

Ivan Rubachev
HSE, 2019

Semi-supervised Learning with Deep Generative Models

Diederik P. Kingma^{*}, Danilo J. Rezende[†], Shakir Mohamed[†], Max Welling^{*}

^{*}Machine Learning Group, Univ. of Amsterdam, {D.P.Kingma, M.Welling}@uva.nl

[†]Google Deepmind, {danilor, shakir}@google.com

<https://arxiv.org/abs/1406.5298>

Semi-supervised learning

- ▶ Self training
- ▶ Transductive SVMs
- ▶ Graph based methods
- ▶ ...
- ▶ *Generative models*



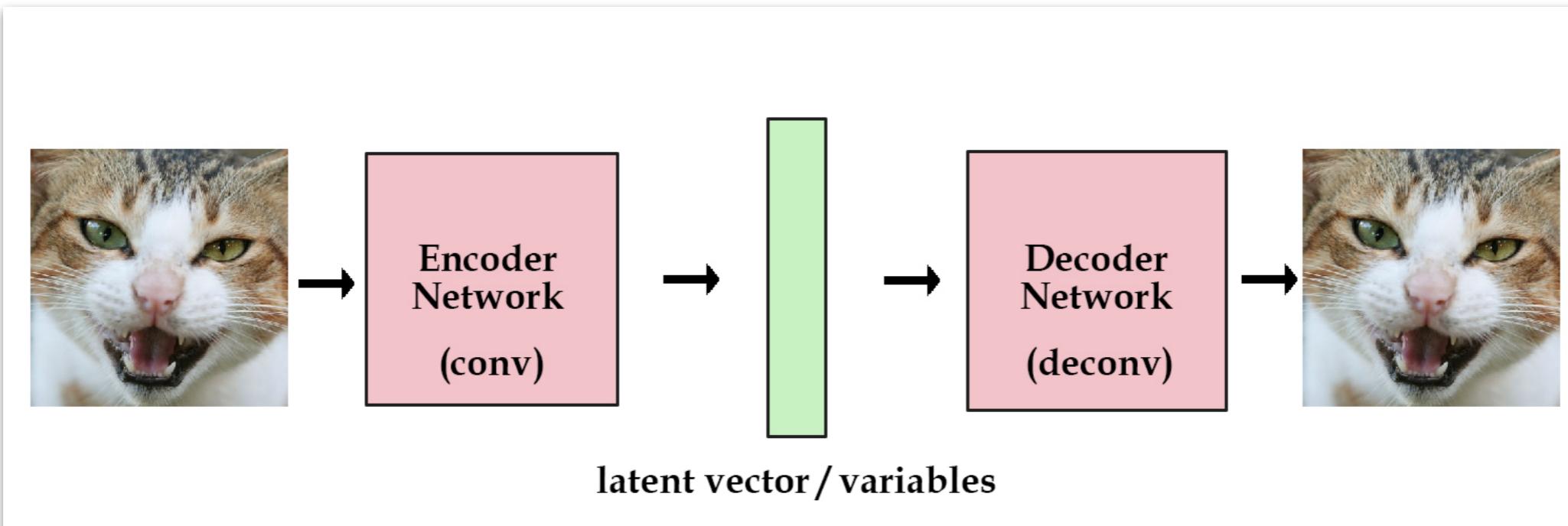
[cake](#)

