

Lecture 8: Autoencoders

Learning representation – no supervision

- How do we know what is relevant and what is not?
- **Main idea of unsupervised representation learning:** a new representation should be “smaller” but allow to recover the original one.
- *Why smaller?* (because we need to eliminate irrelevant, because we need something small)

Recap: PCA

We now want to pick the d optimal components

$$\{v_1, \dots, v_d\} \in \mathbb{R}^N$$

Let $\tilde{X} = [\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_M]$

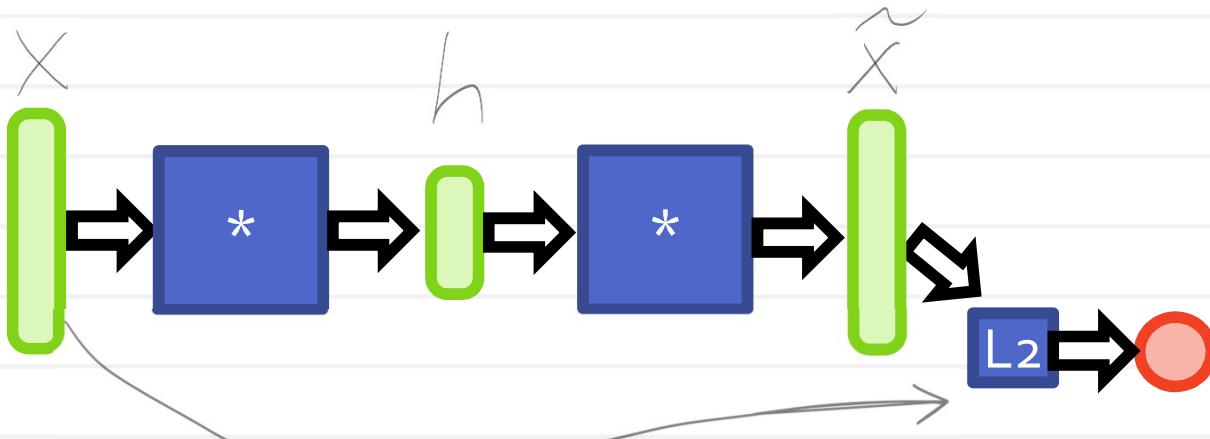
be the reconstruction of the original vectors.

We want to minimize the reconstruction error:

$$\sum_{i=1}^M \|x_i - \tilde{x}_i\|_2^2 \rightarrow \min$$

$$\|\mathbf{X} - \tilde{\mathbf{X}}\|_F^2 \rightarrow \min$$

Neural network version (autoencoder)



PCA:

$$h = U^T X$$

$$\tilde{X} = Uh$$

$$\min_U \|X - \tilde{X}\|_2^2$$

$$U$$

$$\text{s.t. } \langle u^i, u^j \rangle = S_{ij}$$

Autoencoder:

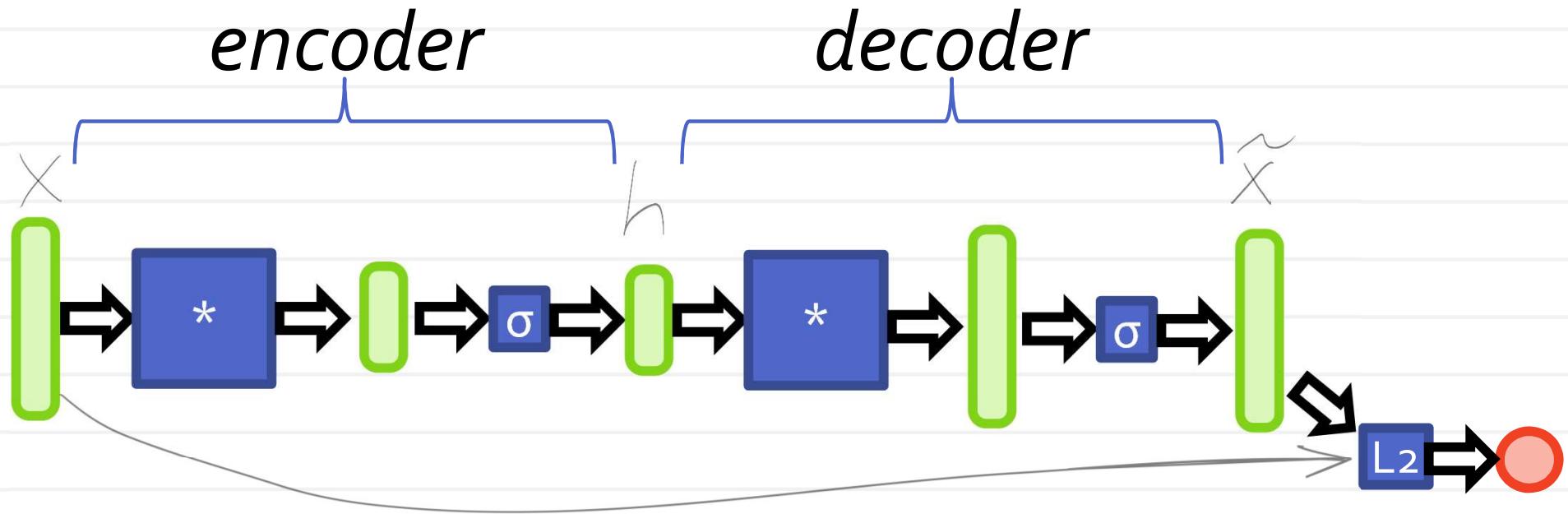
$$h = W X$$

$$\tilde{X} = W^T h$$

$$\min_W \|X - \tilde{X}\|_2^2$$

[Bengio et al. 07]

Autoencoder: adding non-linearity

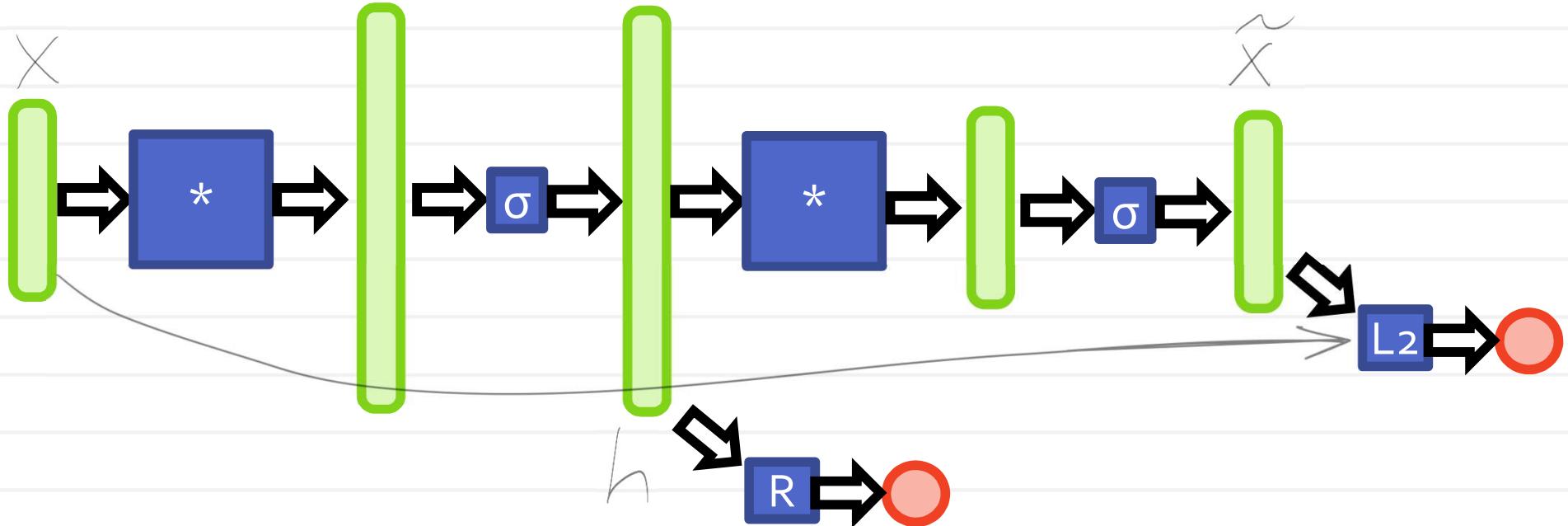


$$h = w x$$

$$\tilde{x} = w' h$$

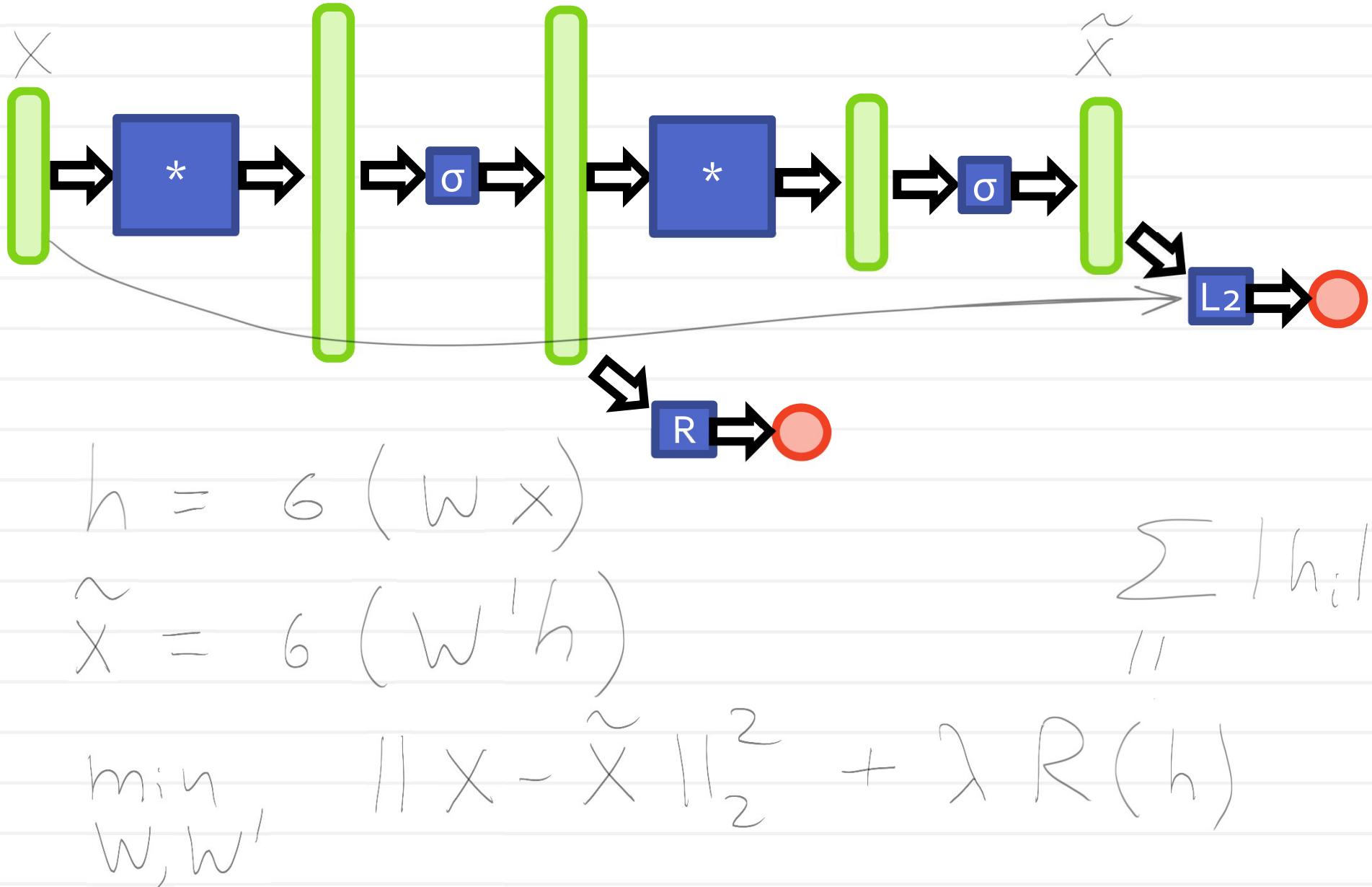
$$\min_{w, w'} \|x - \tilde{x}\|_2^2$$

Learning bigger/richer representation



- Naïve approach will learn an identity function
- Solution 1: regularize hidden representation

Sparse autoencoders



Sparse autoencoders

$$\min_{W W'} \|X - \tilde{X}\|_2^2 + \lambda R(h)$$

estimating from e.g. minibatch

$$R(h) = \sum_i KL(\text{Binomial}(p_i) || \text{Binomial}(\hat{p}_i))$$

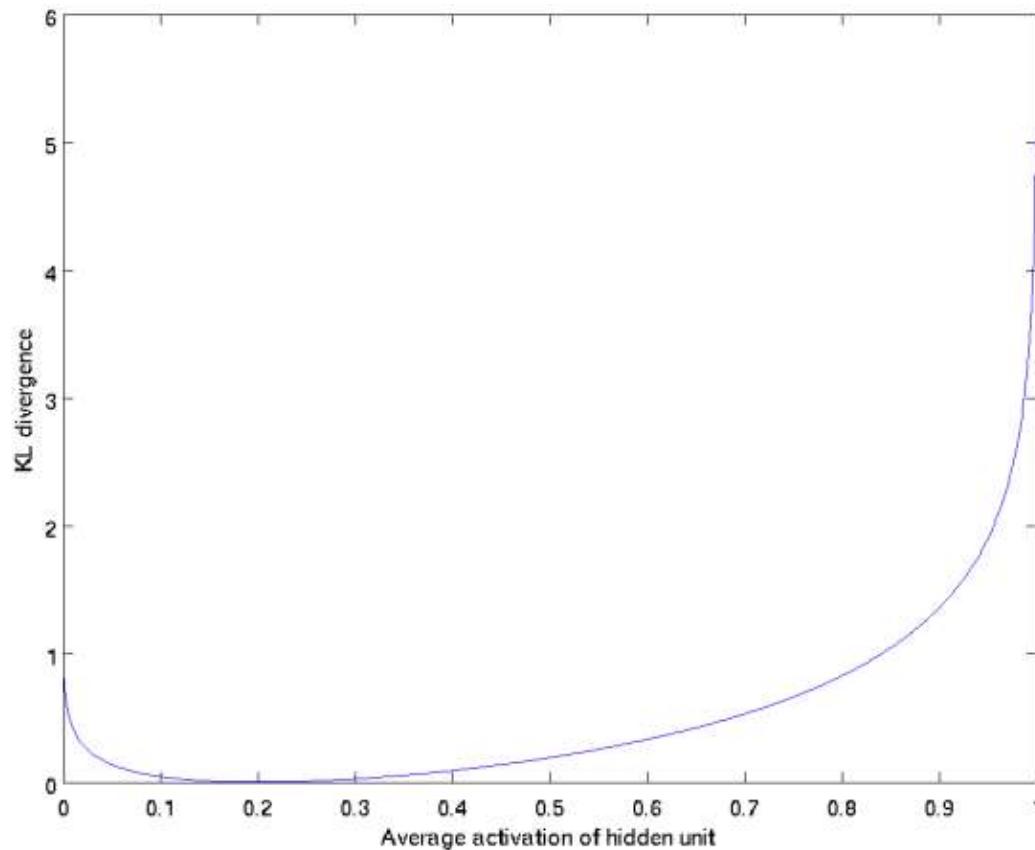
$$= \sum_i p_i \log \frac{p_i}{\hat{p}_i} + (1-p_i) \log \frac{1-p_i}{1-\hat{p}_i}$$

$$\frac{\partial R}{\partial \hat{p}_i} = -\frac{p_i}{\hat{p}_i} + \frac{1-p_i}{1-\hat{p}_i}$$

[Ng 2011]

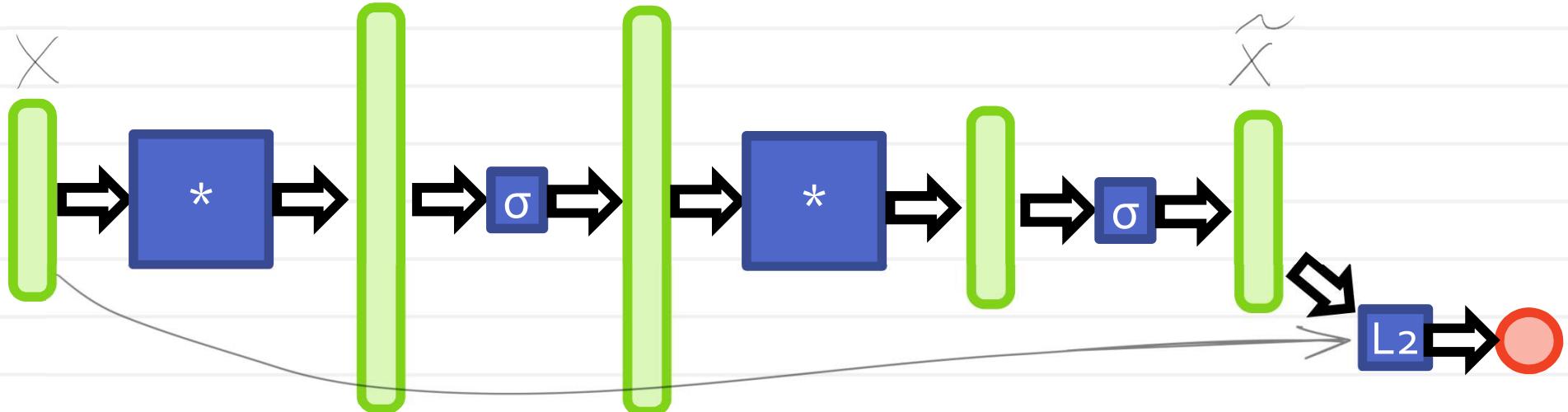
Sparse autoencoders

$$R(h) = \sum_i KL(\text{Binomial}(p_i) \parallel \text{Binomial}(\hat{p}_i))$$
$$= \sum_i p_i \log \frac{p_i}{\hat{p}_i} + (1-p_i) \log \frac{1-p_i}{1-\hat{p}_i}$$



[Ng 2011]

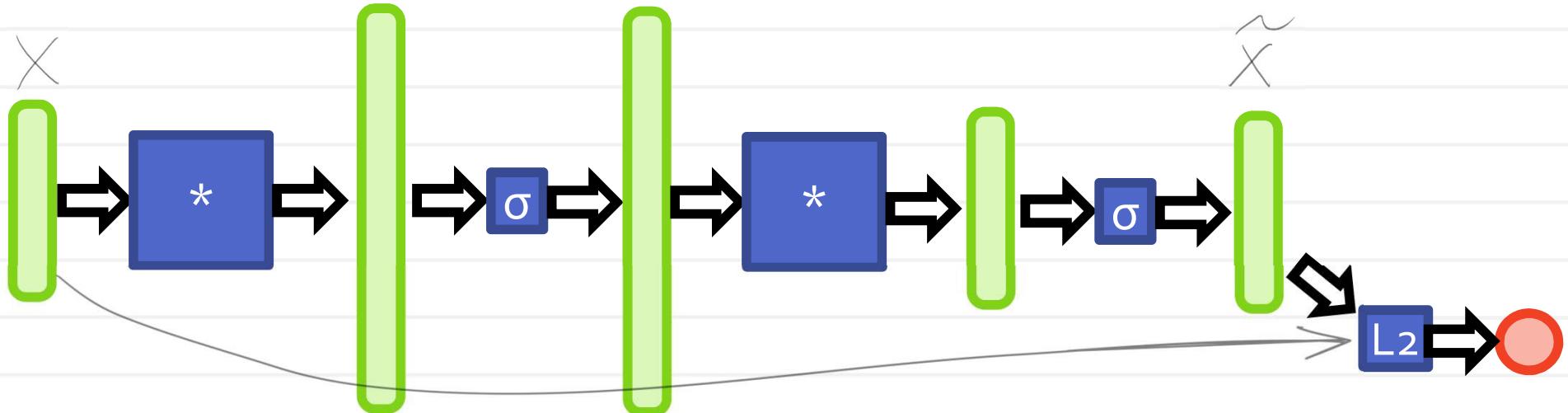
Learning bigger representation



- Naïve approach will learn an identity function
- Solution 2: corrupt the input and reconstruct uncorrupted

