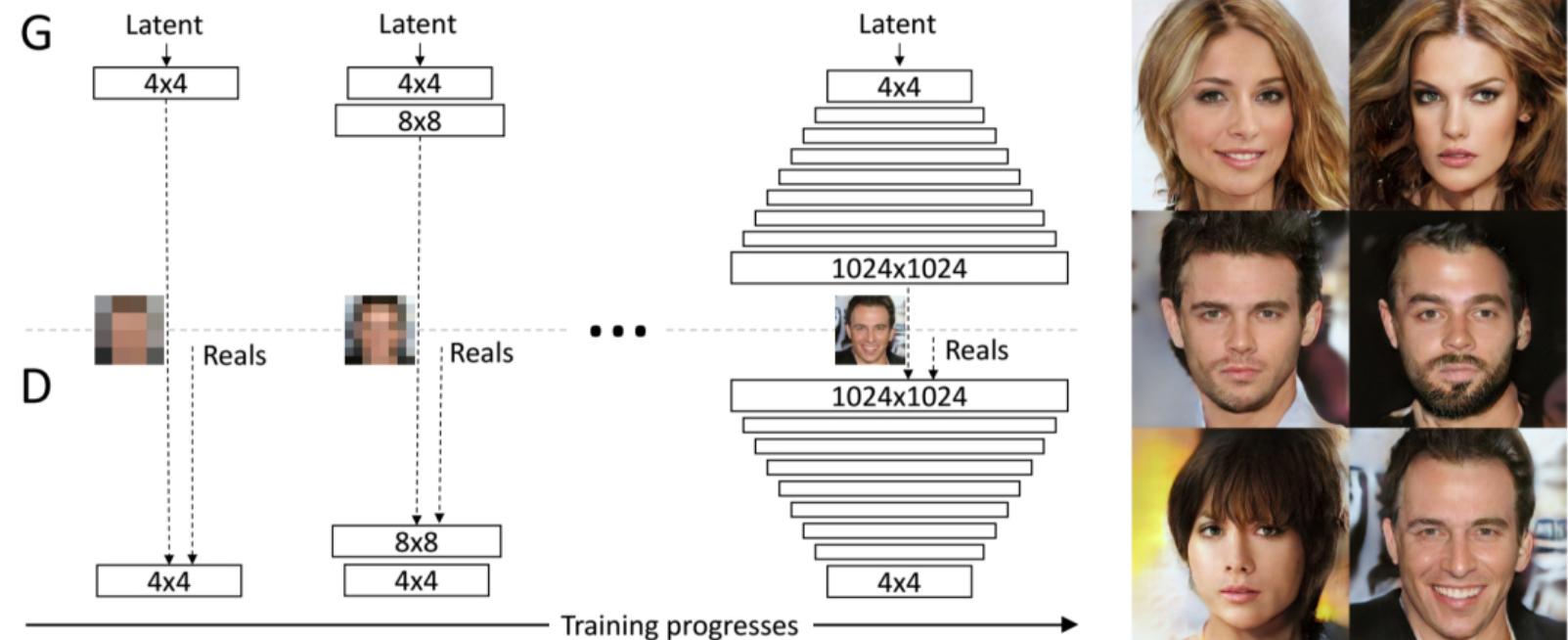


Figure 1: Our training starts with both the generator (G) and discriminator (D) having a low spatial resolution of 4×4 pixels. As the training advances, we incrementally add layers to G and D, thus increasing the spatial resolution of the generated images. All existing layers remain trainable throughout the process. Here $N \times N$ refers to convolutional layers operating on $N \times N$ spatial resolution. This allows stable synthesis in high resolutions and also speeds up training considerably. On the right we show six example images generated using progressive growing at 1024×1024 .





BigGAN

- GANs benefit dramatically from **scaling**
- 2x – 4x more parameters
- 8x larger batch size
- Simple architecture changes that improve scalability