Remainder: Autoencoders

$$\phi : \mathcal{X} \rightarrow \mathcal{F}$$

$$\psi : \mathcal{F} \rightarrow \mathcal{X}$$

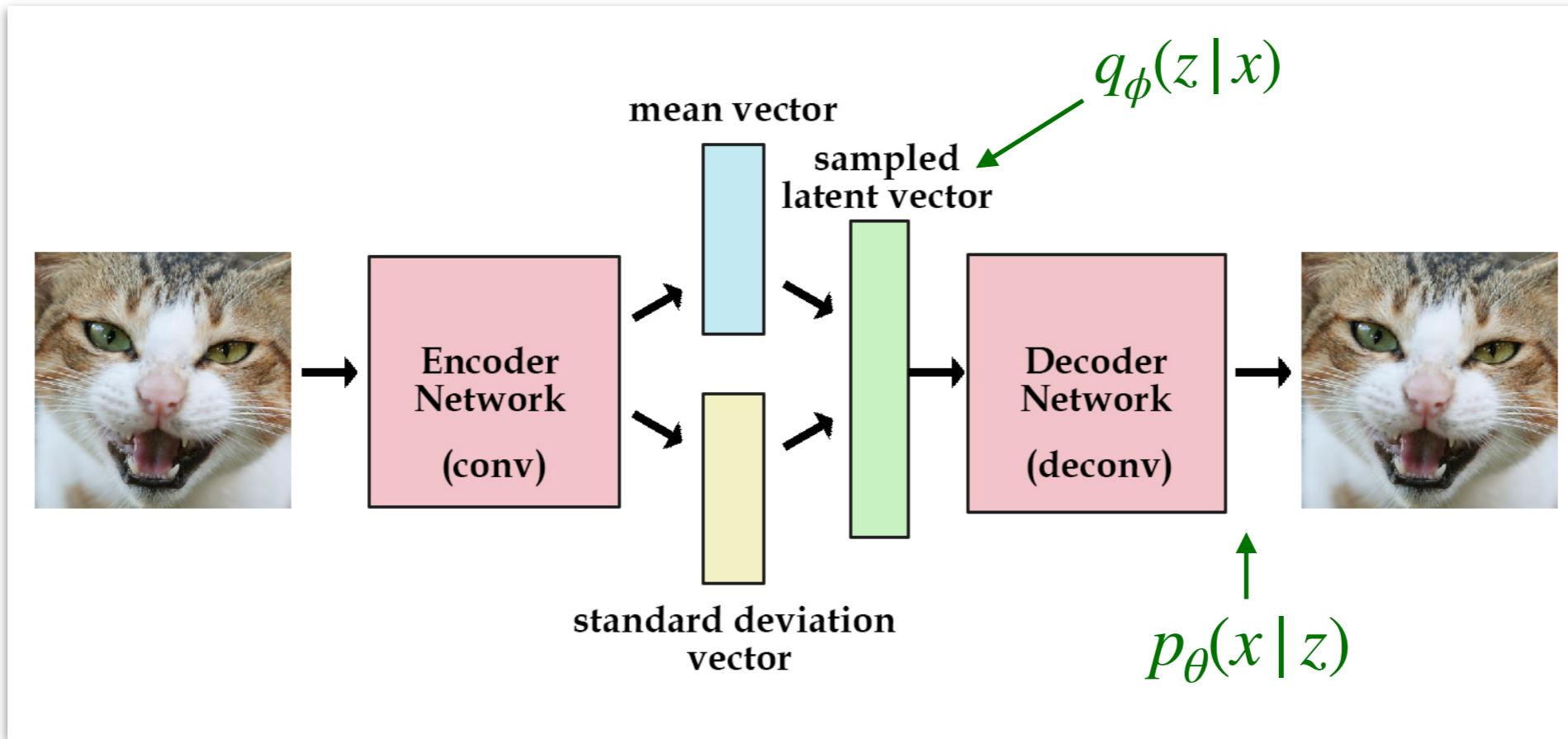
$$\phi, \psi = \arg \min_{\phi, \psi} \|X - (\psi \circ \phi)X\|^2$$



<http://kvfrans.com/variational-autoencoders-explained/>

<https://towardsdatascience.com/generating-images-with-autoencoders-77fd3a8dd368>

Remainder: VAE



Maximize ELBO: $\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - KL[q_\phi(\mathbf{z}|\mathbf{x})\|p_\theta(\mathbf{z})]$

Remainder: VAE

More details

$$X = \{x^{(i)}\}_{i=1}^N - \text{data}$$

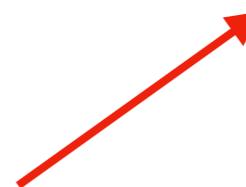
z – latent variables, with prior distribution $p(z)$

x is generated from latent variables $p_\theta(x | z)$

Want to find:

- ▶ Posterior $p_\theta(z | x)$
- ▶ Parameters θ

$$p_\theta(z | x) = \frac{p_\theta(x | z)p_\theta(z)}{\int p_\theta(x | z)p_\theta(z)dz}$$



May be intractable

Remainder: VAE

The variational bound

$q_\phi(z|x)$ — approximation of a posterior $p_\theta(z|x)$

$$\begin{aligned}\log p_\theta(x^{(i)}) &= KL\left(q_\phi(z|x^{(i)}) \| p_\theta(z|x^{(i)})\right) + \mathcal{L}(\theta, \phi; x^{(i)}) \\ &\geq \mathcal{L}(\theta, \phi; x^{(i)}) = \mathbb{E}_{q_\phi(z|x)} \left[-\log q_\phi(z|x) + \log p_\theta(x, z) \right] \\ &= -KL\left(q_\phi(z|x^{(i)}) \| p_\theta(z)\right) + \mathbb{E}_{q_\phi(z|x^{(i)})} \left[\log p_\theta(x^{(i)}|z) \right]\end{aligned}$$