[Vincent et al. 2010]

Denoising autoencoder



$$h = w(x + h)$$

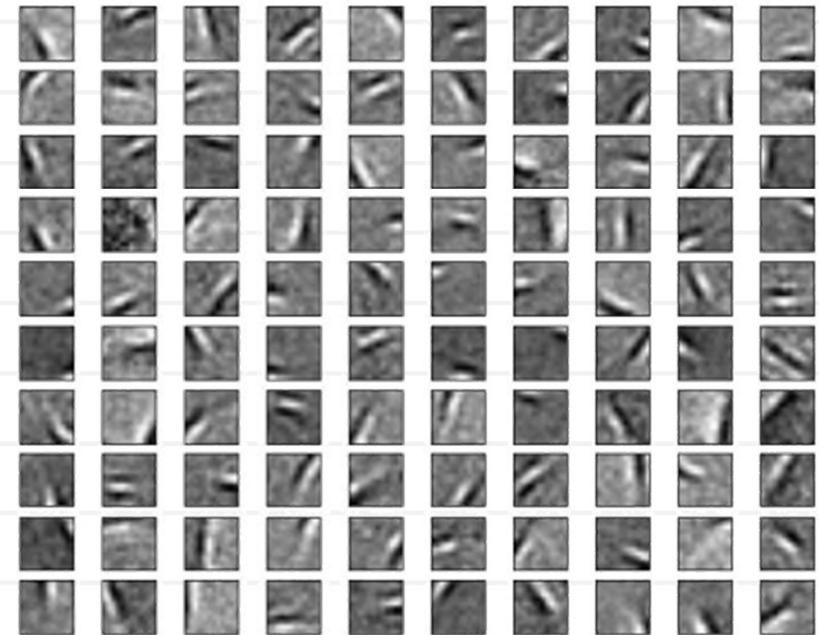
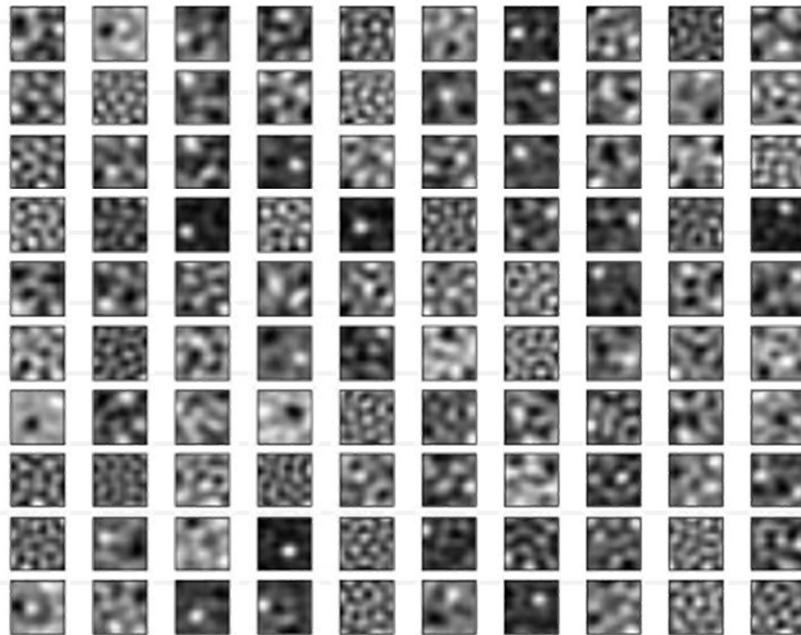
$$\tilde{x} = w^T h$$

$$\min_{w, w'} \sum_i \|x_i - \tilde{x}_i\|_2^2$$

[Vincent et al. 2010]

Denoising autoencoder

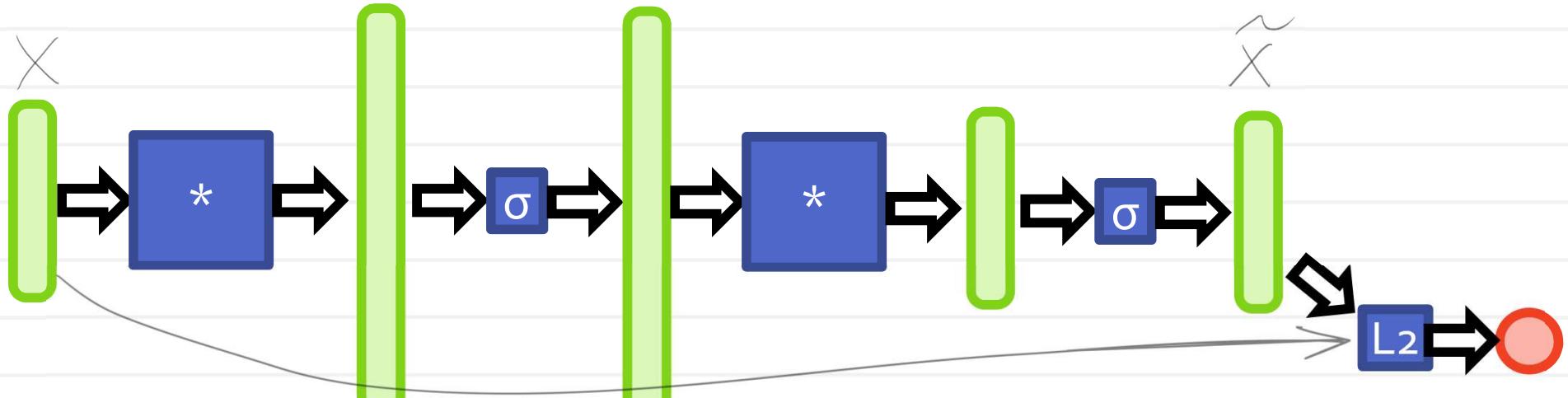
12x12 patches, 200 hidden units



Autoencoder with
weight-decay

Denoising
autoencoder
(Gaussian noise)
[Vincent et al. 2010]

Contractive autoencoders



$$h = 6 (w \times)$$

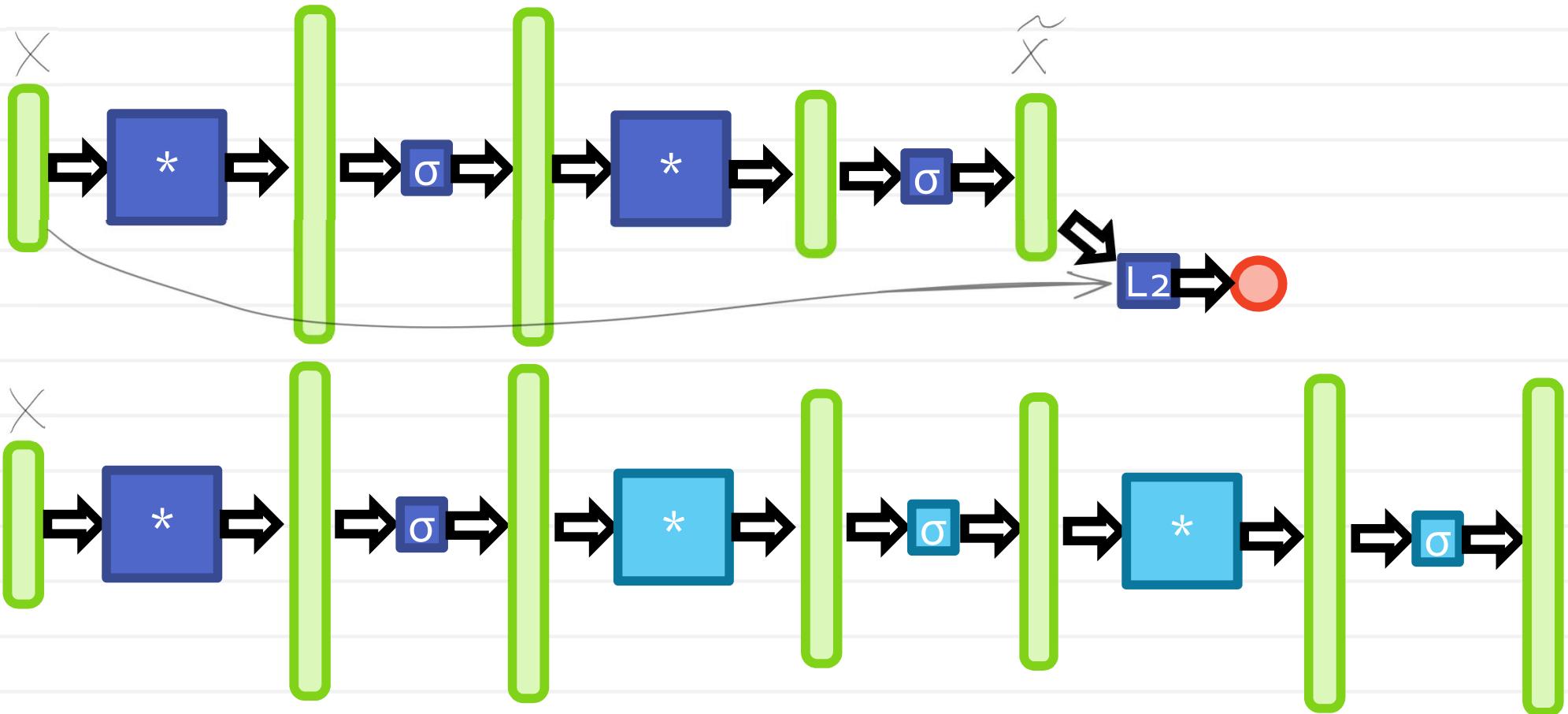
$$\tilde{X} = 6 (w^T h)$$

$$\min_{w, w'} \|X - \tilde{X}\|_2^2 + \lambda \left\| \frac{\partial h_i(x)}{\partial x_j} \right\|_F^2$$

$$\frac{\partial h_i}{\partial x_j} = (h_i(1-h_i)) w_{ij}$$

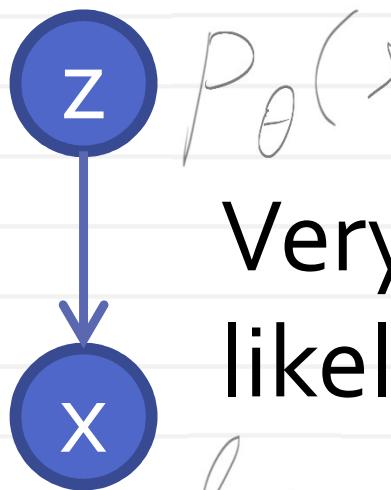
[Rifai et al. 2011]

Training deep autoencoders



- Historically: layerwise pretraining
- Most modern implementations: end-to-end
- Commonly including convolutional and up-convolutional layers

Maximum likelihood



$$P_{\theta}(x | z)$$

Very natural idea: train θ by maximum likelihood:

$$\log P_{\theta}(x) = \log \int_z P_{\theta}(x, z) dz \rightarrow \max$$

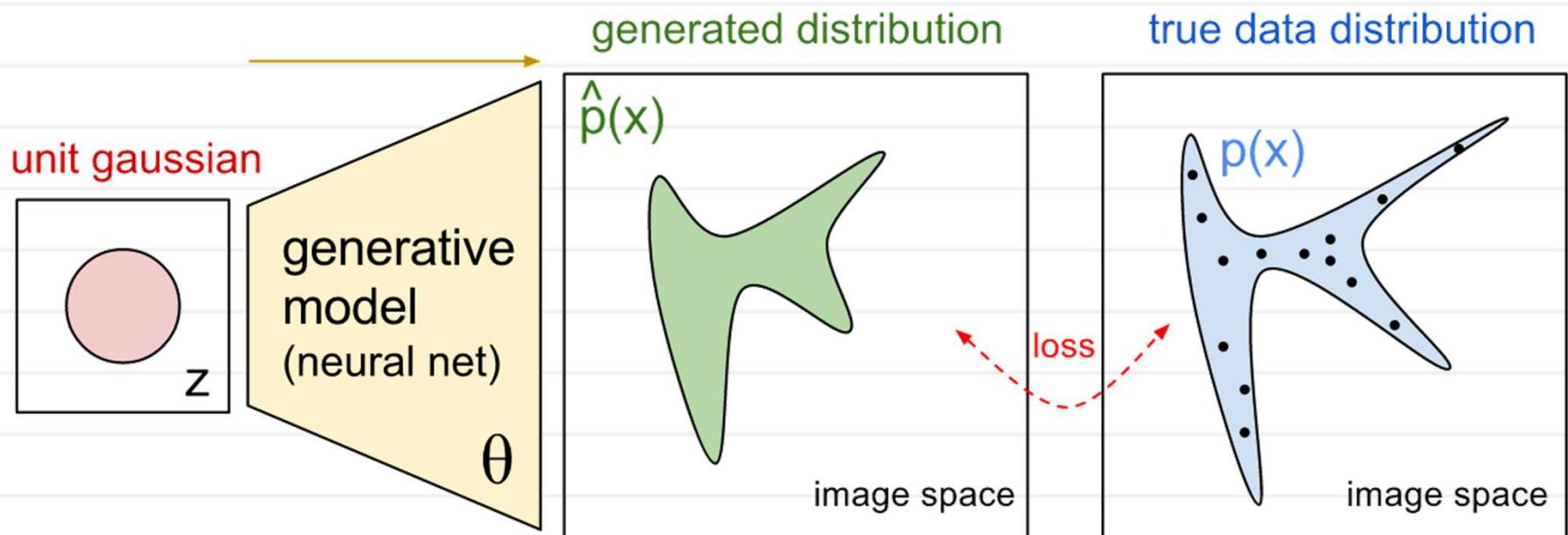
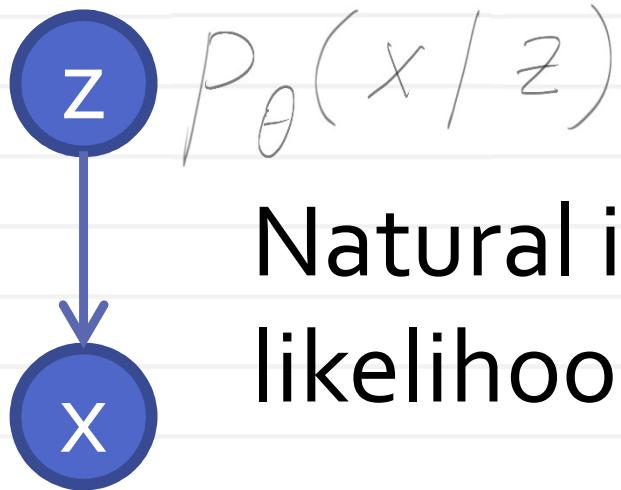


Image credit: OpenAI blog

Maximum likelihood via VAE



Natural idea: train θ by maximum likelihood.

$$\log p(x) = \log \int_z p(x, z) dz \rightarrow \max$$

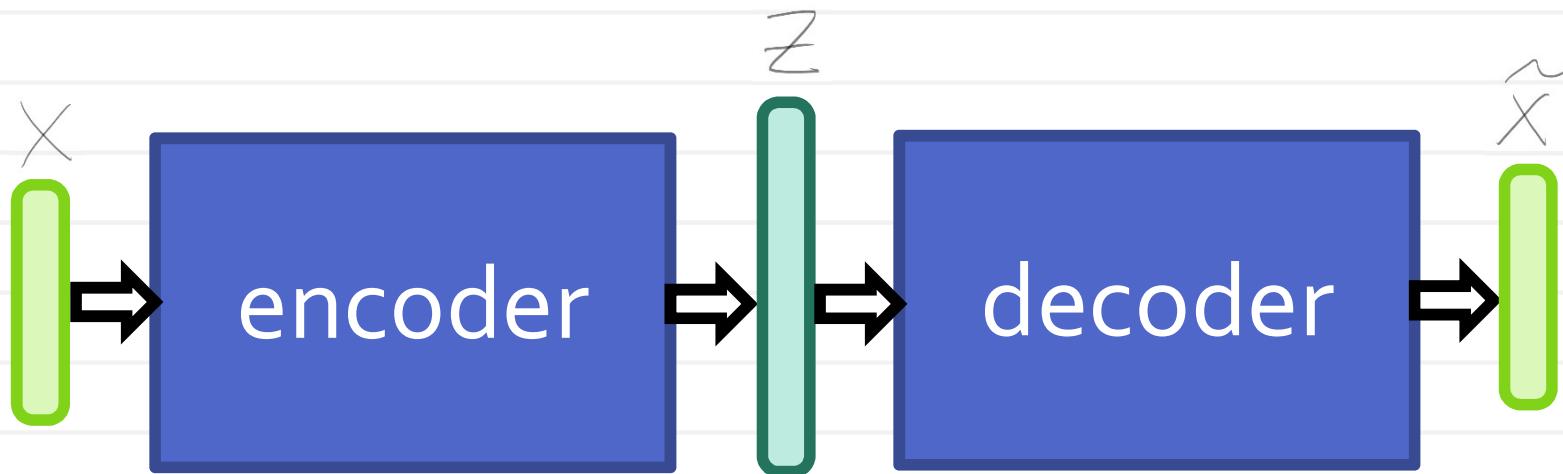
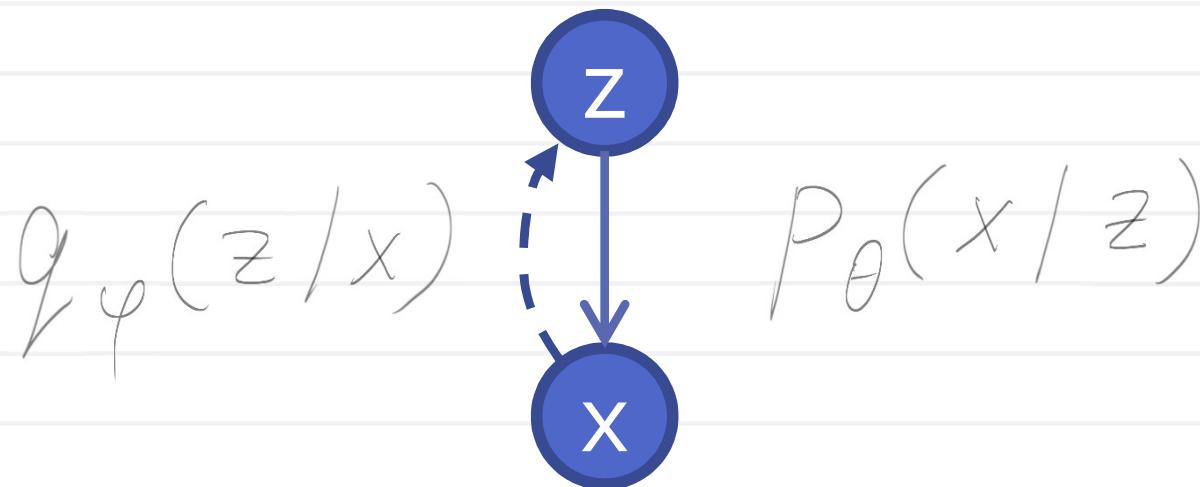
The integral is intractable.