BigGAN

- GANs benefit dramatically from scaling
- 2x
- 8x
- Sir





Outline

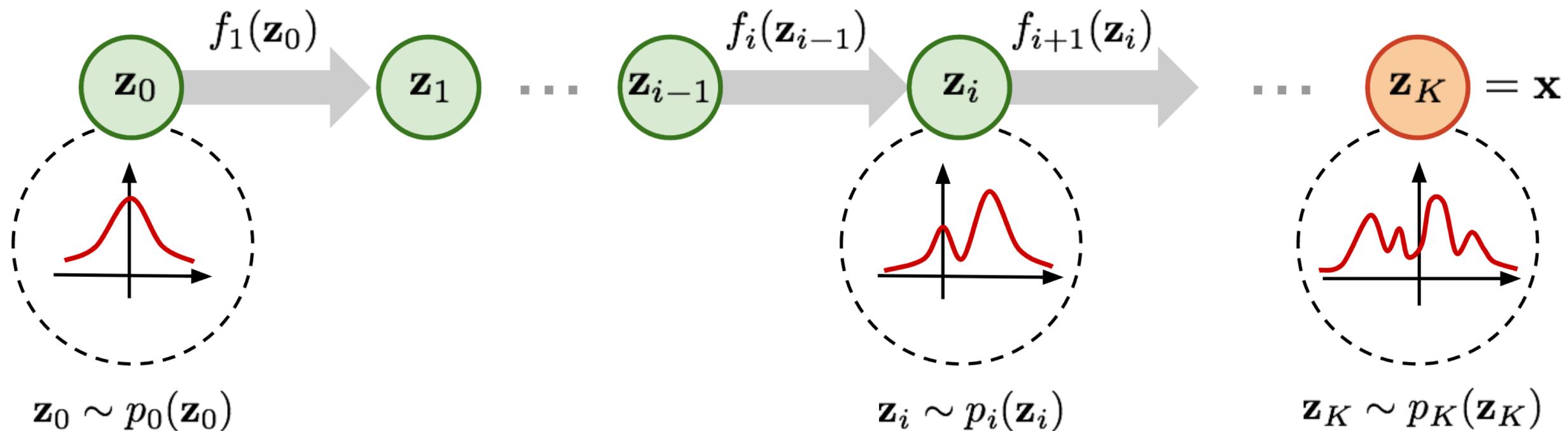
- Generative Adversarial Networks (GANs)
 - GANs Progress
 - Vanilla GAN, Wasserstein GAN, Progressive GAN, BigGAN
- Normalizing Flow (NF)
 - Basic Concepts
 - GLOW
- Integrating Domain Knowledge into Deep Learning





Normalizing Flow (NF)

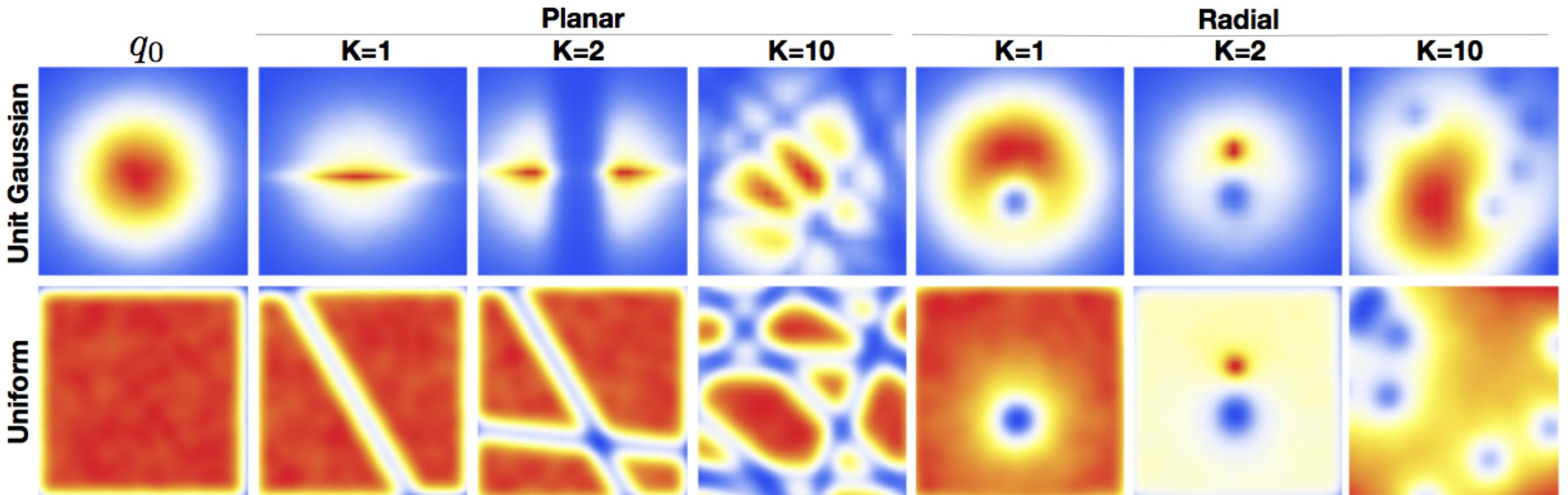
- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**





Normalizing Flow (NF)

- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**





Normalizing Flow (NF)

- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**

$$\mathbf{z} \sim p(\mathbf{z})$$

$$\mathbf{x} = f(\mathbf{z})$$

Transformation function f

inference: $\mathbf{z} = f^{-1}(\mathbf{x})$

- > • Invertible

density: $p(\mathbf{x}) = p(\mathbf{z}) \left| \det \frac{d\mathbf{z}}{d\mathbf{x}} \right|$