Optimizing the lower bound

(equivalent to minimizing $KL\left(q_\phi(z|x) \| p(z|x)\right)$):

$$\max_{\phi, \theta} \mathcal{L}(\theta; \phi; x)$$

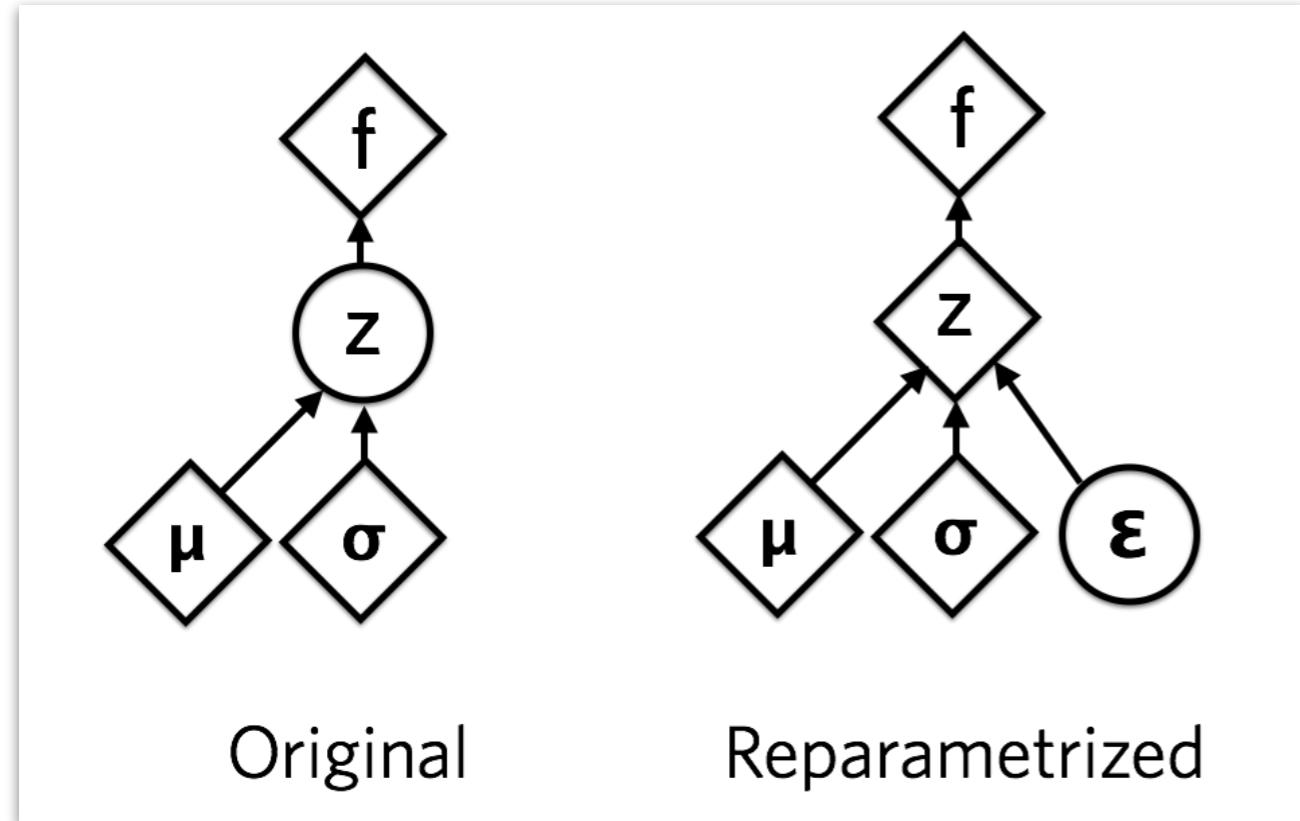
$$\nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z})}[f(\mathbf{z})] = \mathbb{E}_{q_\phi(\mathbf{z})} [f(\mathbf{z}) \nabla_{q_\phi(\mathbf{z})} \log q_\phi(\mathbf{z})] \simeq \frac{1}{L} \sum_{l=1}^L f(\mathbf{z}) \nabla_{q_\phi(\mathbf{z}^{(l)})} \log q_\phi(\mathbf{z}^{(l)})$$