VAE solution: introduce an *encoder*

$$q_{\theta_0}(z|x) \sim p(z|x)$$

Derive an approximation to ML objective

Variational autoencoder as a graphical model

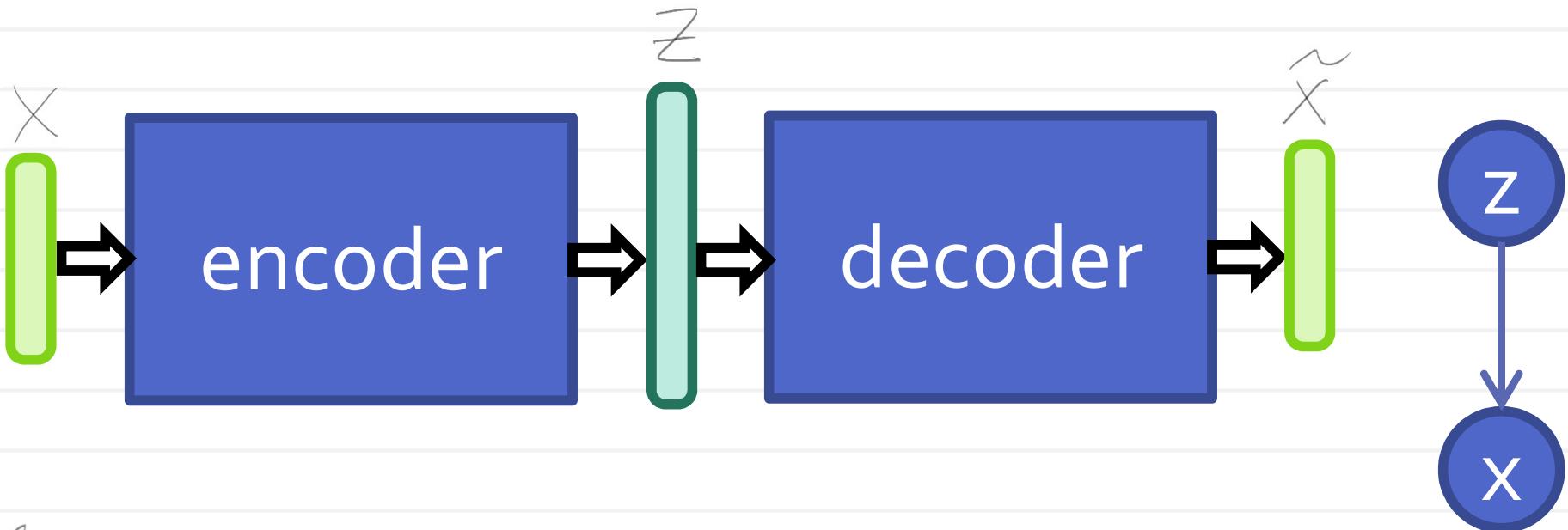


[Kingma & Welling 14]

Variational lower bound

$$\begin{aligned}\log p(x) &= \log \int_z p(x, z) dz = \\ &= \log \int_z p(x, z) \frac{q(z|x)}{q(z|x)} dz = \\ &= \log E_{q(z|x)} \frac{p(x, z)}{q(z|x)} = \\ &\quad \log E_{q(z|x)} \frac{p(x|z) p(z)}{q(z|x)} \geq \\ &\geq E_{q(z|x)} \log p(x|z) + E_{q(z|x)} \log \frac{p(z)}{q(z|x)} = \\ &= E_{q(z|x)} \log p(x|z) - KL(q(z|x) || p(z))\end{aligned}$$

Variational lower-bound



$$\log p(x) \geq$$

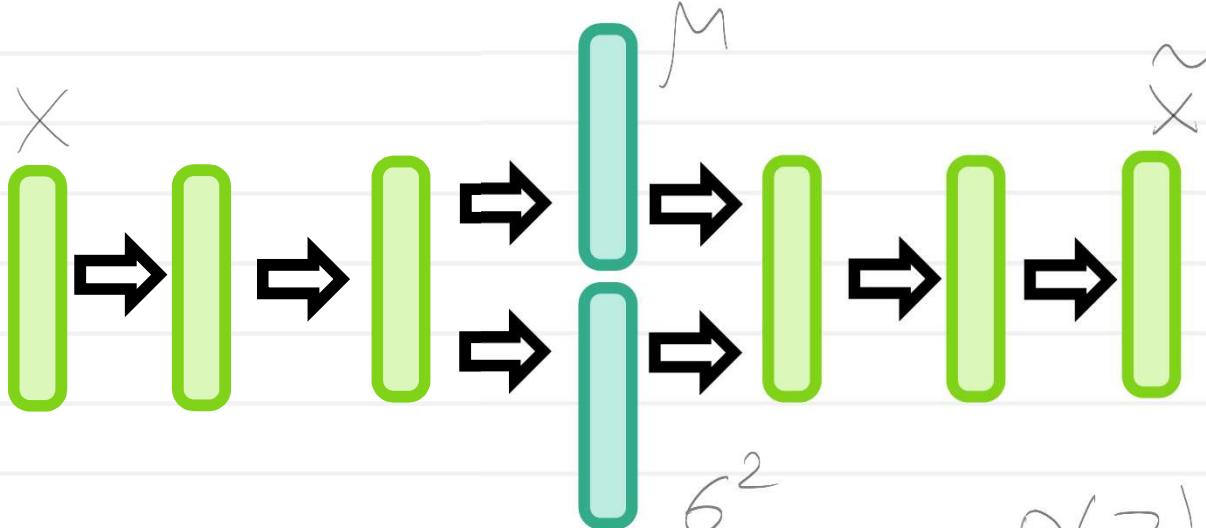
$$-KL(g_{\varphi}(z|x) || p(z)) +$$

regularization

$$E_{g_{\varphi}(z|x)} [\log p_{\theta}(x|z)]$$

[Kingma & Welling 14] ~"denoising auto-encoder"

Minimizing regularization term



$$\theta^2$$

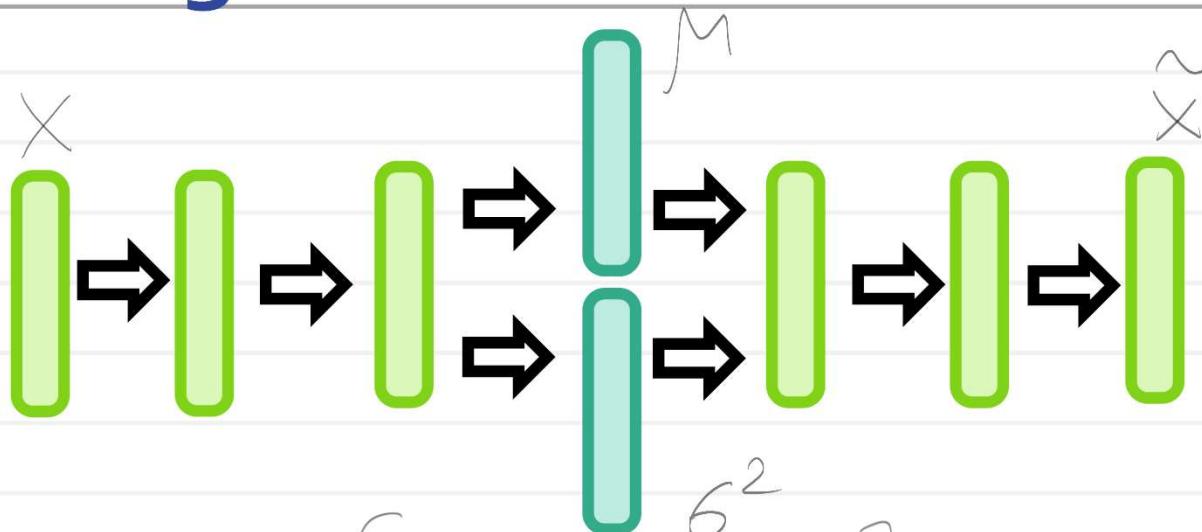
$$p(z) = \prod N(\theta, 1)$$

$$KL(g_{\phi}(z/x) || p(z)) =$$

$$= \frac{1}{2} \sum \left(\mu_j^2 + \theta_j^2 - 1 - \log \theta_j^2 \right)$$

[Kingma & Welling 14]

Optimizing reconstruction term



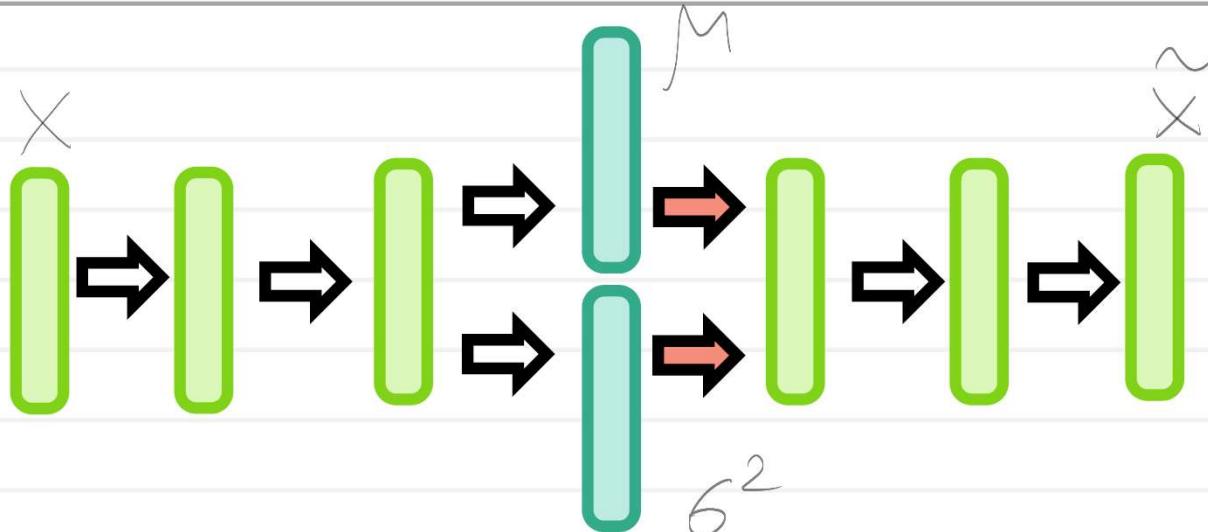
$$E_{\mathcal{G}_\theta(z|x)} \left[\log p_\theta(x|z) \right] \xrightarrow{6^2} \max$$

$$p_\theta(x|z) \propto \exp(-\|x - \hat{x}\|^2)$$

$$E_{\mathcal{G}_\theta(z|x)} \| d_\theta(z) - x \|_2^2 \rightarrow \min$$

[Kingma & Welling 14]

Reparameterization trick



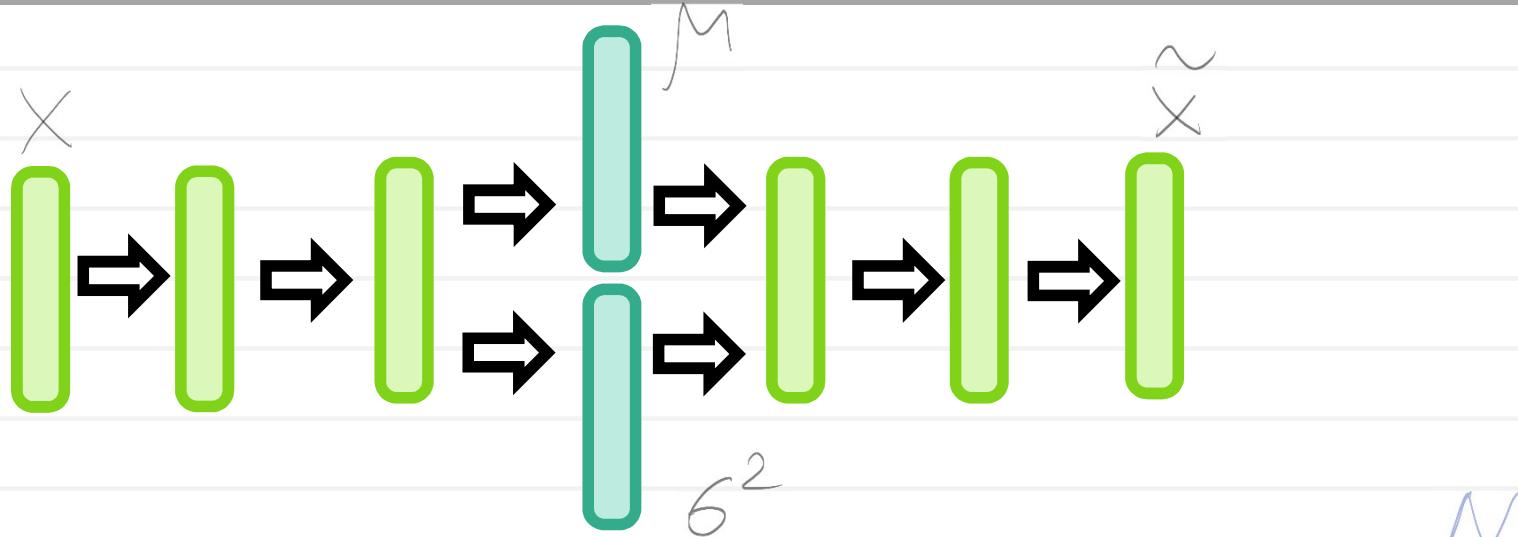
$$E_{\mathcal{D}_\theta}(z|x) \left\| d_\theta(z) - x \right\|_2^2 \rightarrow \min$$

$$\tilde{z} = \mu + \sigma \odot \mathcal{E} \sim \mathcal{N}(\cdot; \cdot)$$

$$L = \left\| d_\theta(\tilde{z}) - x \right\|_2^2 \rightarrow \min$$

[Kingma & Welling 14]

Reparameterization trick



$$\tilde{z} = \mu + \sigma^2 \circ \epsilon \sim \mathcal{N}(\mu; \sigma^2)$$

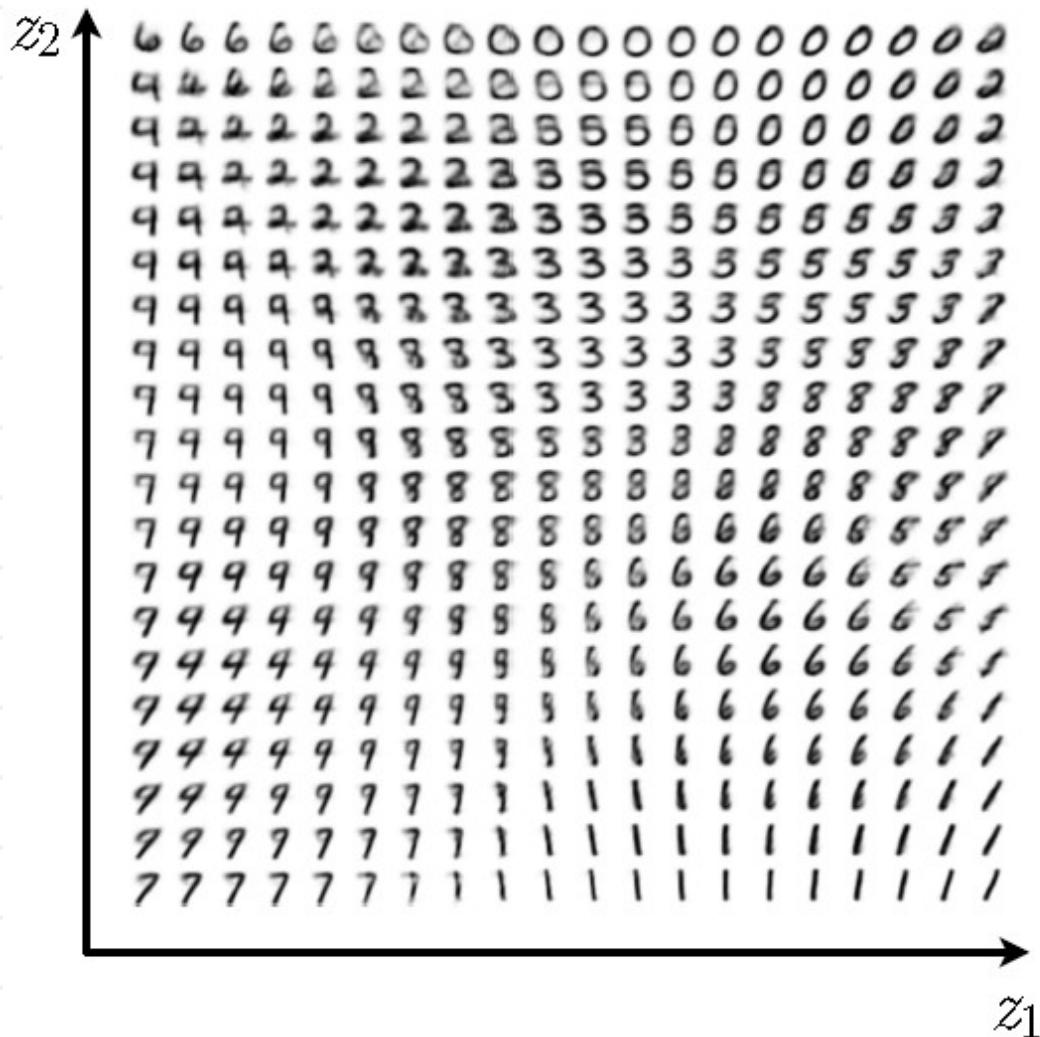
$$\frac{dL}{d\mu} = \frac{dL}{d\tilde{z}}$$

$$\frac{dL}{d\sigma^2} = \frac{dL}{d\tilde{z}} \circ \epsilon$$

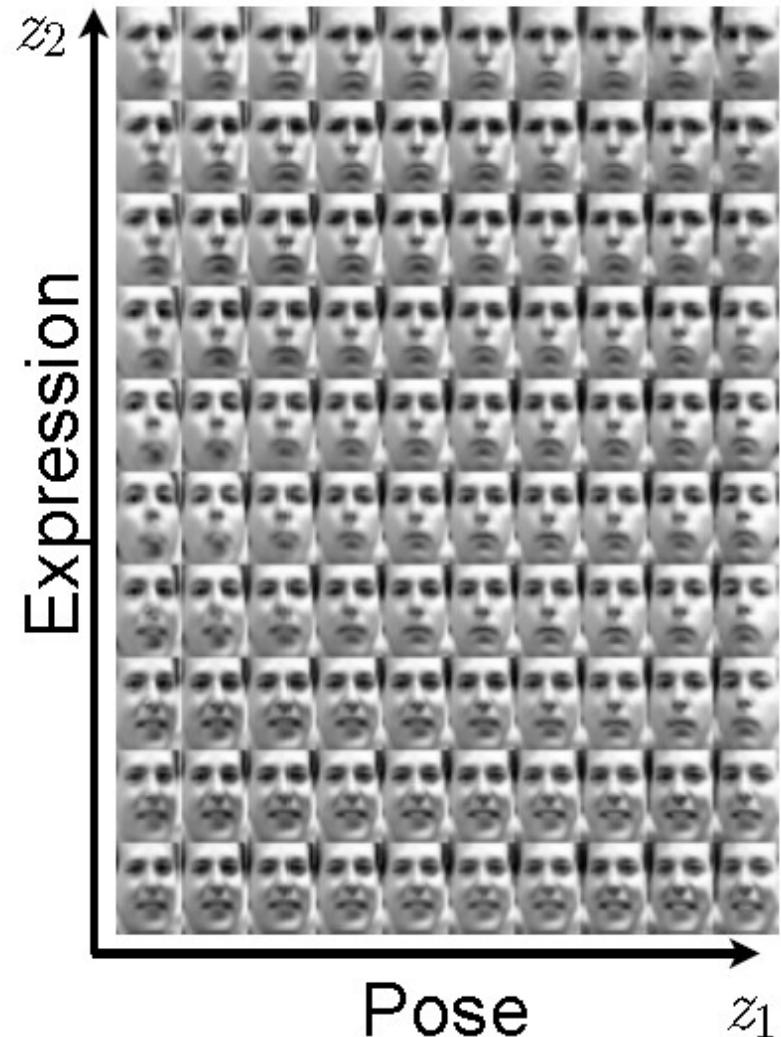
[Kingma & Welling 14]