- > • Jacobian determinant easy to compute
e.g., choose $df_i^{-1}/d\mathbf{x}$ to be a triangular matrix

$$= p(f^{-1}(\mathbf{x})) \left| \det \frac{df^{-1}}{d\mathbf{x}} \right|$$

$\det \frac{df^{-1}}{d\mathbf{x}}$ -- Jacobian determinant





Normalizing Flow (NF)

- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**

$$\mathbf{z}_0 \sim p(\mathbf{z}_0)$$

$$\mathbf{x} = \mathbf{z}_K = f_K \circ f_{K-1} \circ \cdots \circ f_1(\mathbf{z}_0)$$

Transformation function f_i

inference: $\mathbf{z}_i = f_i^{-1}(\mathbf{z}_{i-1})$

- > • Invertible

density: $p(\mathbf{z}_i) = p(\mathbf{z}_{i-1}) \left| \det \frac{d\mathbf{z}_{i-1}}{d\mathbf{z}_i} \right|$

- > • Jacobian determinant easy to compute
e.g., choose $d\mathbf{f}_i^{-1}/d\mathbf{z}_i$ to be a triangular matrix





Normalizing Flow (NF)

- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**

$$\mathbf{z}_0 \sim p(\mathbf{z}_0)$$

$$\mathbf{x} = \mathbf{z}_K = f_K \circ f_{K-1} \circ \cdots \circ f_1(\mathbf{z}_0)$$

Transformation function f_i

inference: $\mathbf{z}_i = f_i^{-1}(\mathbf{z}_{i-1})$

- > • Invertible

density: $p(\mathbf{z}_i) = p(\mathbf{z}_{i-1}) \left| \det \frac{d\mathbf{z}_{i-1}}{d\mathbf{z}_i} \right|$

- > • Jacobian determinant easy to compute
e.g., choose $df_i^{-1}/d\mathbf{z}_i$ to be a triangular matrix

training: maximizes data log-likelihood

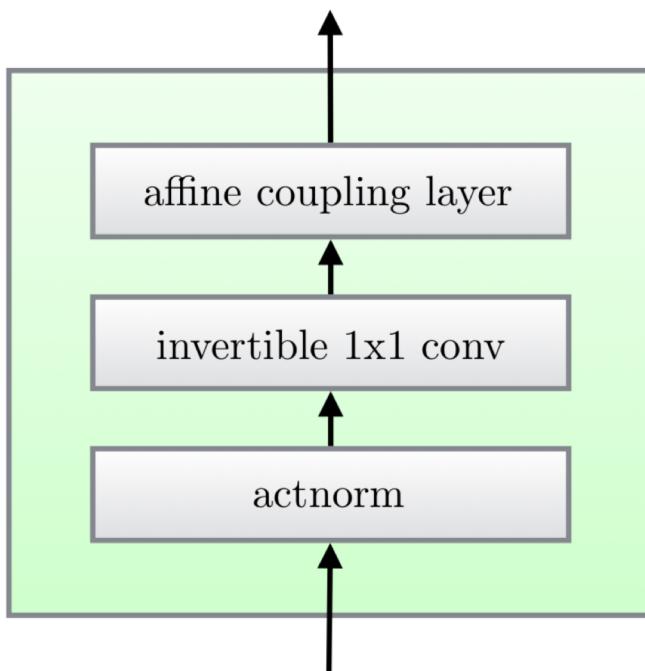
$$\log p(\mathbf{x}) = \log p(\mathbf{z}_0) + \sum_{i=1}^K \log \left| \det \frac{d\mathbf{z}_{i-1}}{d\mathbf{z}_i} \right|$$





GLOW

- [Kingma and Dhariwal., 2018]



One step of flow in the Glow model





Outline

- Generative Adversarial Networks (GANs)
 - GANs Progress
 - Vanilla GAN, Wasserstein GAN, Progressive GAN, BigGAN
- Normalizing Flow (NF)
 - Basic Concepts
 - GLOW
- Integrating Domain Knowledge into Deep Learning





Deep Learning

- Heavily rely on massive labeled data
- Uninterpretable
- Hard to encode human intention and domain knowledge





How Humans Learn

- Learn from **concrete** examples (as DNNs do)
- Learn from **abstract** knowledge (definitions, logic rules, etc) [Minsky 1980; Lake et al., 2015]

Examples:

add	→ added
accept	→ accepted
ignore	→ ignored
end	→ ended
block	→ blocked
love	→ loved

V.S.

Rule:

regular verbs –d/-ed

...