$\frac{\partial}{\partial \{\phi, \theta\}} \mathcal{L}$ — how?

$$\mathbf{z}^{(l)} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})$$

Remainder: VAE

Reparametrization trick



$$p(z) = \mathcal{N}(z | 0, I)$$

$$p_\theta(x | z) = f(x; z; \theta)$$

$$q_\phi(z | x) = \mathcal{N} \left(\mu_\phi(x), \text{diag}(\sigma_\phi^2(x)) \right)$$

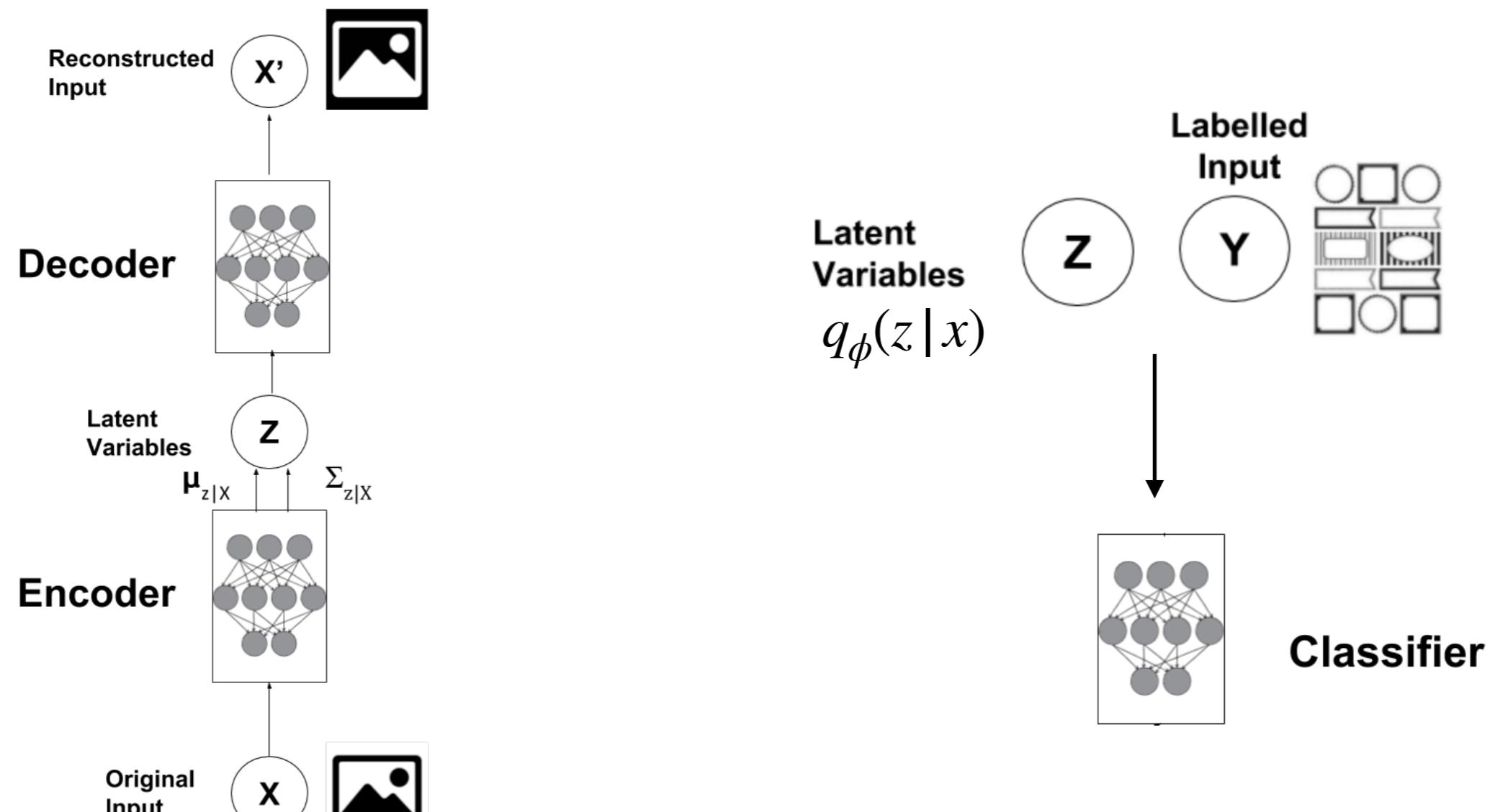
$$z = \mu_\phi(x) + \sigma_\phi(x) \odot \epsilon,$$

$$\epsilon \sim \mathcal{N}(0,1)$$

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [f(\mathbf{z})] = \mathbb{E}_{p(\epsilon)} \left[f(g_\phi(\epsilon, \mathbf{x}^{(i)})) \right] \simeq \boxed{\frac{1}{L} \sum_{l=1}^L f(g_\phi(\epsilon^{(l)}, \mathbf{x}^{(i)}))}$$

Can compute gradients!