VAE-learned manifolds

MNIST:



Frey Face dataset:



[Kingma & Welling 14]

VAE-learned manifolds

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

(a) 2-D latent space

(b) 5-D latent space

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

(c) 10-D latent space

(d) 20-D latent space

[Kingma & Welling 14]

Collapsing components

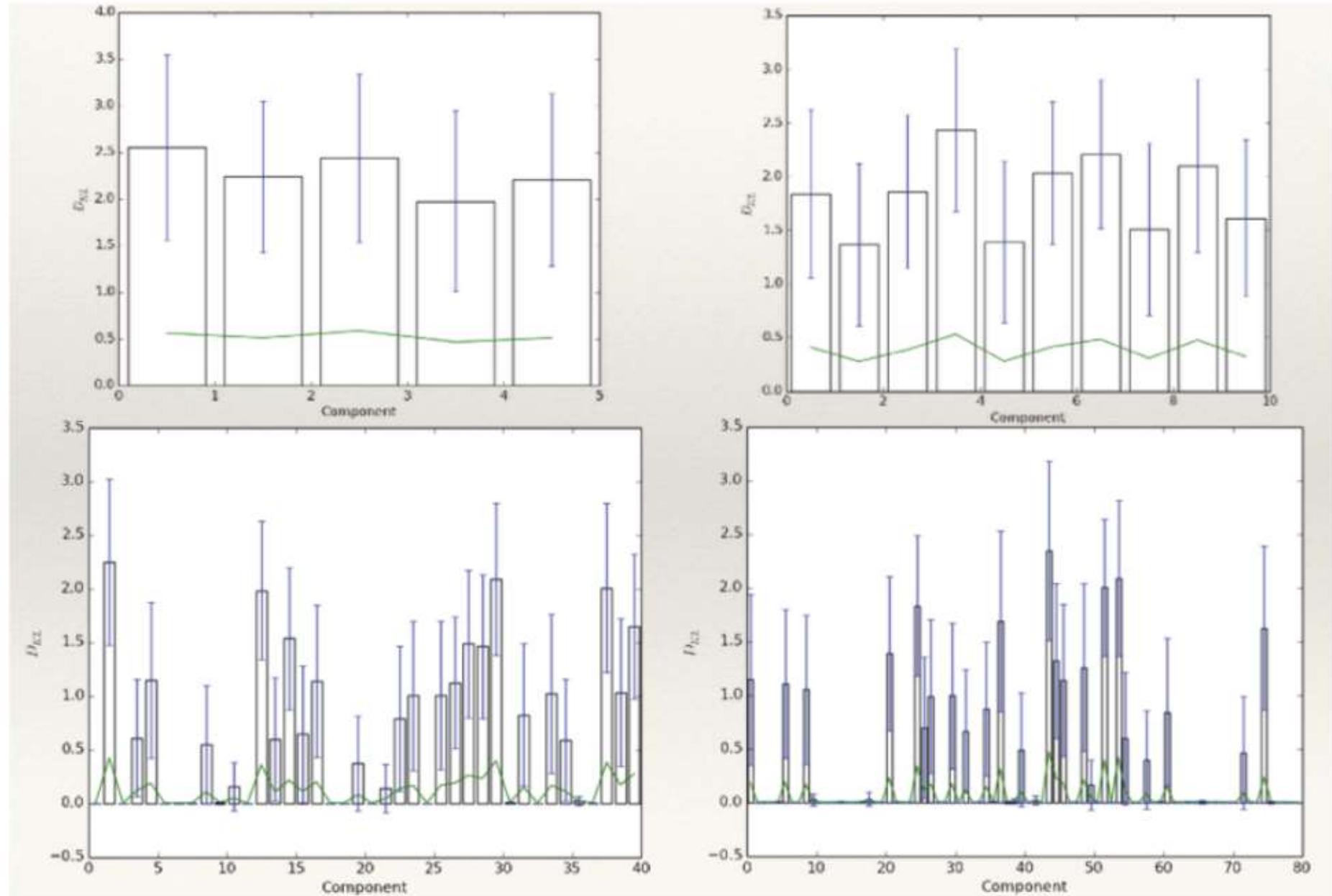


Figure from Laurent Dinh & Vincent Dumoulin

Collapsing components

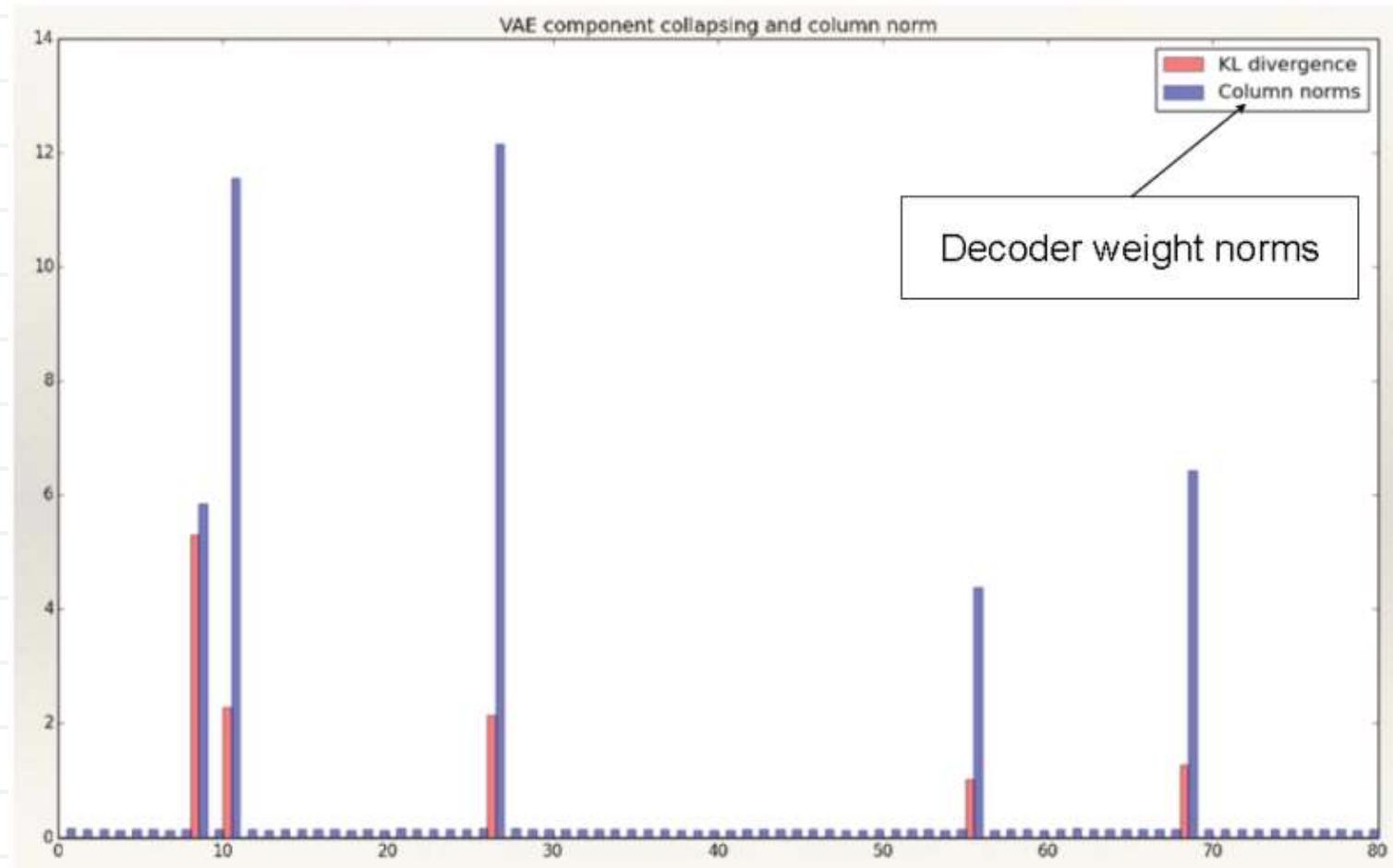


Figure from Laurent Dinh & Vincent Dumoulin

- Components that are shrunk to $N(0, 1)$ are not used by decoder

Sample quality comparison



VAE



GAN

Source:
OpenAI
blog

Sample quality comparison

VAE:

Source:
OpenAI blog



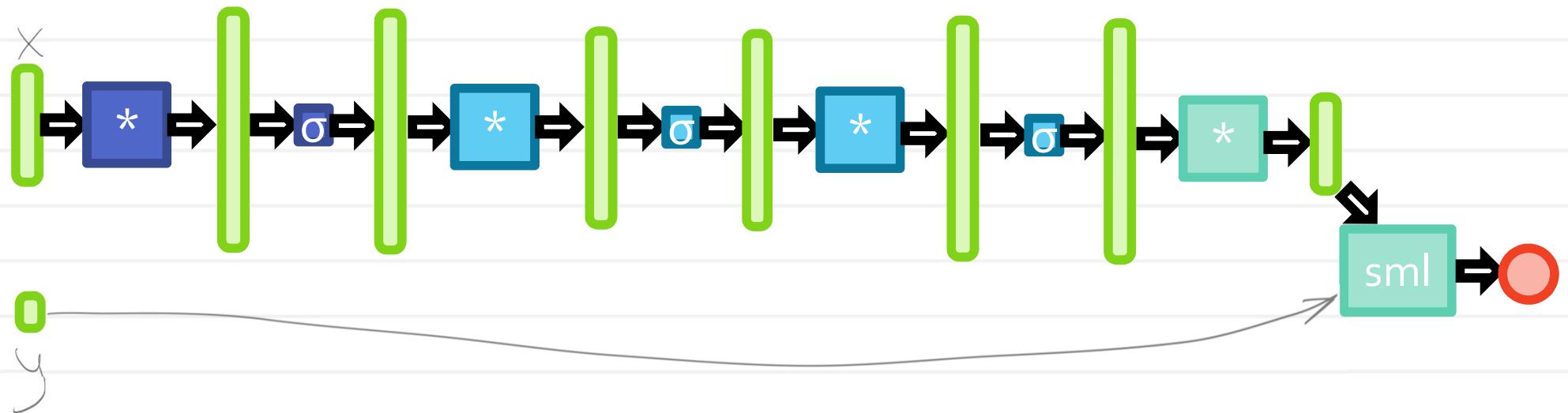
Sample quality comparison

DCGAN:



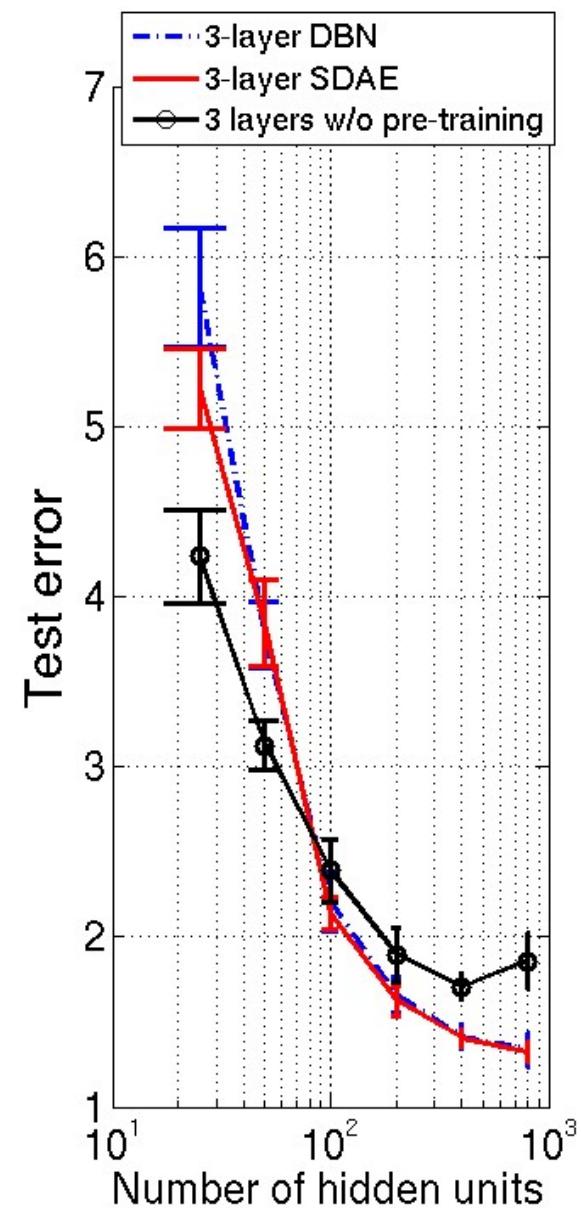
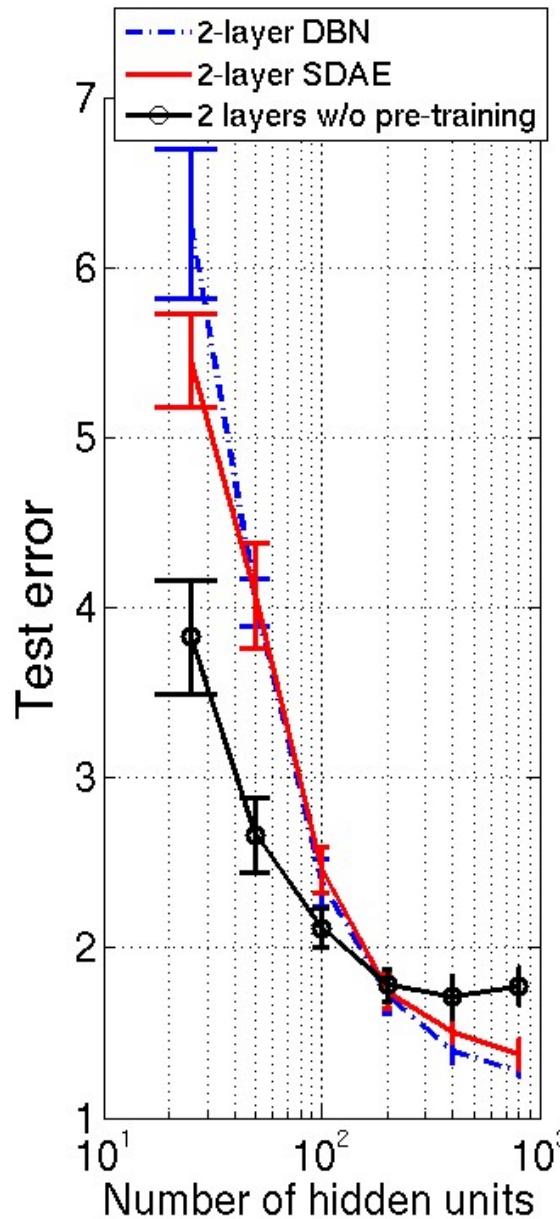
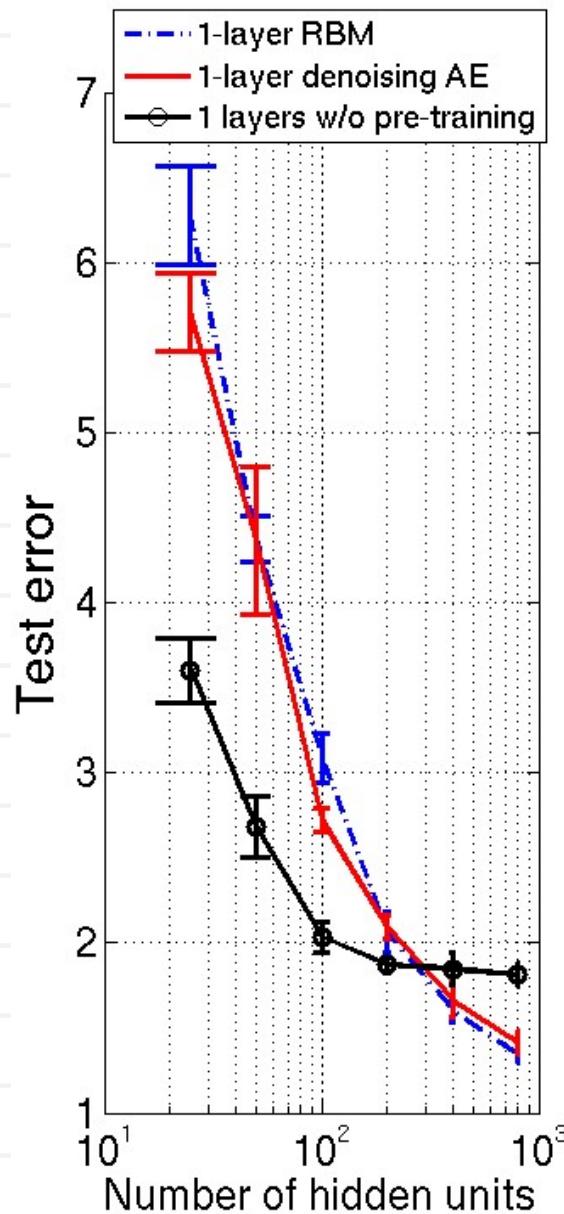
Source:
OpenAI blog

Using auto-encoders for pre-training



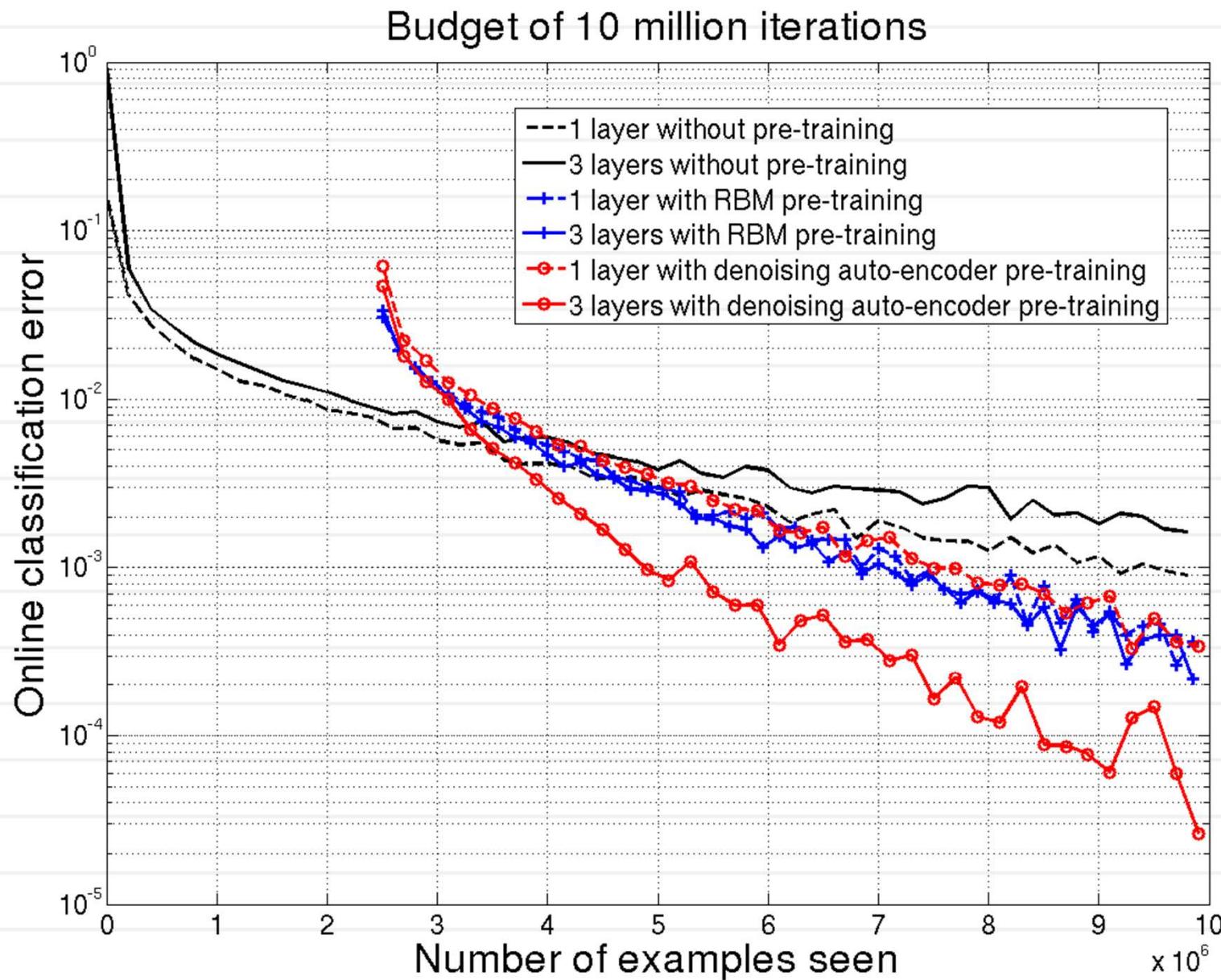
- Autoencoder is used to initialize a deep supervised feedforward architecture