Integrating Domain Knowledge into Deep Learning

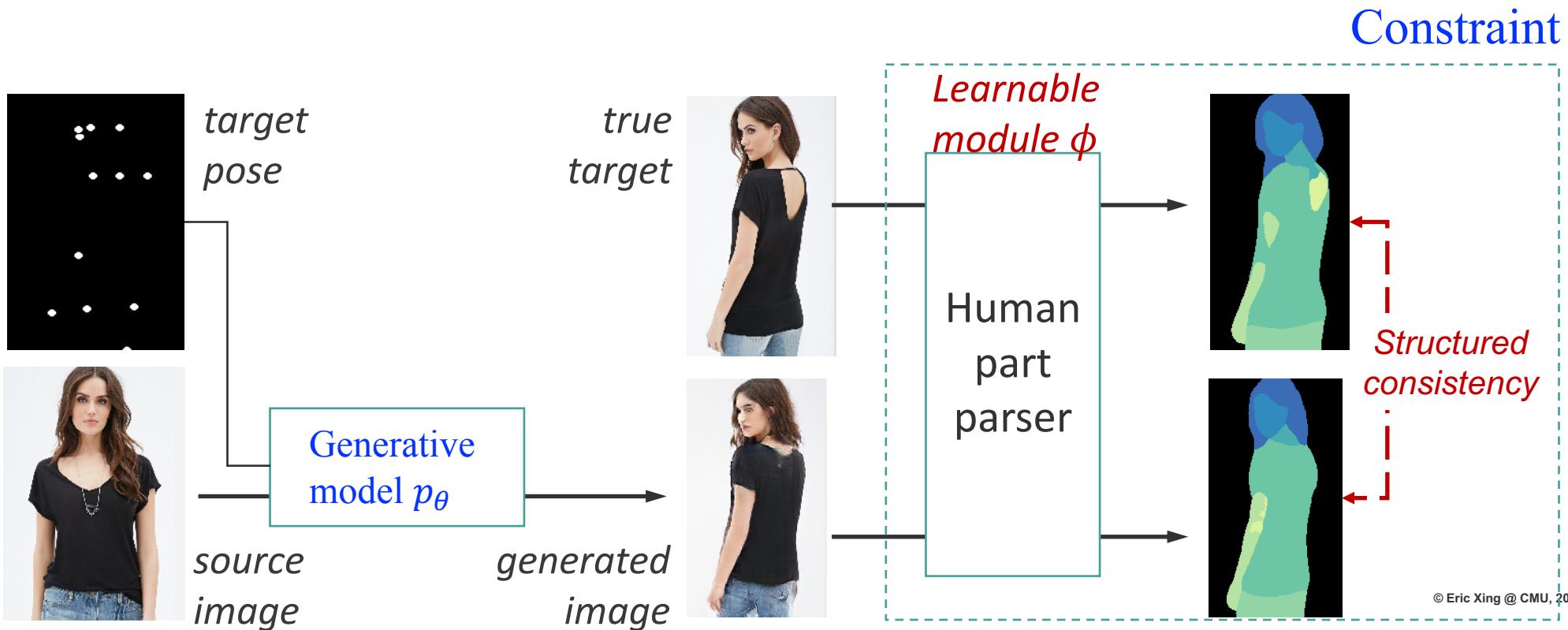
- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
 - Conditional model, $p_{\theta}(\mathbf{x}|\text{inputs})$
 - Generative model, e.g., \mathbf{x} is an image
 - Discriminative model, e.g., \mathbf{x} is a sentence label





Integrating Domain Knowledge into Deep Learning

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f_{\phi}(\mathbf{x}) \in \mathbb{R}$
 - Higher f_{ϕ} value, better \mathbf{x} w.r.t. the knowledge





Integrating Domain Knowledge into Deep Learning

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f_{\phi}(\mathbf{x}) \in \mathbb{R}$
 - Higher f_{ϕ} value, better \mathbf{x} w.r.t. the knowledge
- Sentiment classification
 - “This was a terrific movie, but the director could have done better”
- Logical Rules:
 - Sentence S with structure A-but-B \Rightarrow sentiment of B dominates!





Learning with Constraints

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f_{\phi}(\mathbf{x}) \in \mathbb{R}$
 - Higher f_{ϕ} value, better \mathbf{x} w.r.t. the knowledge
- One way to impose the constraint is to maximize: $\mathbb{E}_{p_{\theta}}[f_{\phi}(\mathbf{x})]$
- Objective:

$$\min_{\theta} \mathcal{L}(\theta) - \alpha \mathbb{E}_{p_{\theta}}[f_{\phi}(\mathbf{x})]$$

Regular objective (e.g.,
cross-entropy loss, etc.)