Latent-feature discriminative model (m1)



Step 1 (all data)

Step 2 (labeled data)

Generative semi-supervised model (m2)

$$p(y) = \text{Cat}(y | \pi)$$

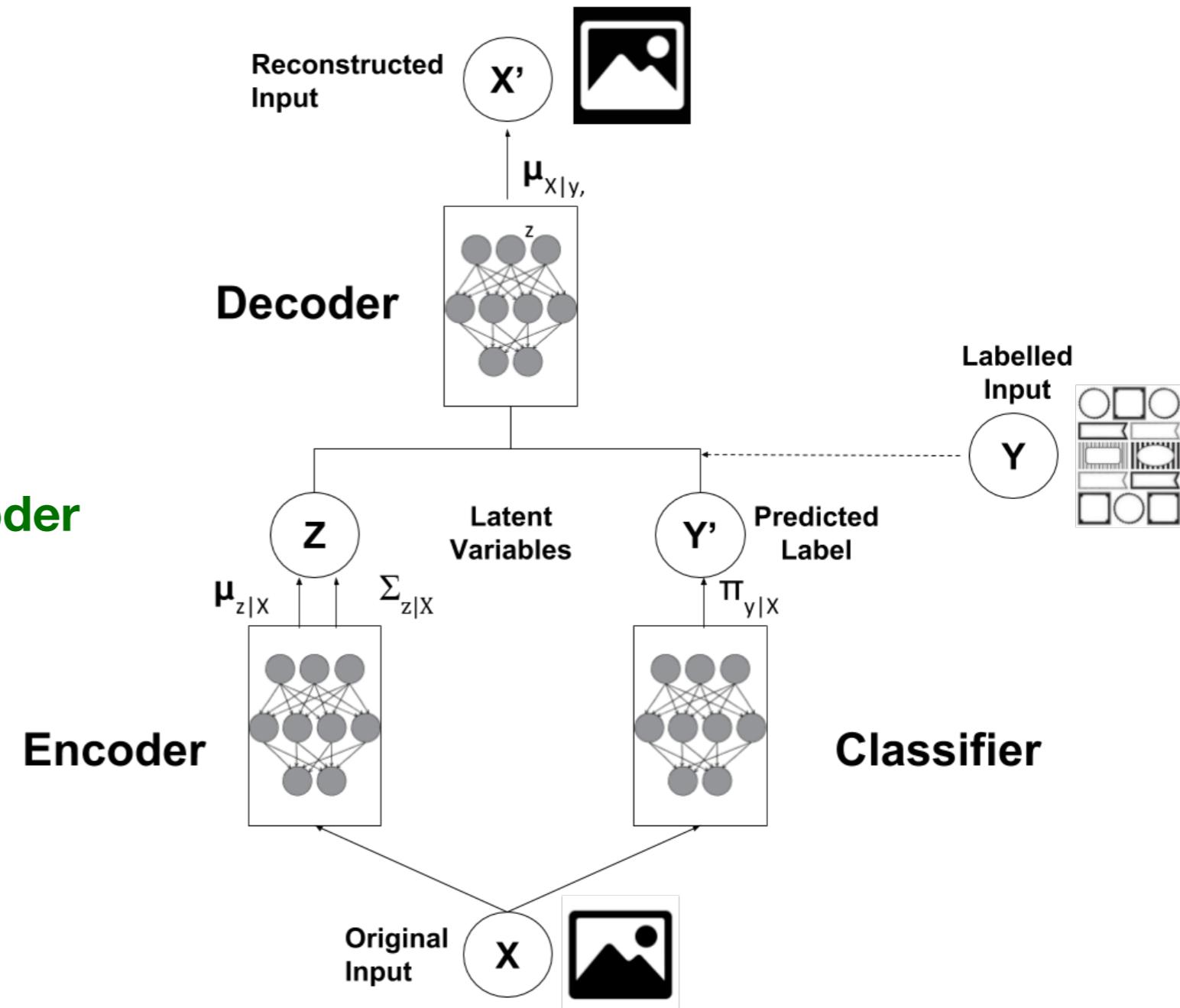
$$p(z) = \mathcal{N}(z | 0, I)$$

$p_\theta(x | y, z) = f(x; y; z; \theta)$ **Decoder**

$$q_\phi(z, y | x) = q_\phi(z | x)q_\phi(y | x)$$

$q_\phi(z | x) = \mathcal{N}(\mu_\phi(x), \text{diag}(\sigma_\phi^2(x)))$ **Encoder**

$q_\phi(y | x) = \text{Cat}(y | \pi_\phi(x))$ **Classifier**



Generative semi-supervised model (m2)

Unlabeled data:

$$\begin{aligned}\log p(x) &\geq \mathbb{E}_{q_\phi(y,z|x)}[\log p_\theta(x|y,z) + \log p_\theta(y) + \log p(z) - \log q_\phi(y,z|x)] \\ &= \sum_y q_\phi(y|x) \mathcal{L}(x,y) + \mathcal{H}(q_\phi(y|x)) = \mathcal{U}(x)\end{aligned}$$

where $\mathcal{L}(x,y) = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|y,z) + \log p_\theta(y)] - KL(q_\phi(z|x)||p_\theta(z))$

Labeled data:

$$\log p(x,y) \geq \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|y,z) + \log p_\theta(y) + \log p_\theta(z) - \log q_\phi(z|x)] = \mathcal{L}(x,y)$$

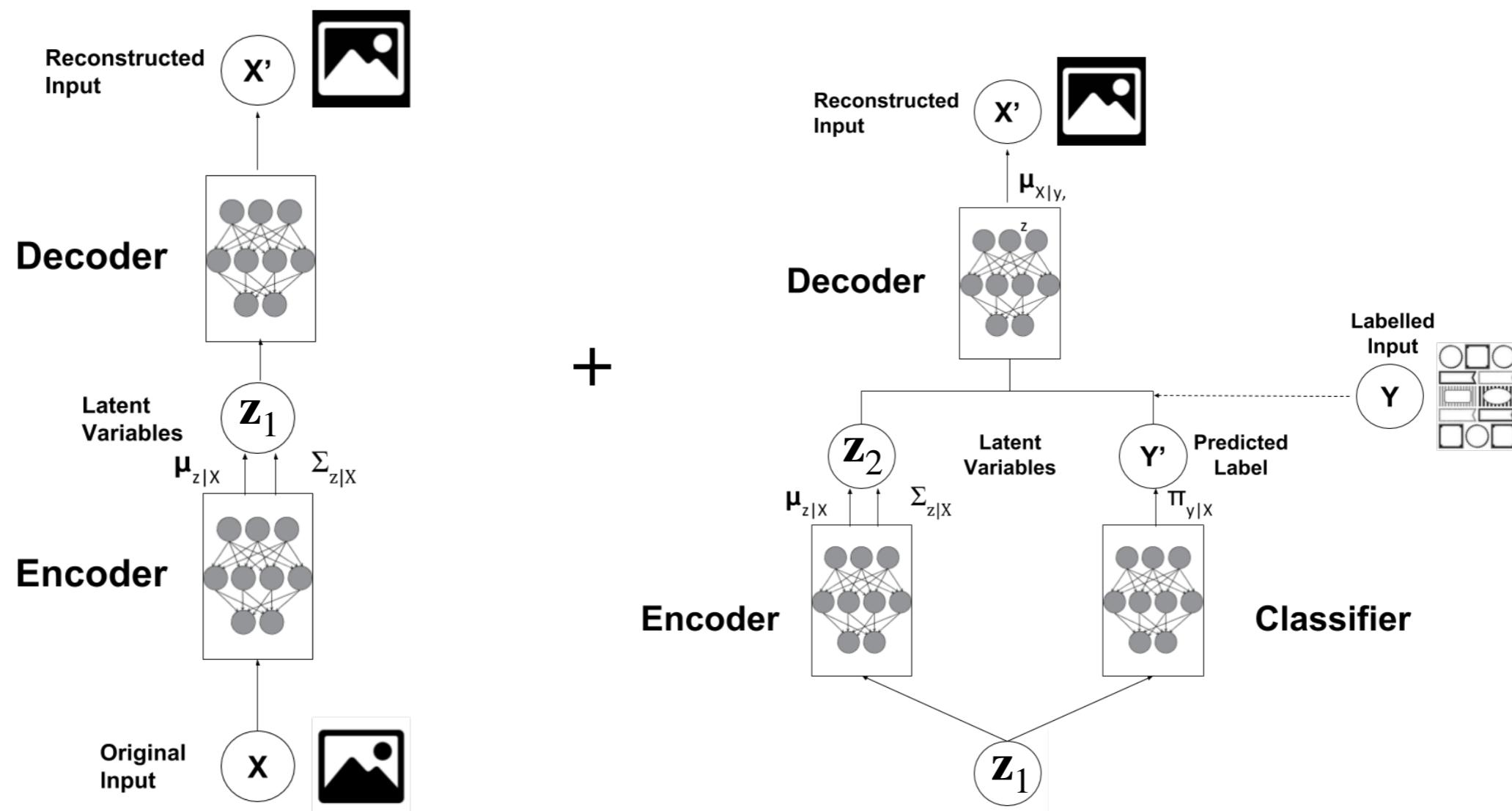
LOSS:

$$\mathcal{J} = - \sum_{x,y \in \text{labeled}} \mathcal{L}(x,y) - \sum_{x \in \text{unlabeled}} \mathcal{U}(x)$$

$$\mathcal{J}^\alpha = \mathcal{J} + \alpha \cdot \mathbb{E}_{p_l(x,y)} \left[-\log q_\phi(y|x) \right]$$

Stacked m1+m2

$$p_{\theta}(x, y, z_1, z_2) = p(y)p(z_2)p_{\theta}(z_1 | y, z_2)p_{\theta}(x | z_1)$$



Error rate on MNIST

- ▶ 50,000 training examples
- ▶ N labeled examples

Table 1: Benchmark results of semi-supervised classification on MNIST with few labels.

N	NN	CNN	TSVM	CAE	MTC	AtlasRBF	M1+TSVM	M2	M1+M2
100	25.81	22.98	16.81	13.47	12.03	8.10 (\pm 0.95)	11.82 (\pm 0.25)	11.97 (\pm 1.71)	3.33 (\pm 0.14)
600	11.44	7.68	6.16	6.3	5.13	–	5.72 (\pm 0.049)	4.94 (\pm 0.13)	2.59 (\pm 0.05)
1000	10.7	6.45	5.38	4.77	3.64	3.68 (\pm 0.12)	4.24 (\pm 0.07)	3.60 (\pm 0.56)	2.40 (\pm 0.02)
3000	6.04	3.35	3.45	3.22	2.57	–	3.49 (\pm 0.04)	3.92 (\pm 0.63)	2.18 (\pm 0.04)

Error rate on SVHN

- ▶ 70,000 training examples
- ▶ 10 classes



Table 2: Semi-supervised classification on the SVHN dataset with 1000 labels.

KNN	TSVM	M1+KNN	M1+TSVM	M1+M2
77.93 (± 0.08)	66.55 (± 0.10)	65.63 (± 0.15)	54.33 (± 0.11)	36.02 (± 0.10)

Error rate on NORB

- ▶ 194,400 training examples
- ▶ 5 classes



Table 3: Semi-supervised classification on the NORB dataset with 1000 labels.

KNN	TSVM	M1+KNN	M1+TSVM
78.71 (\pm 0.02)	26.00 (\pm 0.06)	65.39 (\pm 0.09)	18.79 (\pm 0.05)

Conditional generation

2 2 2 2 2 2 2 2 2 2	3 3 3 3 3 3 3 3 3 3	4 4 4 4 4 4 4 4 4 4
2 2 2 2 2 2 2 2 2 2	3 3 3 3 3 3 3 3 3 3	4 4 4 4 4 4 4 4 4 4
2 2 2 2 2 2 2 2 2 2	3 3 3 3 3 3 3 3 3 3	4 4 4 4 4 4 4 4 4 4
2 2 2 2 2 2 2 2 2 2	3 3 3 3 3 3 3 3 3 3	4 4 4 4 4 4 4 4 4 4
2 2 2 2 2 2 2 2 2 2	3 3 3 3 3 3 3 3 3 3	4 4 4 4 4 4 4 4 4 4
2 2 2 2 2 2 2 2 2 2	3 3 3 3 3 3 3 3 3 3	4 4 4 4 4 4 4 4 4 4
2 2 2 2 2 2 2 2 2 2	3 3 3 3 3 3 3 3 3 3	4 4 4 4 4 4 4 4 4 4
2 2 2 2 2 2 2 2 2 2	3 3 3 3 3 3 3 3 3 3	4 4 4 4 4 4 4 4 4 4
2 2 2 2 2 2 2 2 2 2	3 3 3 3 3 3 3 3 3 3	4 4 4 4 4 4 4 4 4 4
2 2 2 2 2 2 2 2 2 2	3 3 3 3 3 3 3 3 3 3	4 4 4 4 4 4 4 4 4 4