Pretraining by unsupervised learning



[Erhan et al. 2010]

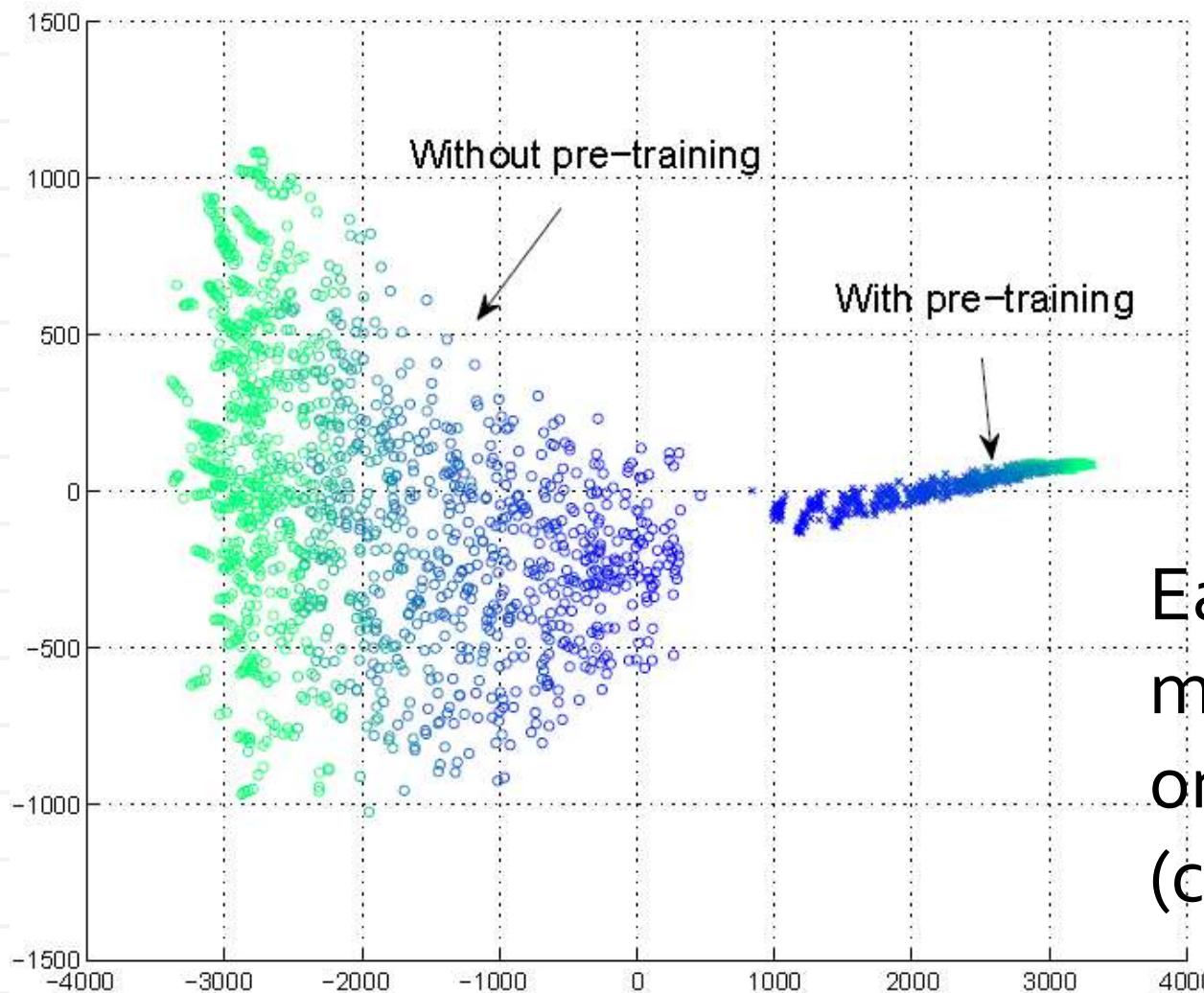
Pretraining by unsupervised learning



[Erhan et al. 2010]

Pretraining by unsupervised learning

Isomap of model space (greener = later)

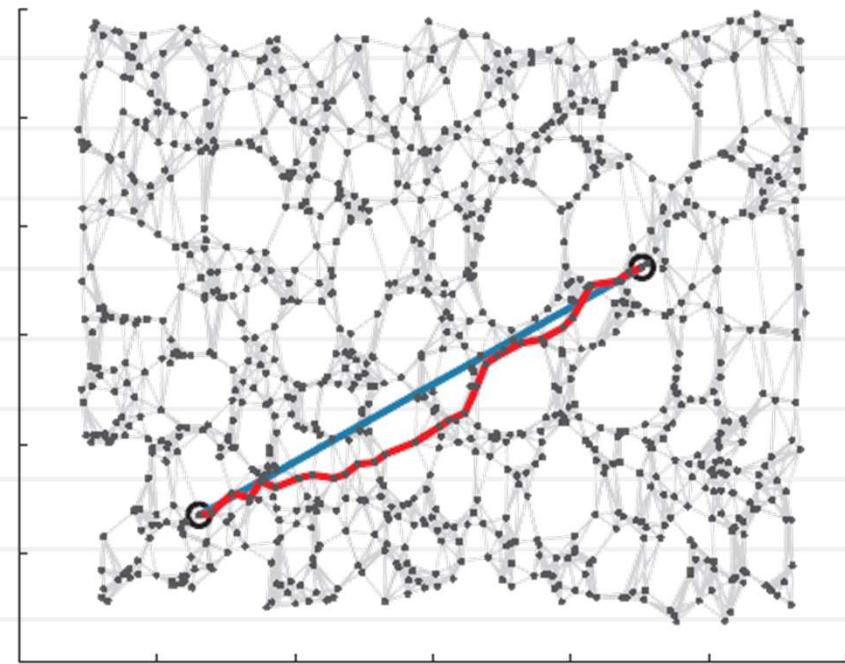
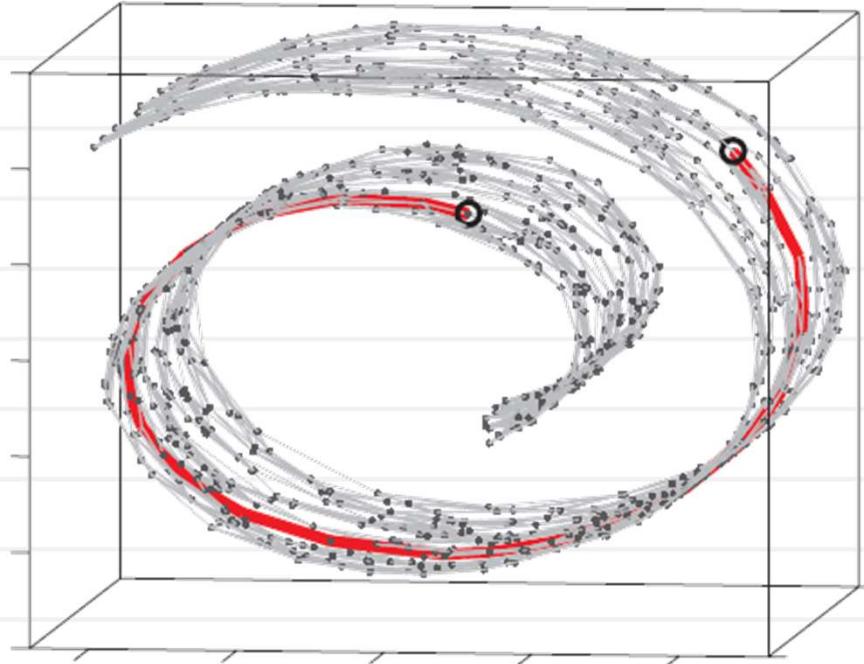


Each network is
mapped to its output
on a validation set
(concat. posteriors)

[Erhan et al. 2010]

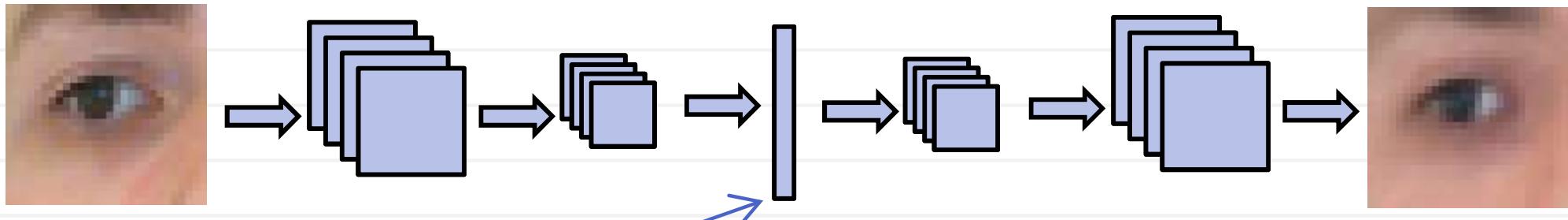
Isomap manifold learning

- Build a NN graph (unless your data are already a graph!)
- $d_{ij} = \text{shortest_path}(V_i, V_j)$
- Do MDS (project while preserving d_{ij})

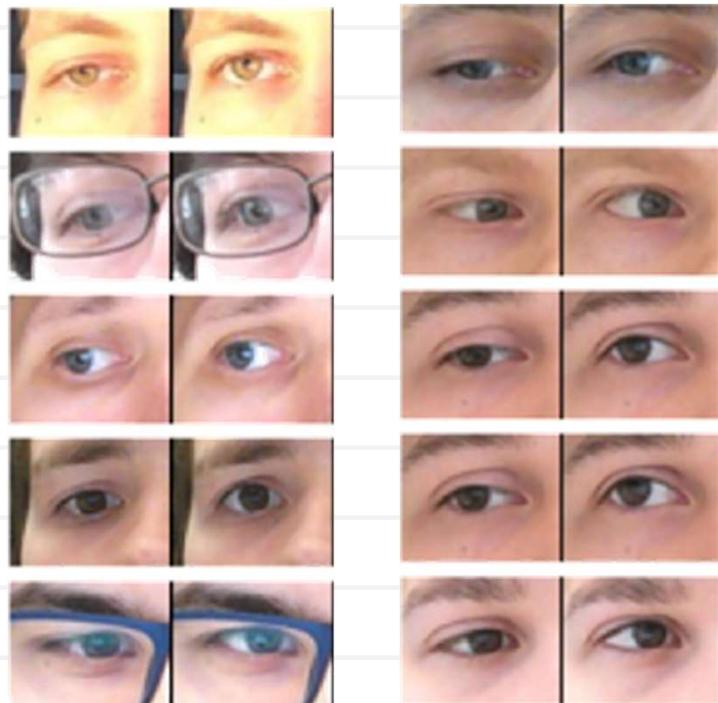


[Tenenbaum et al. *Science*'2000]

Smart image editing

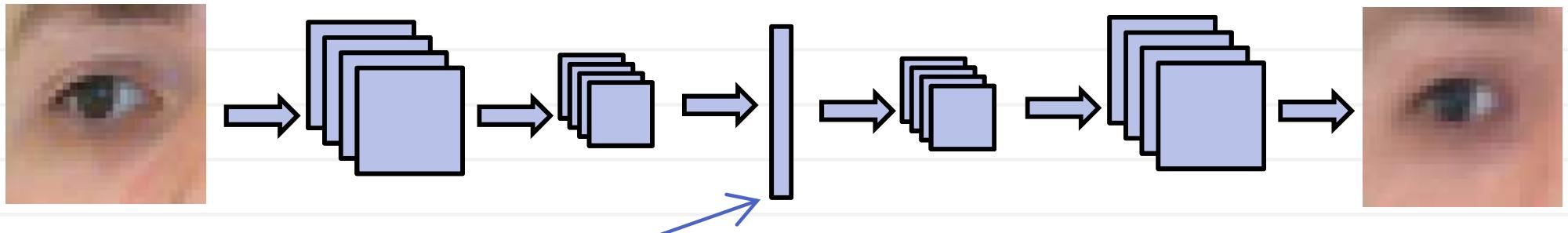


Low-dim space, where we can estimate semantically meaningful directions

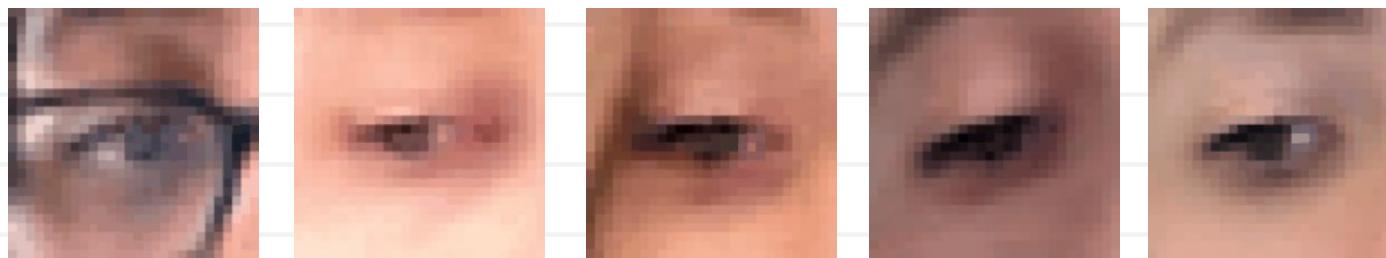


Given a bunch of pairs we can estimate a vector for gaze redirection

Latent space image editing



Low-dim space, where we can estimate semantically meaningful directions

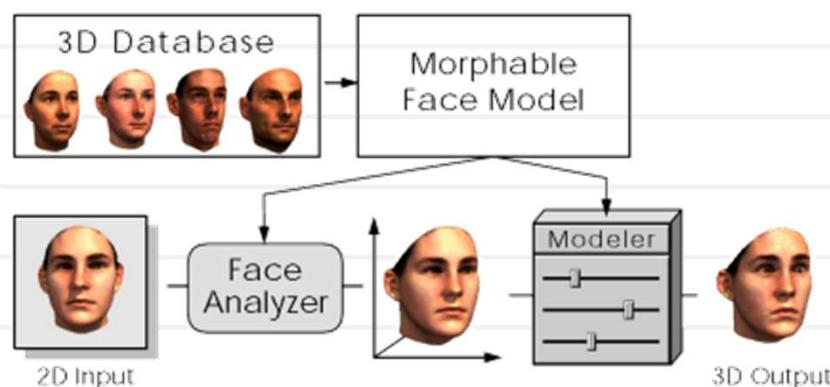
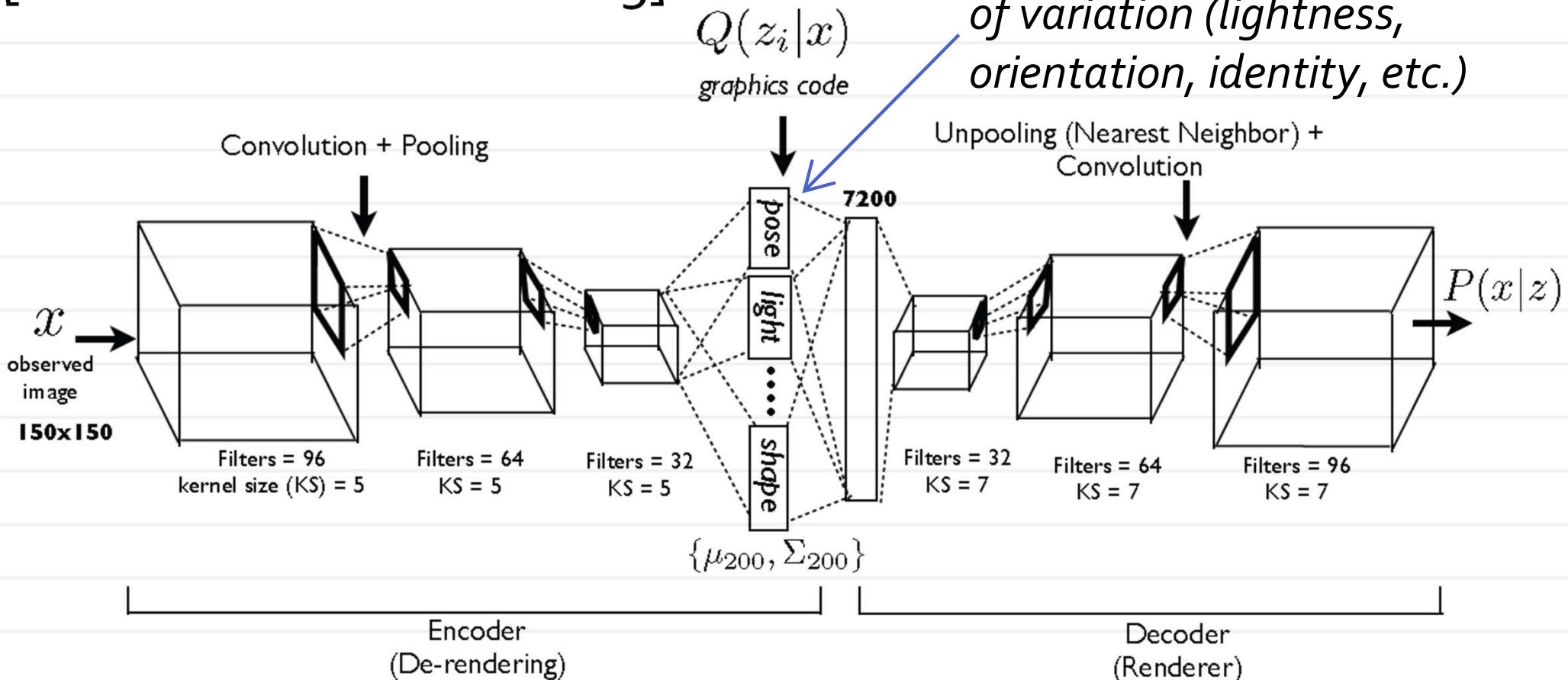


Thanks to L.Yekimov
(+F.Chervinsky, D.Kononenko, D.Sungatullina, Y.Ganin)

Deep Inverse Graphics

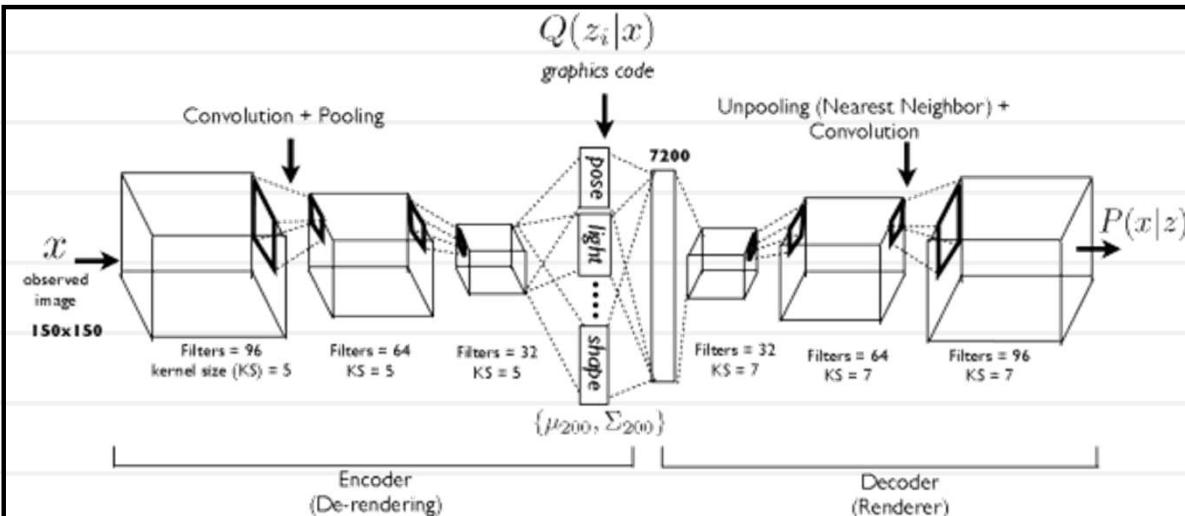
[Kulkarni et al. NIPS 2015]

Different variables z are assigned to different factors of variation (lightness, orientation, identity, etc.)