Regularization:
imposing constraints
(difficult to compute)





Learning with Constraints

- Consider a statistical model $\mathbf{x} \sim p_\theta(\mathbf{x})$
- Consider a constraint function $f_\phi(\mathbf{x}) \in \mathbb{R}$

$$\min_{\theta} \mathcal{L}(\theta) - \alpha \mathbb{E}_{p_\theta} [f_\phi(\mathbf{x})]$$
$$\mathcal{L}(\theta, q) = \text{KL}(q(\mathbf{x}) || p_\theta(\mathbf{x})) - \lambda \mathbb{E}_q [f_\phi(\mathbf{x})]$$

Posterior Regularization
[Ganchev et al., 2010]

- Introduce variational distribution q
 - Impose constraint on q
 - Encourage q to stay close to p
- Objective

$$\min_{\theta, q} \mathcal{L}(\theta) + \alpha \mathcal{L}(\theta, q)$$





Learning with Constraints

$$\min_{\theta, q} \mathcal{L}(\theta) + \alpha \mathcal{L}(\theta, q)$$

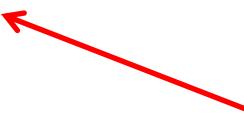
$$\mathcal{L}(\theta, q) = \text{KL}(q(x) \parallel p_\theta(x)) - \lambda \mathbb{E}_q[f_\phi(x)]$$

- EM algorithm for solving the problem

- E-step

$$q^*(x) = p_\theta(x) \exp\{\lambda f_\phi(x)\}/Z$$

- M-step



Higher value -- higher probability
under q – “soft constraint”

$$\min_{\theta} \mathcal{L}(\theta) - \mathbb{E}_{q^*}[\log p_\theta(x)]$$





Logical Rule Constraints

- Consider a supervised learning: $p_{\theta}(y|x)$
- Input-Target space (X, Y)
- First-order logic rules: (r, λ)
 - $r(X, Y) \in [0, 1]$, could be soft
 - λ is the confidence level of the rule
- Given l rules:

- E-step:
$$q^*(y|x) = p_{\theta}(y|x) \exp \left\{ \sum_l \lambda_l r_l(y, x) \right\} / Z$$

- M-step:

$$\min_{\theta} \mathcal{L}(\theta) - \mathbb{E}_{q^*} [\log p_{\theta}(y|x)]$$



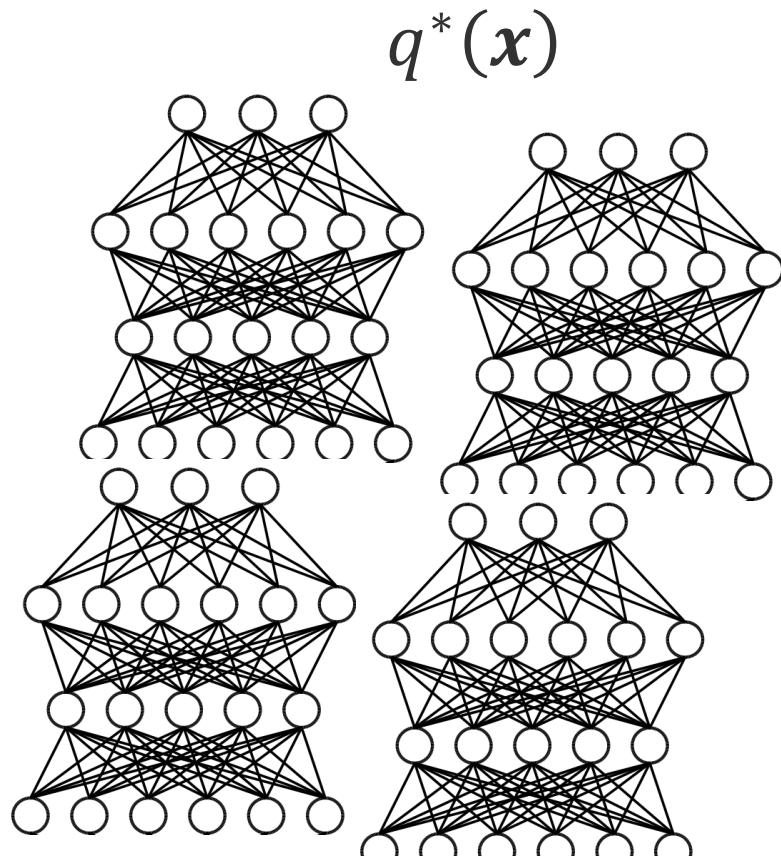
Knowledge distillation [Hinton et al., 2015; Bucilu et al., 2006]