Different styles generated by varying latent variables

Conditional generation

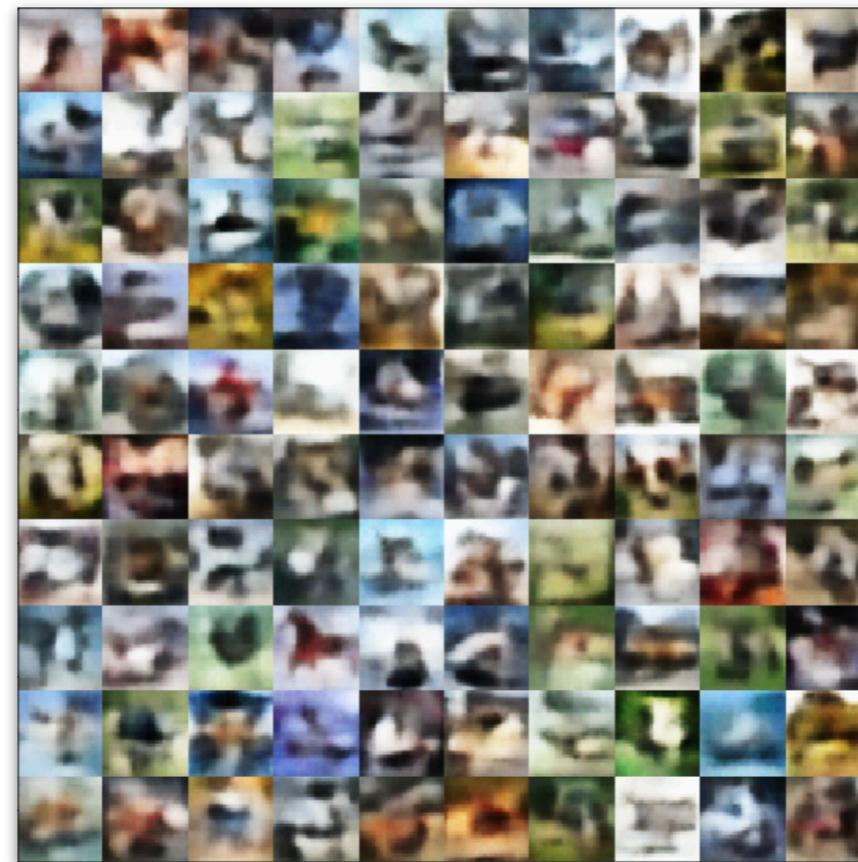
4 0 1 2 3 4 5 6 7 8 9
9 0 1 2 3 4 5 6 7 8 9
5 0 1 2 3 4 5 6 7 8 9
4 0 1 2 3 4 5 6 7 8 9
2 0 1 2 3 4 5 6 7 8 9
7 0 1 2 3 4 5 6 7 8 9
5 0 1 2 3 4 5 6 7 8 9
1 0 1 2 3 4 5 6 7 8 9
7 0 1 2 3 4 5 6 7 8 9
1 0 1 2 3 4 5 6 7 8 9

40	11	2	3	4	5	6	17	18	19	10
15	1	2	3	4	5	6	7	8	9	10
36	10	20	30	40	50	60	70	80	90	00
7	11	2	3	4	5	6	7	8	9	10
13	11	12	13	14	15	16	17	18	19	10
30	11	2	3	4	5	6	7	8	9	10
61	11	21	31	41	51	61	71	81	91	01
20	10	20	30	40	50	60	70	80	90	00
28	21	22	23	24	25	26	27	28	29	20
22	21	22	23	24	25	26	27	28	29	20

Analogies generated by varying classes

Conclusion

- ▶ Almost as efficient as other Neural Network approaches
- ▶ Linear scale with # of classes
- ▶ Not SOTA today (2019)



https://github.com/bjlkeng/sandbox/blob/master/notebooks/vae-semi_supervised_learning/vae-m2-fit.ipynb