[Blanz and Vetter 1999]
morphable model

Ensuring semantic meaning of dimensions



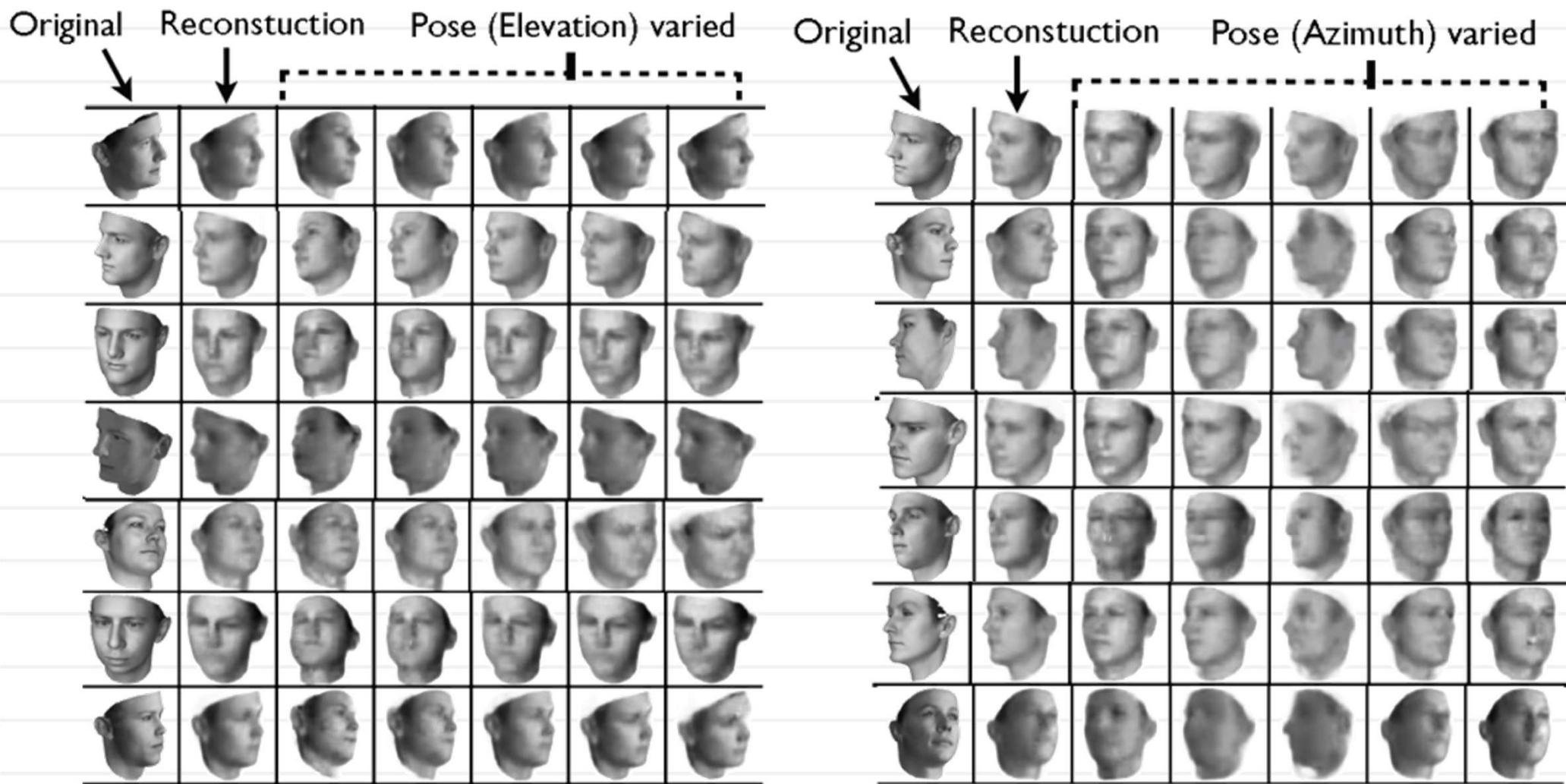
Training with clamping:

- Draw minibatches with some factors fixed
- During forwardprop replace corresponding z with averages
- During backprop replace loss gradients with

$$z_j^i = \bar{z}_j$$

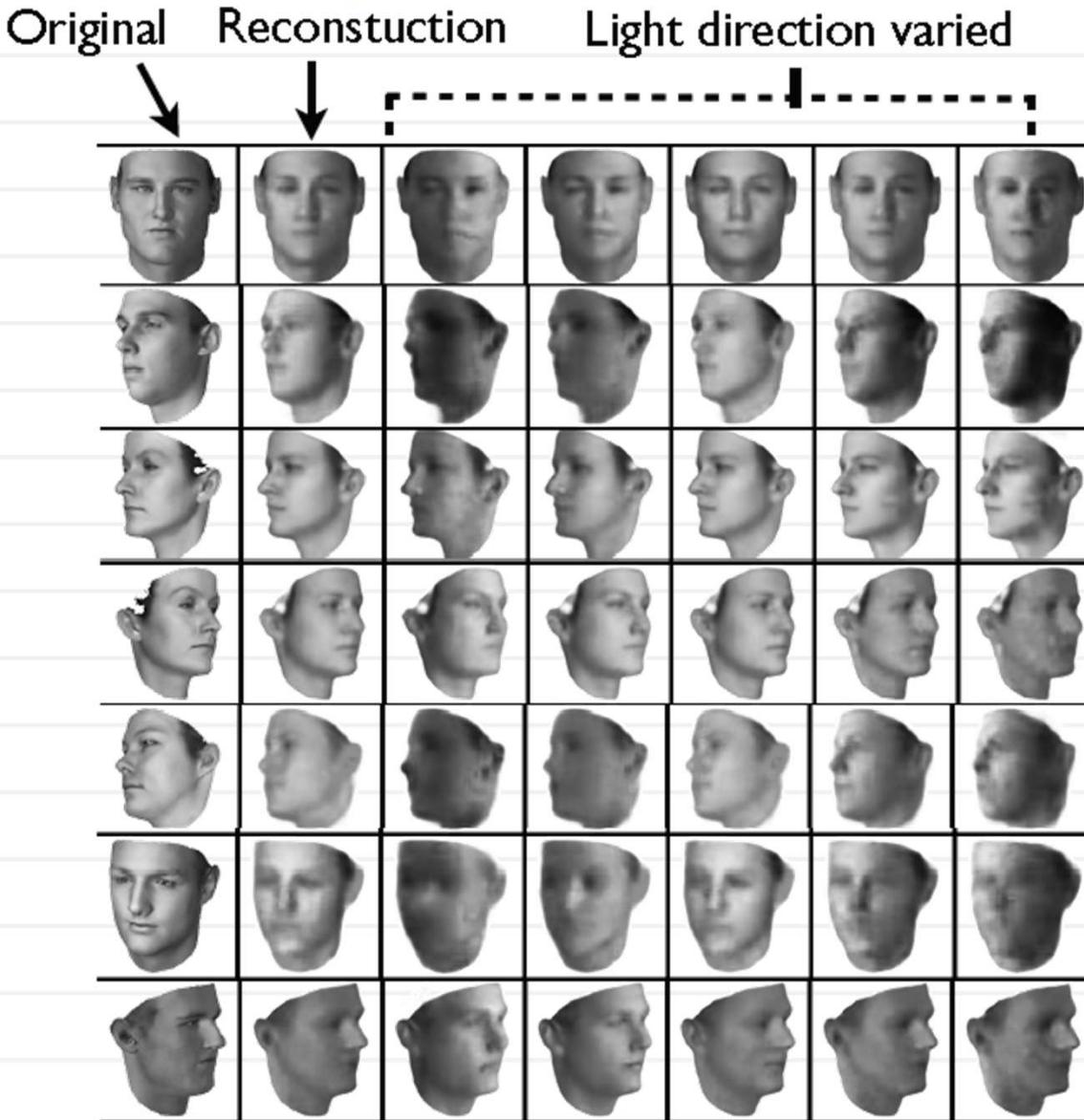
$$\frac{\partial L}{\partial z_j^i} = z_j^i - \bar{z}_j$$

Deep Inverse Graphics



[Kulkarni et al. NIPS 2015]

Deep Inverse Graphics

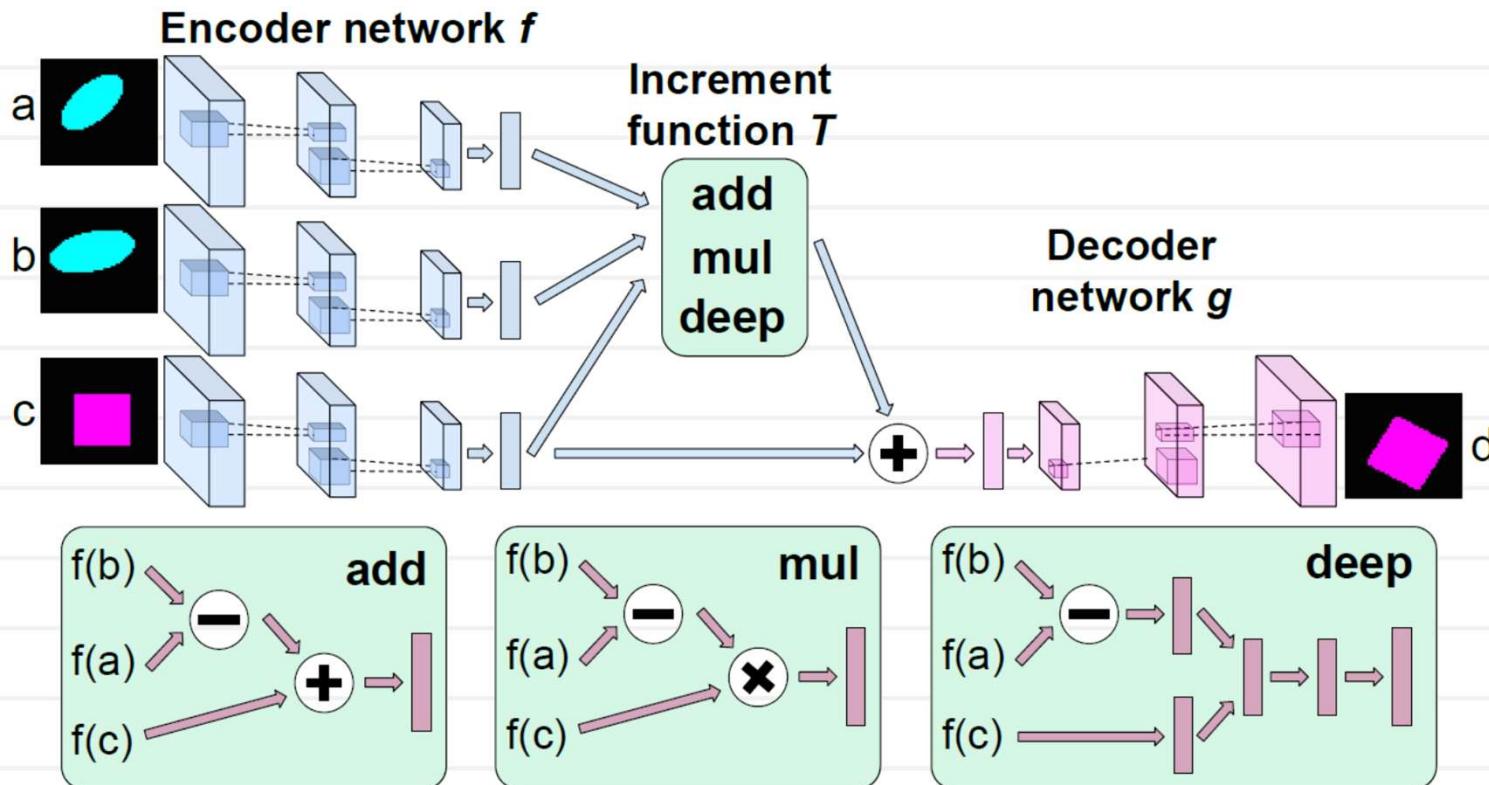


[Kulkarni et al. NIPS 2015]

“Deep visual analogy making”

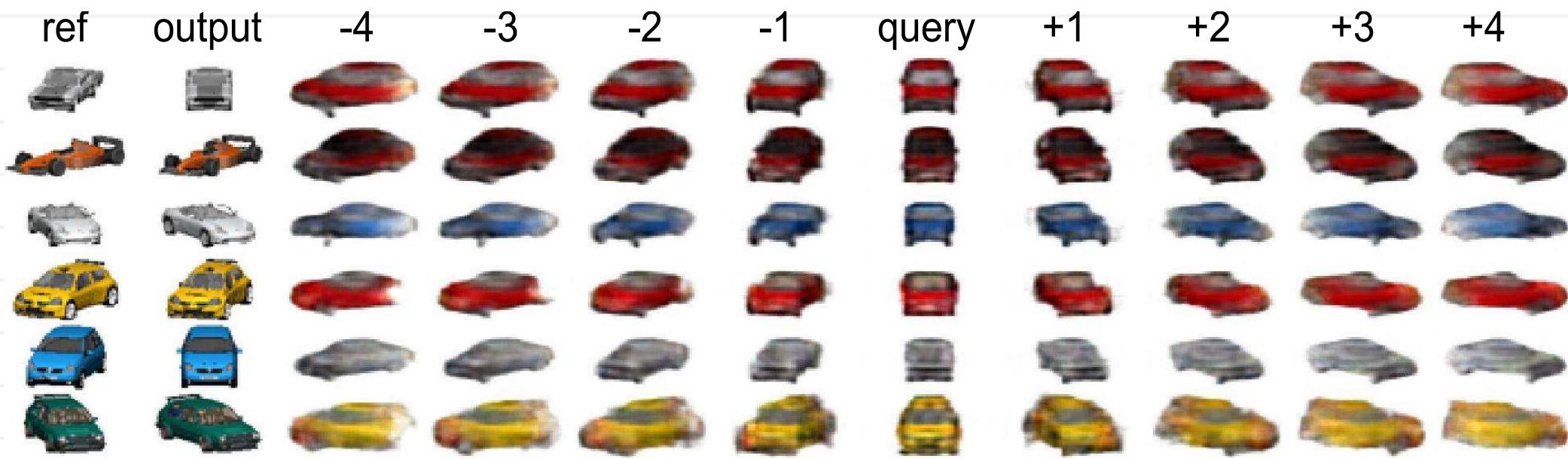
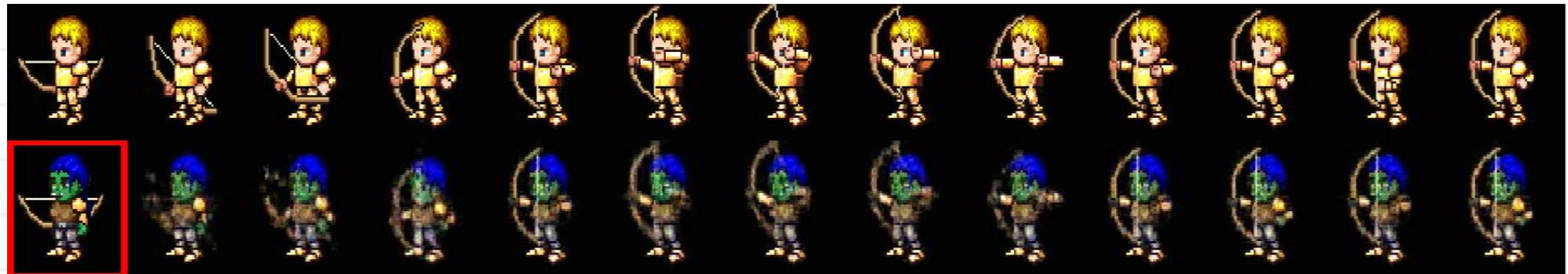


Fully-supervised setting:



[Reed et al. NIPS 2015]

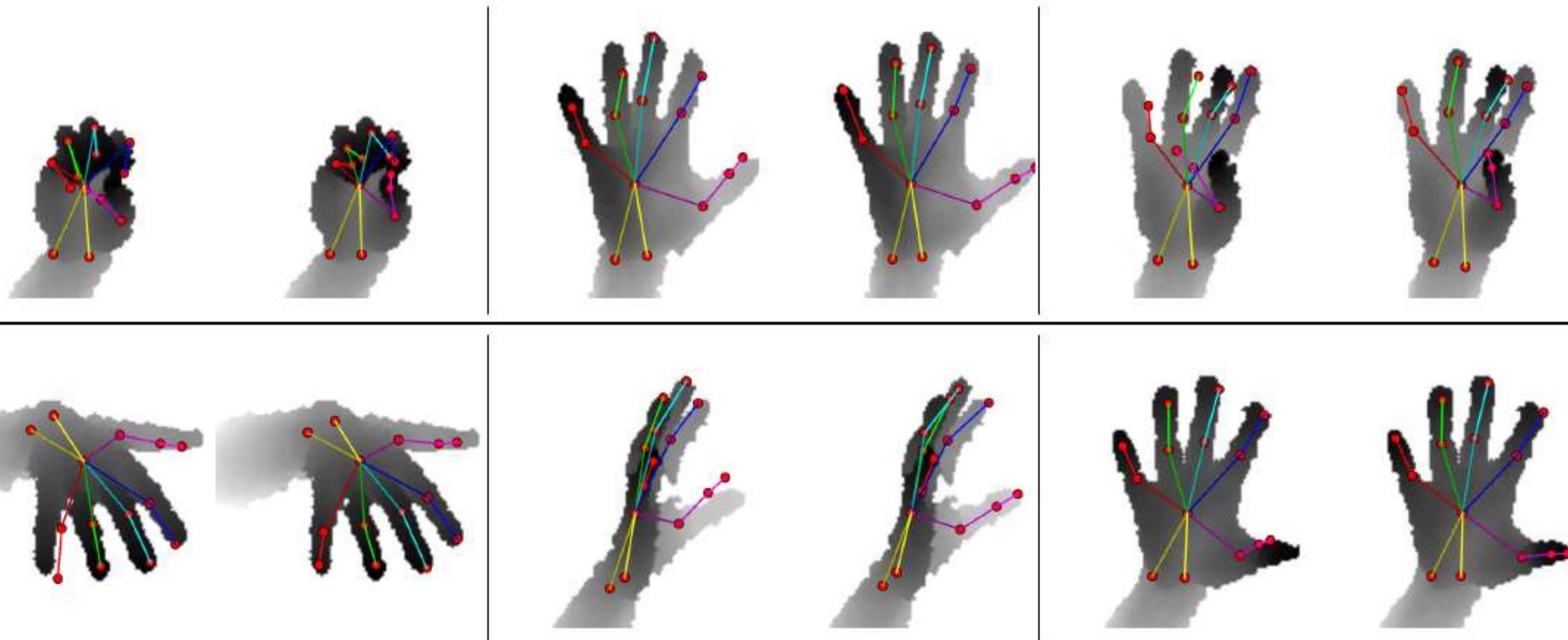
“Deep visual analogy making”



see also [Reed et al. NIPS 2015]