© Eric Xing @ CMU, 2005-2019





Knowledge Distillation

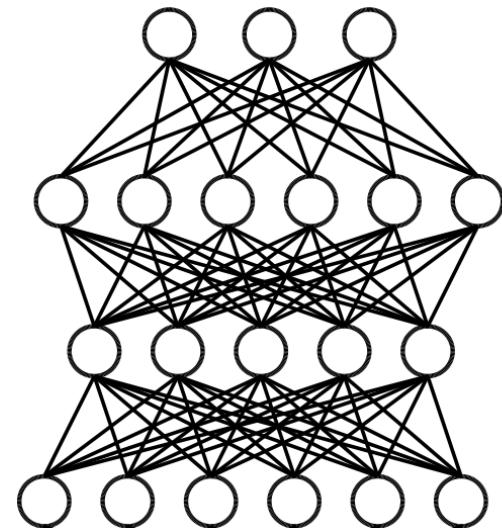


Teacher
(Ensemble)

Match soft predictions of the teacher network and student network



$$p_{\theta}(x)$$



Student





Rule Knowledge Distillation

$$\min_{\theta} \mathcal{L}(\theta) - \mathbb{E}_{q^*} [\log p_\theta(y|x)]$$

- Neural network $p_\theta(y|x)$
 - Train to imitate the outputs of the rule-regularized teacher network
 - At iteration t :

$$\theta^{(t+1)} = \arg \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(y_n, \sigma_\theta(x_n))$$

true hard label soft prediction of $p_\theta(y|x)$

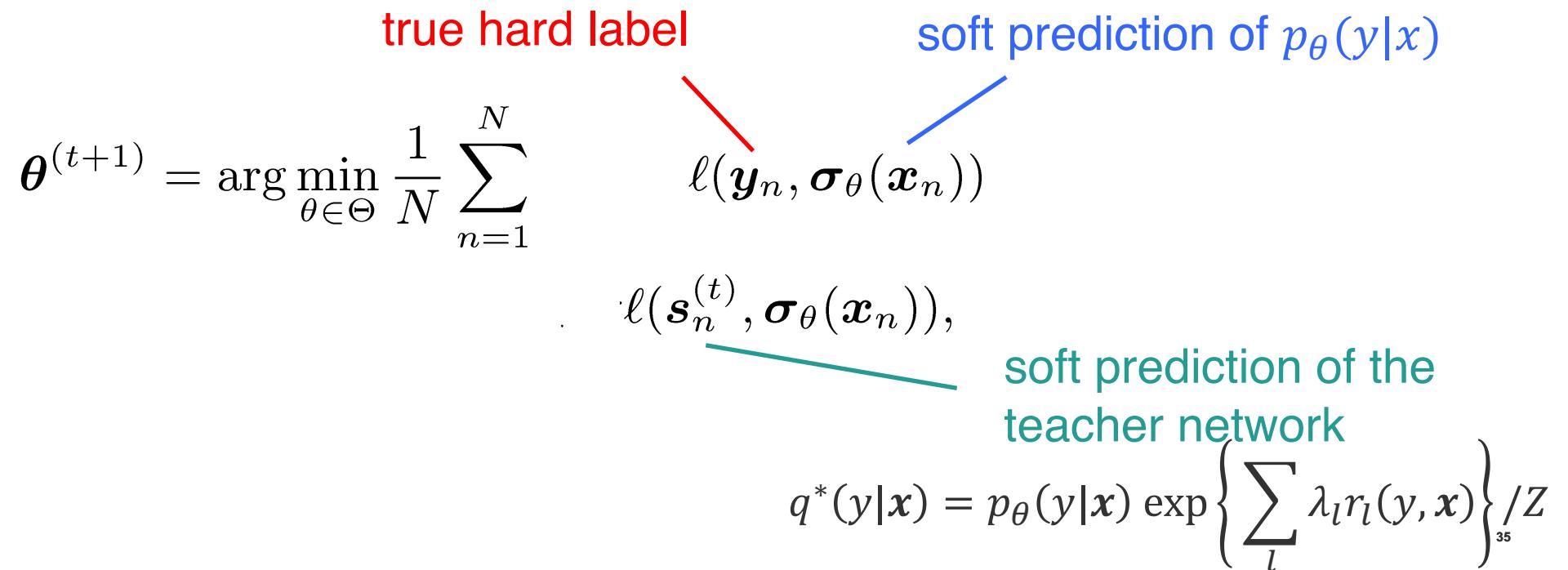




Rule Knowledge Distillation

$$\min_{\theta} \left[\mathcal{L}(\theta) - \mathbb{E}_{q^*} [\log p_\theta(y|x)] \right]$$

- Neural network $p_\theta(y|x)$
- Train to imitate the outputs of the rule-regularized teacher network
- At iteration t :





Rule Knowledge Distillation

$$\min_{\theta} \mathcal{L}(\theta) - \mathbb{E}_{q^*} [\log p_\theta(y|x)]$$

- Neural network $p_\theta(y|x)$
 - Train to imitate the outputs of the rule-regularized teacher network
 - At iteration t :

$$\theta^{(t+1)} = \arg \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N (1 - \pi) \ell(y_n, \sigma_\theta(x_n)) + \pi \ell(s_n^{(t)}, \sigma_\theta(x_n)),$$

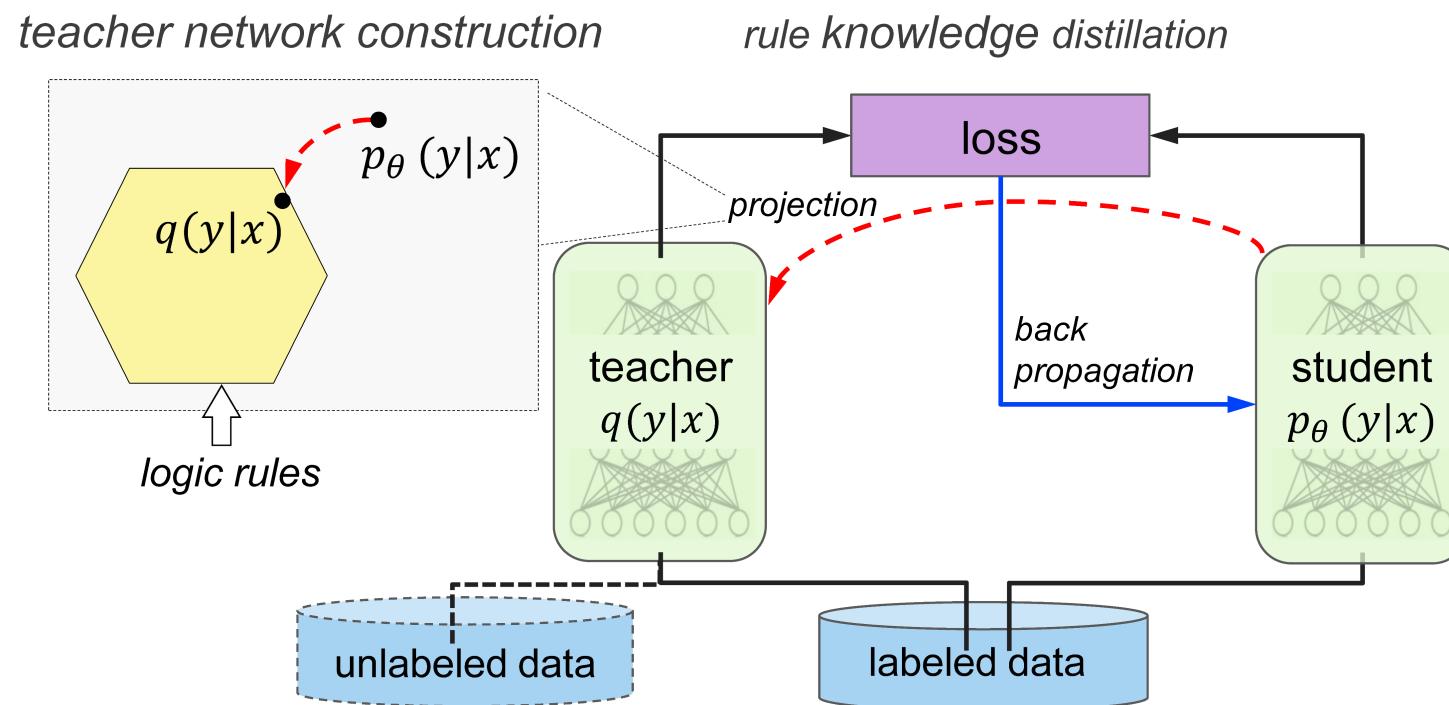
true hard label soft prediction of $p_\theta(y|x)$
balancing parameter soft prediction of the teacher network

$$q^*(y|x) = p_\theta(y|x) \exp \left\{ \sum_l \lambda_l r_l(y, x) \right\} / Z$$



Rule Knowledge Distillation

- Neural network $p_\theta(y|x)$
- At each iteration
 - Construct a teacher network with “soft constraint”
 - Train DNN to emulate the teacher network





Learning Rules / Constraints

$$q^*(y|x) = p_\theta(y|x) \exp \left\{ \sum_l \lambda_l r_l(y, x) \right\} / Z$$

- Learn the confidence value λ_l for each logical rule [Hu et al., 2016b]