Vision with Feedback loop: background

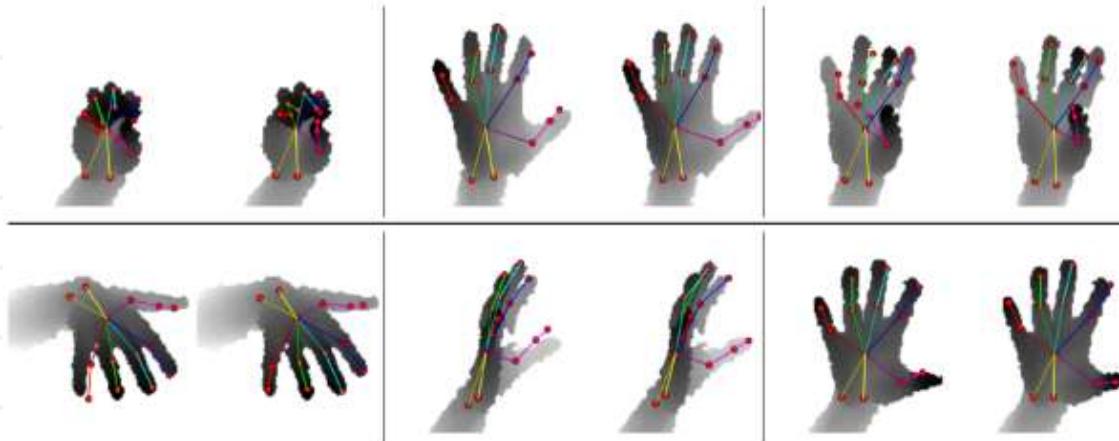
Problem setting:



- Difficult, even humans far from perfect
- Synthetic data is reasonable

[Oberweger, Wohlhart, Lepetit ICCV15]

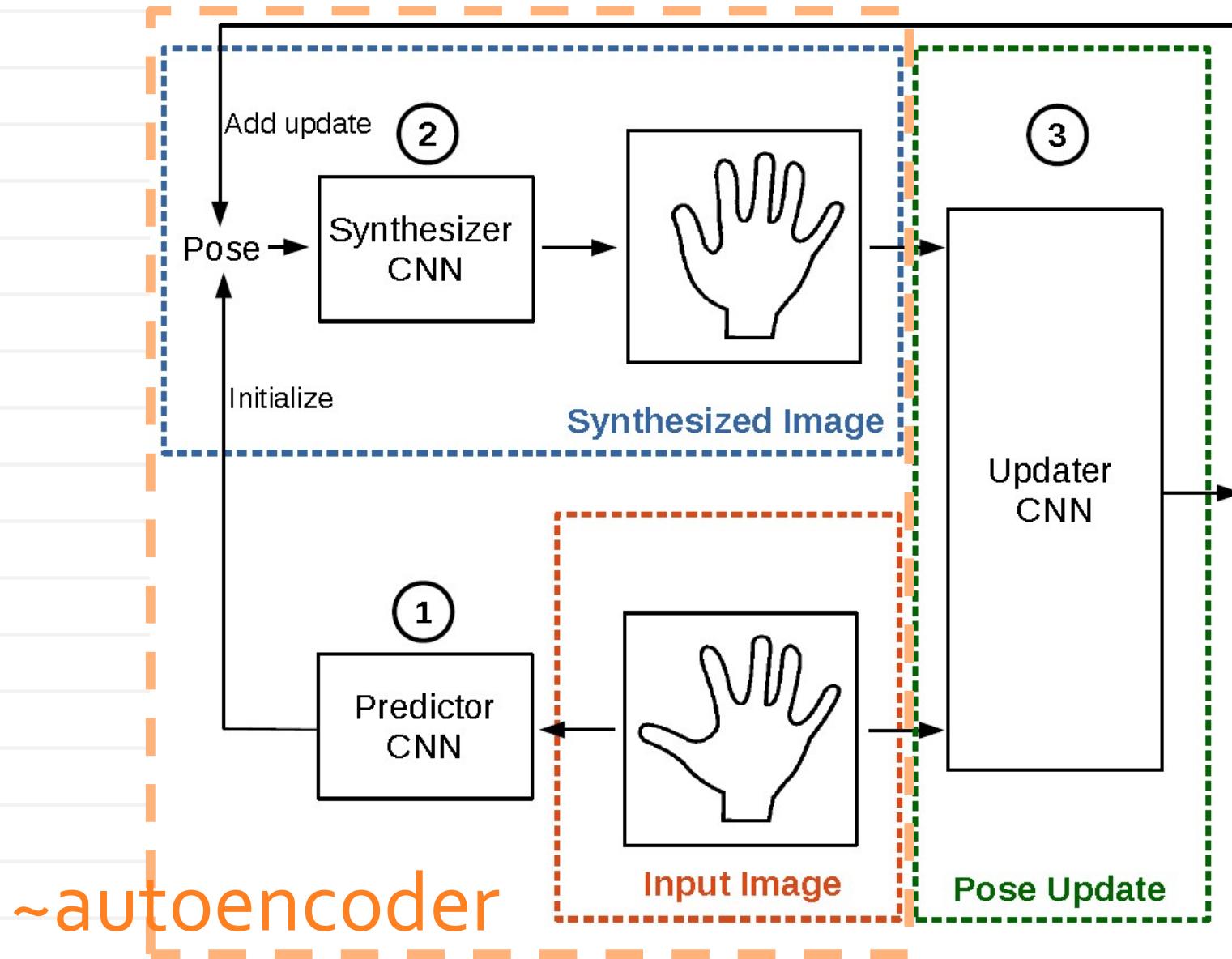
Vision with Feedback loop : background



SoA: $\|I - R(p)\| \rightarrow \min_p$
“gold standard algorithm”

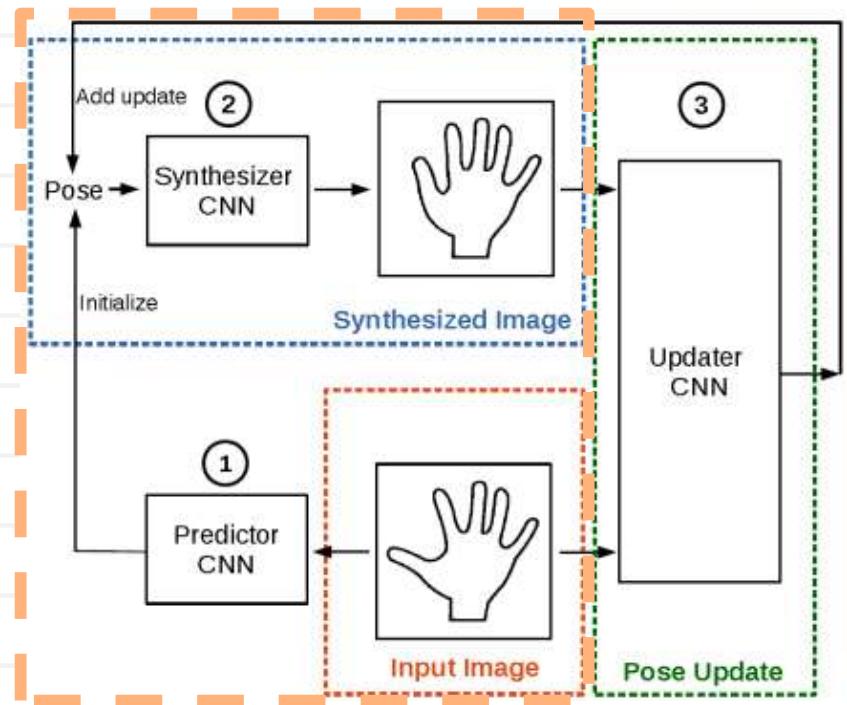
- Slow
- Get stuck in poor minima
- Needs good model
- When properly engineered, very hard to beat

Vision with Feedback loop



[Oberweger, Wohlhart, Lepetit ICCV15]

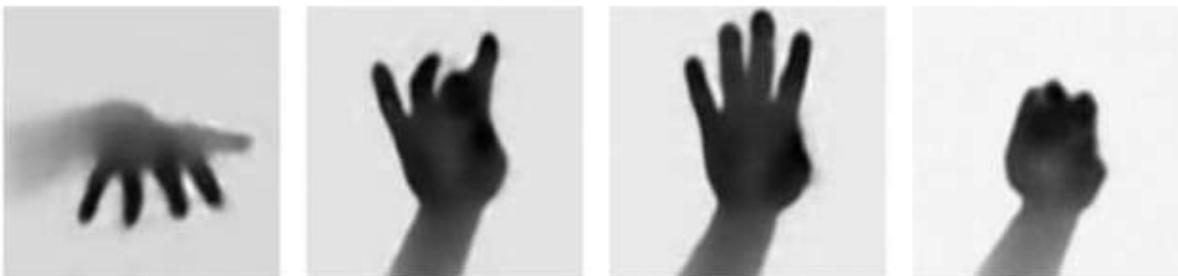
Vision with Feedback loop



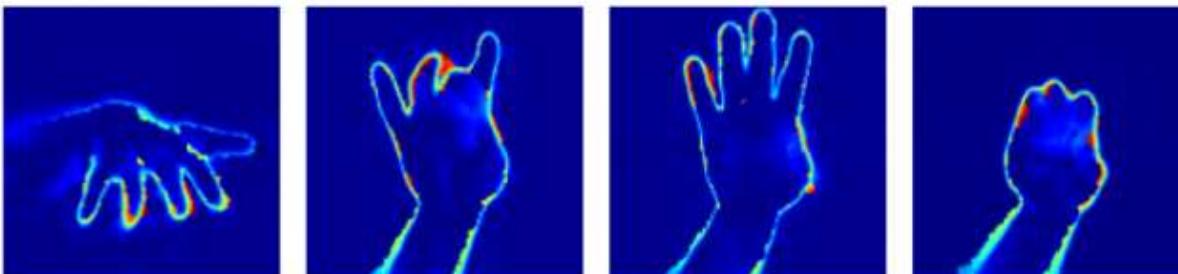
[Oberweger et al. ICCV15]



Autoencoder
reconstruction:

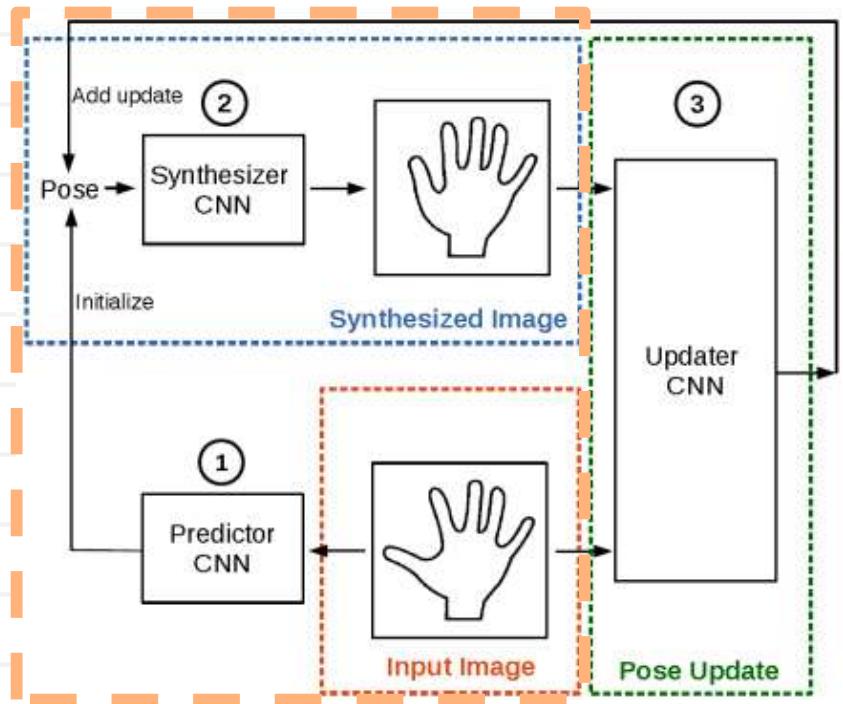


Difference:

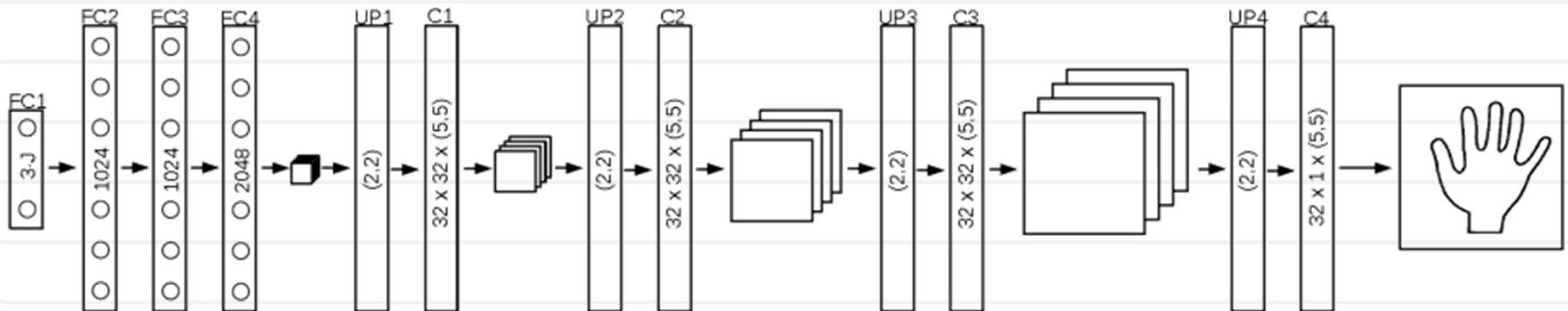


Vision with Feedback loop

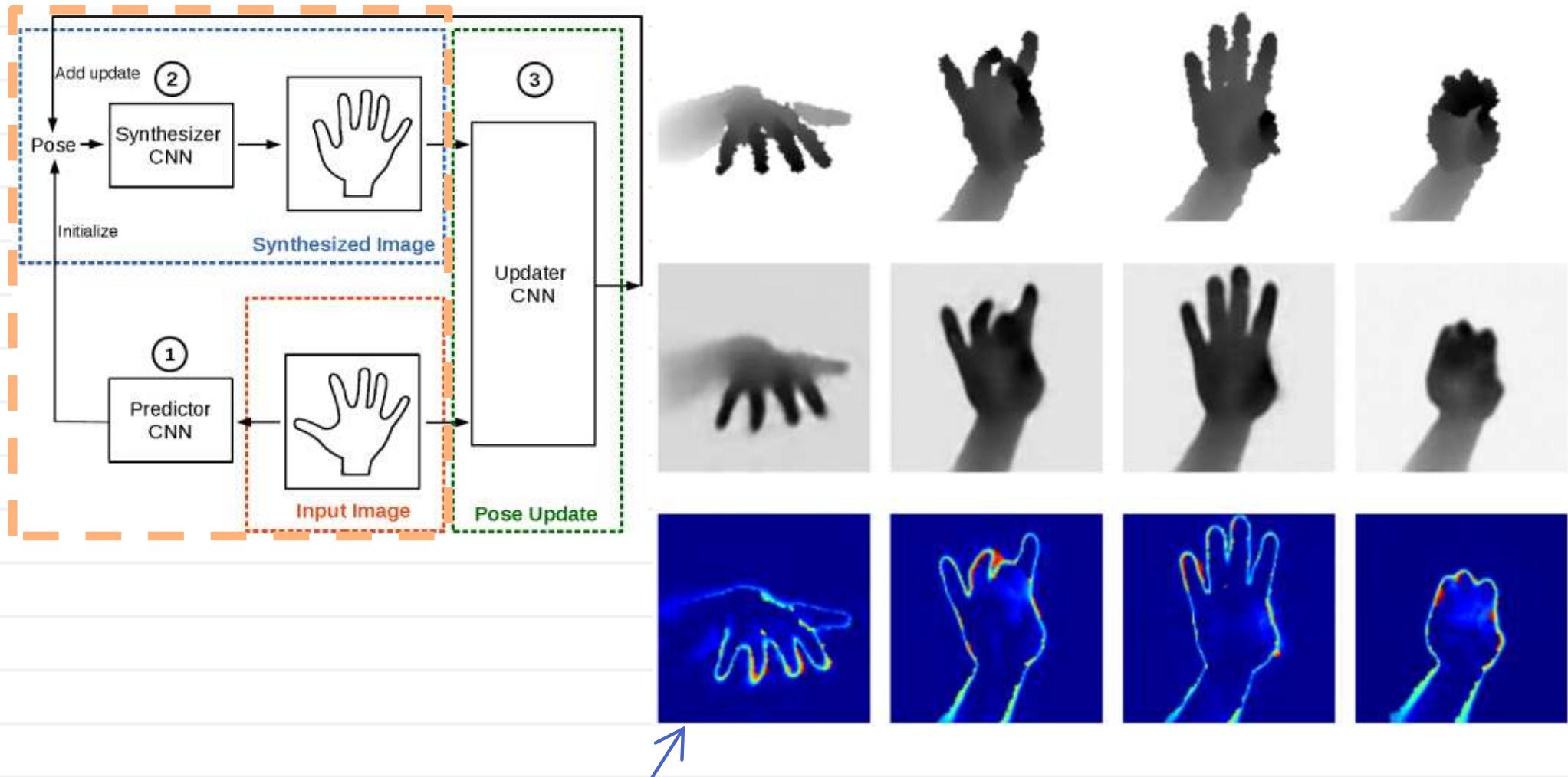
[Oberweger et al. ICCV15]



Synthesizer (decoder):