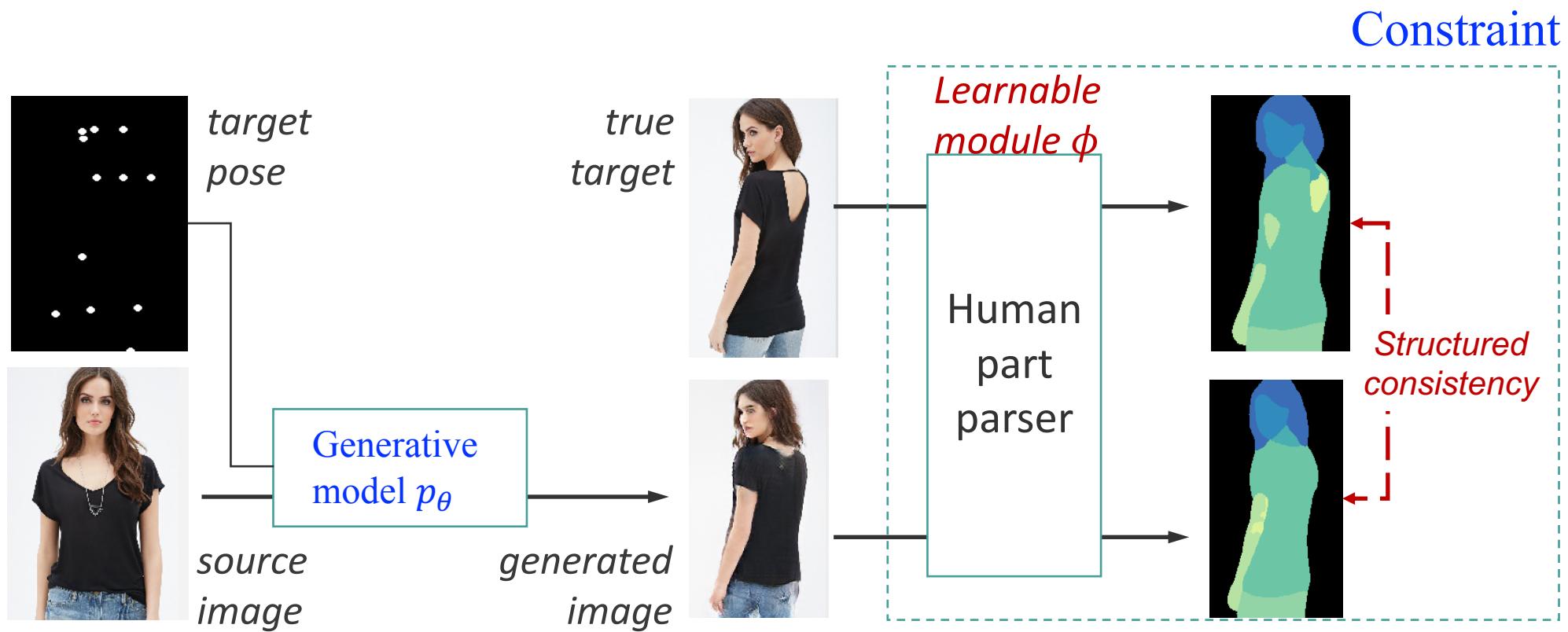
$$q^*(x) = p_\theta(x) \exp \{ \lambda f_\phi(x) \} / Z$$

- More generally, optimize parameters of the constraint $f_\phi(x)$ [Hu et al., 2018]
 - Treat $f_\phi(x)$ as the extrinsic reward function
 - Use MaxEnt Inverse Reinforcement Learning to learn the “reward”





Pose-conditional Human Image Generation





Pose-conditional Human Image Generation



Samples generated by the models

Method	SSIM	Human
1 Ma et al. [38]	0.614	—
2 Pumarola et al. [44]	0.747	—
3 Ma et al. [37]	0.762	—
4 Base model	0.676	0.03
5 With fixed constraint	0.679	0.12
6 With learned constraint	0.727	0.77

Quantitative and Human Evaluation





Template-guided Sentence Generation

- Task: Given a template, generate a complete sentence following the template
- Constraint: force the match between the infilling content of the generated sentence with the true content

template:

“ _____ meant to _____
not to _____.”

true target:

“It was meant to dazzle
not to make sense.”

Constraint

Learnable module ϕ

Infilling content matching

generated:

“It was meant to dazzle
not to make it.”

Generative model p_θ





Template-guided Sentence Generation

Model	Perplexity	Human
1 Base model	30.30	0.19
2 With binary D	30.01	0.20
3 With constraint updated in M-step (Eq.5)	31.27	0.15
4 With learned constraint	28.69	0.24

Samples by the full model are considered as of higher quality in 24% cases.

acting	out of 10 .
the acting	is the acting .
the acting	is also very good .
10	out of 10 .
I will give the movie 7	out of 10 .

Two test examples, including the template, the sample by the base model, and the sample by the constrained model.





Takeaways

- Generative Adversarial Networks (GANs)
 - Wasserstein GAN: new learning objectives
 - Progressive GAN: new training schedule
 - BigGAN: scaling up GAN models
- Normalizing Flow (NF)
 - Chained transformation functions
 - Exact latent inference, density evaluation, sampling
- Integrating Domain Knowledge into Deep Learning
 - Domain knowledge as constraint
 - Learning rules / constraints