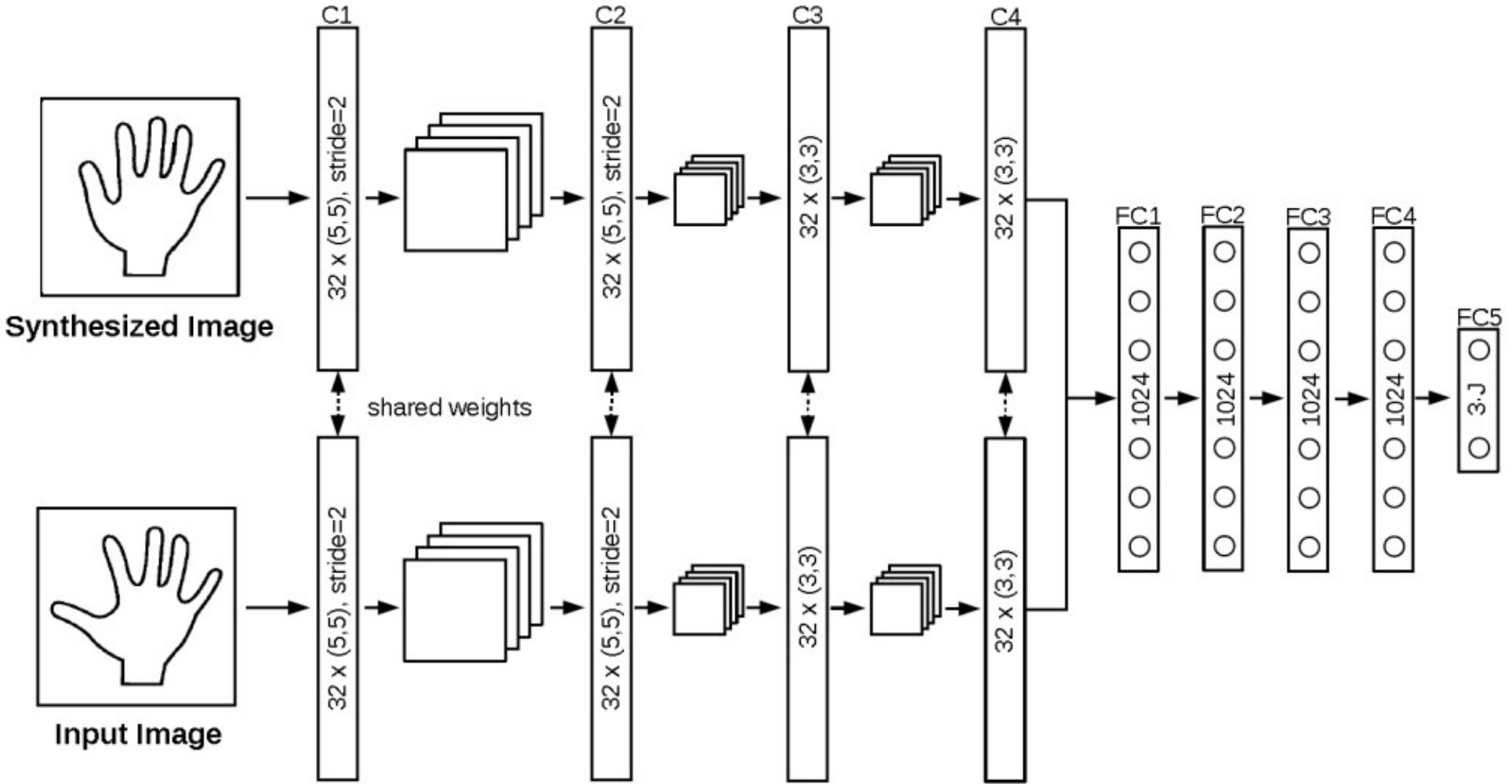
Vision with Feedback loop



Where to improve, how to improve

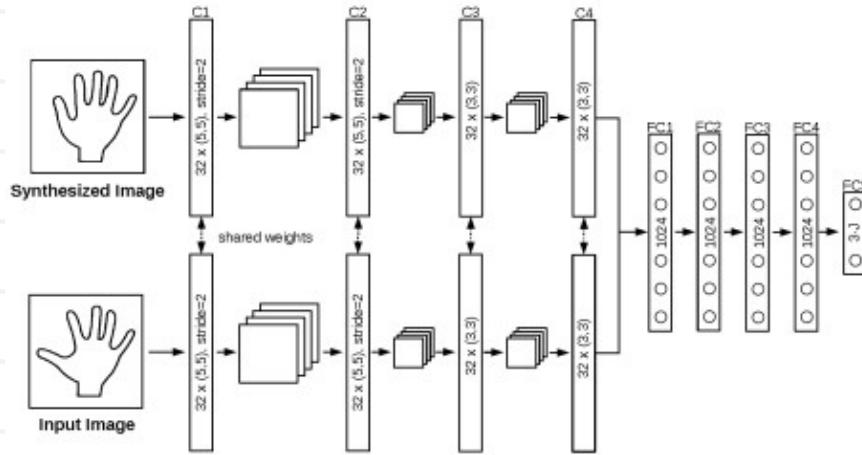
[Oberweger, Wohlhart, Lepetit ICCV15]

Updater network



[Oberweger, Wohlhart, Lepetit ICCV15]

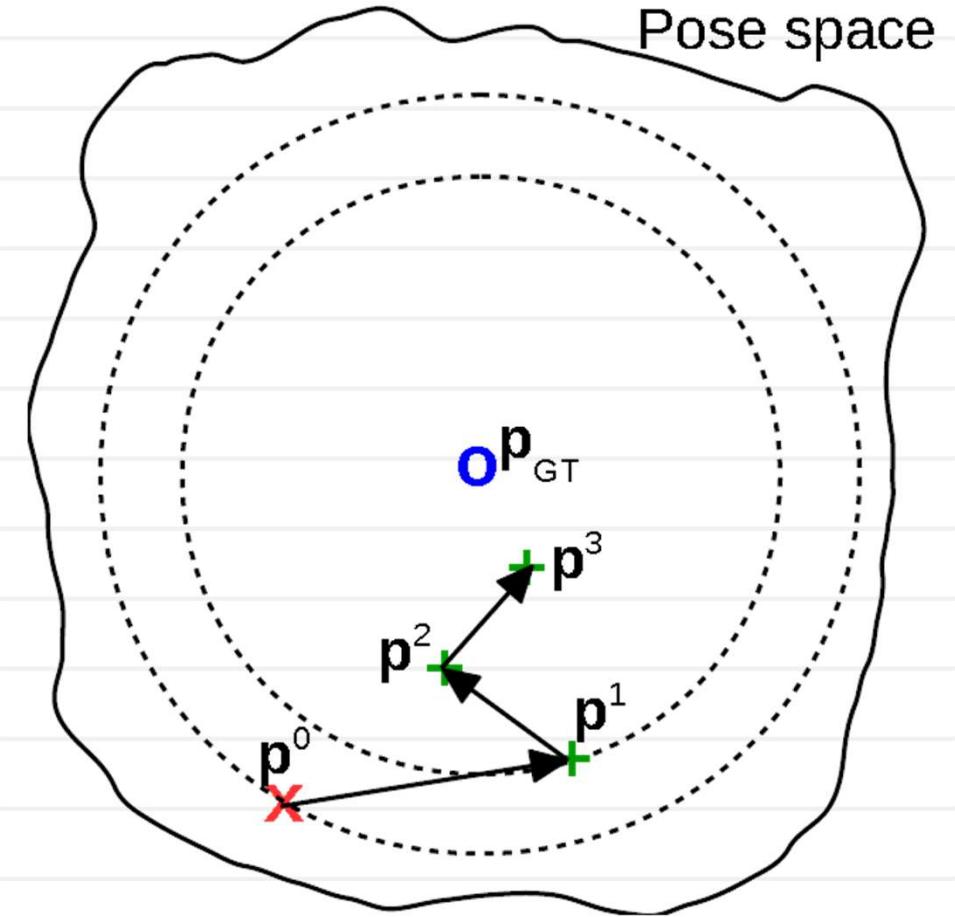
Updater network



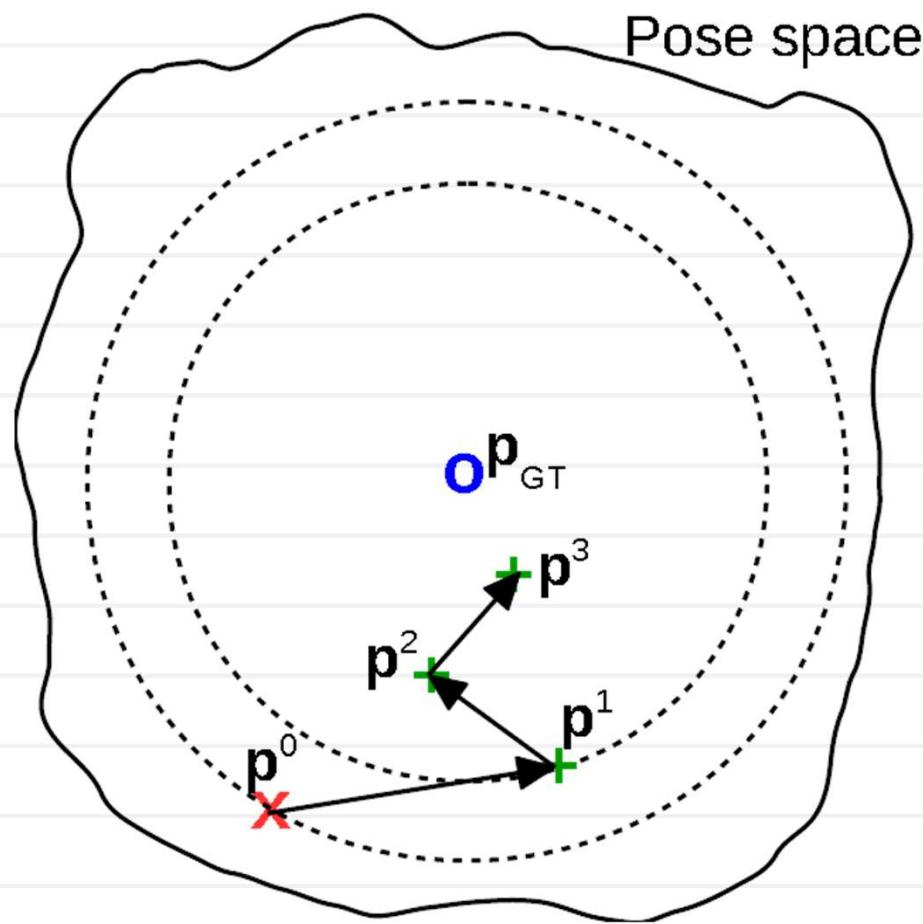
- Training with L₂ loss does not work (too hard)
- Instead:

$$\mathcal{L} = \sum_{(\mathcal{D}, \mathbf{p}) \in \mathcal{T}} \sum_{\mathbf{p}' \in \mathcal{T}_{\mathcal{D}}} \max(0, \|\mathbf{p}'' - \mathbf{p}\|_2 - \lambda \|\mathbf{p}' - \mathbf{p}\|_2)$$

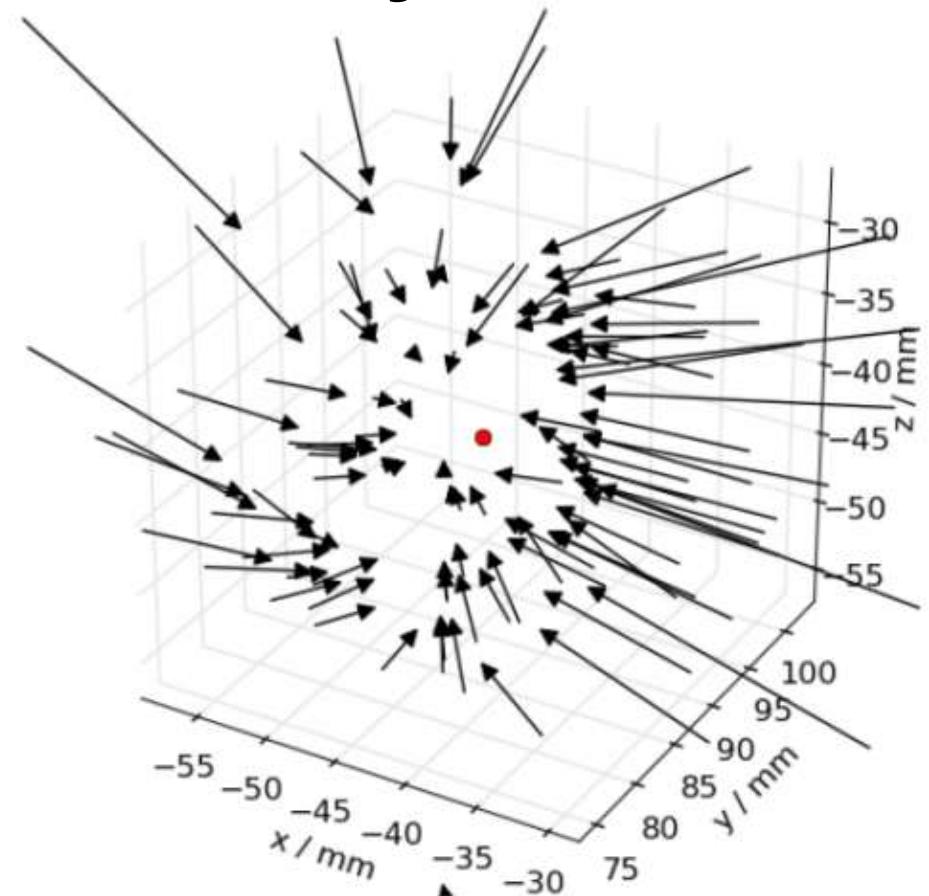
[Oberweger et al. ICCV15]



Resulting updates



One of joints:



$$\mathcal{L} = \sum_{(\mathcal{D}, \mathbf{p}) \in \mathcal{T}} \sum_{\mathbf{p}' \in \mathcal{T}_{\mathcal{D}}} \max(0, \|\mathbf{p}'' - \mathbf{p}\|_2 - \lambda \|\mathbf{p}' - \mathbf{p}\|_2)$$

[Oberweger et al. ICCV15]

Qualitative comparison

“Gold standard”

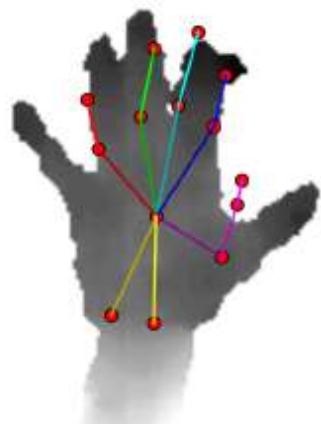
Init

Init

Iter 1

Iter 2

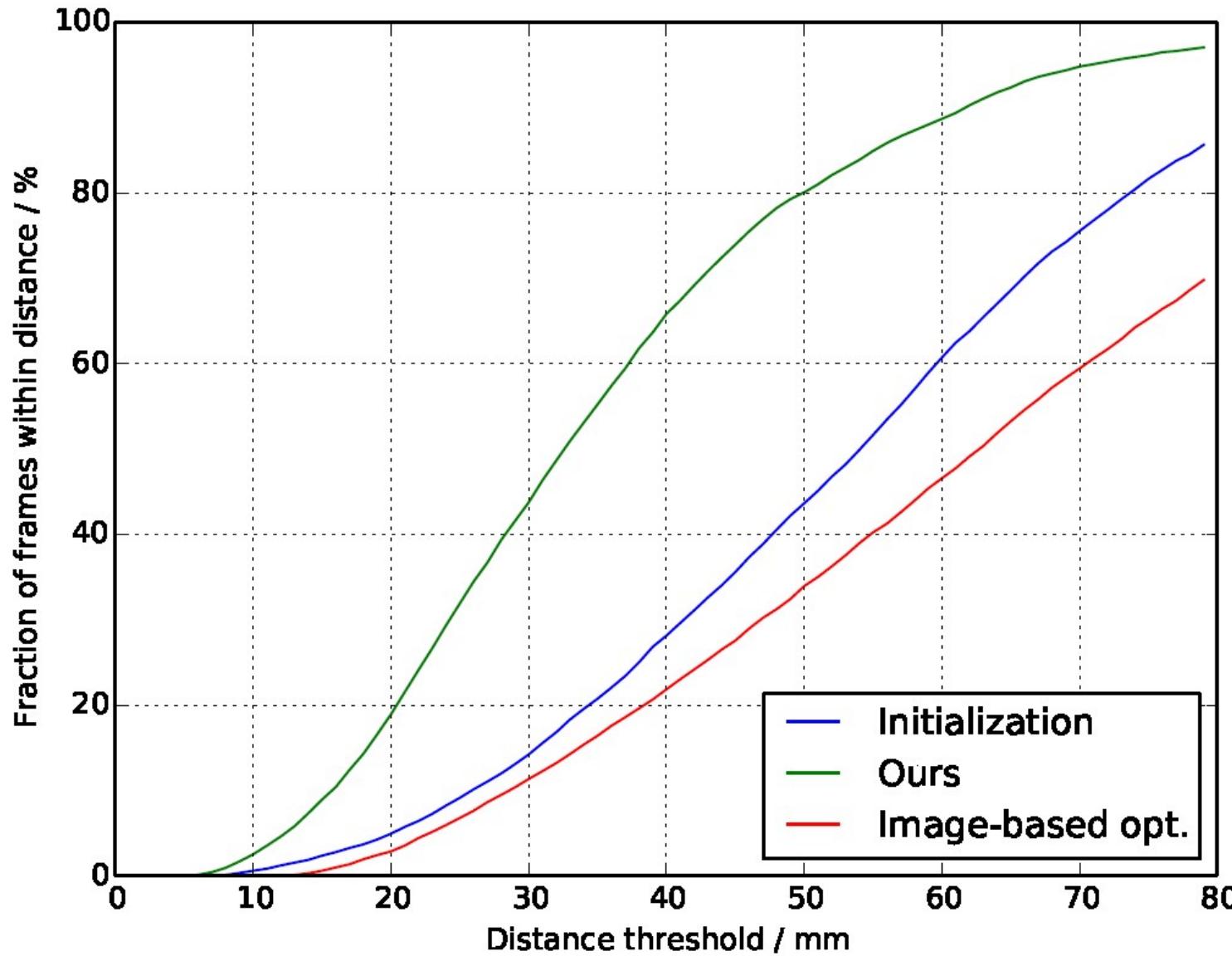
Final



Theirs

[Oberweger, Wohlhart, Lepetit ICCV15]

Quantitative comparison



[Oberweger, Wohlhart, Lepetit ICCV15]

Recap

- Autoencoders are natural ways to tackle unsupervised learning with DL
- Natural ways to regularize: weight decay, contractive, denoising, *variational*
- Useful for pretraining if little annotated data
- Useful for generating and learning transforms
- Can be part of the bigger systems (c.f. last example)

Bibliography

Yoshua Bengio, Pascal Lamblin, Dan Popovici, Hugo Larochelle:
Greedy Layer-Wise Training of Deep Networks. NIPS 2006: 153-160

Andrew Ng, Stanford class CS294a lecture notes, Handout #2

Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio,
Pierre-Antoine Manzagol: Stacked Denoising Autoencoders:
Learning Useful Representations in a Deep Network with a Local
Denoising Criterion. Journal of Machine Learning Research 11:
3371-3408 (2010)

Salah Rifai, Pascal Vincent, Xavier Muller, Xavier Glorot, Yoshua
Bengio: Contractive Auto-Encoders: Explicit Invariance During
Feature Extraction. ICML 2011: 833-840

Bibliography

Diederik P. Kingma, Max Welling:
Auto-Encoding Variational Bayes. ICLR 2014

Aaron Courville, Variational Autoencoders and Extensions,
videolectures.net

Dumitru Erhan, Aaron C. Courville, Yoshua Bengio, Pascal Vincent:
Why Does Unsupervised Pre-training Help Deep Learning?
AISTATS 2010: 201-208

Tenenbaum, Joshua B., Vin De Silva, and John C. Langford. "A
global geometric framework for nonlinear dimensionality
reduction." Science 290.5500 (2000): 2319-2323.

Bibliography

Tejas D. Kulkarni, Will Whitney, Pushmeet Kohli, Joshua B. Tenenbaum: Deep Convolutional Inverse Graphics Network. NIPS 2015

Volker Blanz, Thomas Vetter:
A Morphable Model for the Synthesis of 3D Faces. SIGGRAPH 1999: 187-194

Reed, S.E., Zhang, Y., Zhang, Y. and Lee, H., et al. "Deep Visual Analogy-Making." NIPS 2015.

Markus Oberweger, Paul Wohlhart, Vincent Lepetit:
Training a Feedback Loop for Hand Pose Estimation. ICCV 2015:
3316-3324