

M A D



NYU

COURANT INSTITUTE OF
MATHEMATICAL SCIENCES

MATHEMATICS OF DEEP LEARNING

JOAN BRUNA , CIMS + CDS, NYU, SPRING'18

Lecture 3: The Scattering Transform and beyond

WAVELET COVARIANTS

- If we replace input image by first layer output:

$$\rho(x_0 \star \psi_{j,\theta})(u) = x_1(u, j, \theta)$$

Let $\tilde{x}_0 = R_\alpha x_0$ be a rotation of α degrees.

$$\rho(\tilde{x}_0 \star \psi_{j,\theta})(u) = x_1(R_\alpha u, j, \theta + \alpha)$$

WAVELET COVARIANTS

- If we replace input image by first layer output:

$$\rho(x_0 \star \psi_{j,\theta})(u) = x_1(u, j, \theta)$$

Let $\tilde{x}_0 = R_\alpha x_0$ be a rotation of α degrees.

$$\rho(\tilde{x}_0 \star \psi_{j,\theta})(u) = x_1(R_\alpha u, j, \theta + \alpha)$$

- Similarly, roto-translation acts on x_1 by rotating and translating spatial coordinates and translating orientation coordinates

Let $\tilde{x}_0 = \varphi_{(v,\alpha)} x_0$ be a roto-translation with parameters (v, α) .

$$\rho(\tilde{x}_0 \star \psi_{j,\theta})(u) = x_1(\varphi_v R_\alpha u, j, \theta + \alpha)$$

WAVELET COVARIANTS

- If we replace input image by first layer output:

$$\rho(x_0 \star \psi_{j,\theta})(u) = x_1(u, j, \theta)$$

Let $\tilde{x}_0 = R_\alpha x_0$ be a rotation of α degrees.

$$\rho(\tilde{x}_0 \star \psi_{j,\theta})(u) = x_1(R_\alpha u, j, \theta + \alpha)$$

- Similarly, roto-translation acts on x_1 by rotating and translating spatial coordinates and translating orientation coordinates

Let $\tilde{x}_0 = \varphi_{(v,\alpha)} x_0$ be a roto-translation with parameters (v, α) .

$$\rho(\tilde{x}_0 \star \psi_{j,\theta})(u) = x_1(\varphi_v R_\alpha u, j, \theta + \alpha)$$

- So we can replace convolutions over translation by convolutions over roto-translations.

GROUP CONVOLUTIONS



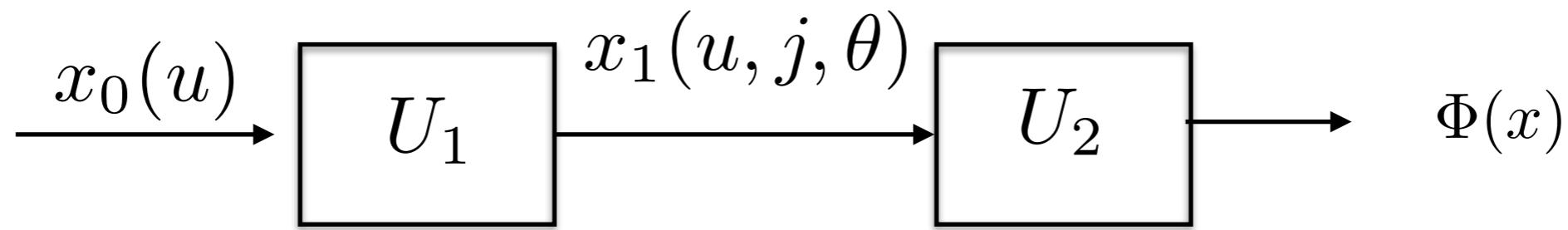
Definition: Let G be a group equipped with a Haar measure $d\mu$, acting on Ω , and $h \in L^1(G)$. The group convolution $x \star_G h$ is defined as

$$x \star_G h(u) = \int_G h(g)x(\varphi_g u)d\mu(g) , \quad x \in L^2(\Omega) .$$

If $x = x_1(u, j, \theta)$ and G are roto-translations, these convolutions recombine different orientation channels.

JOINT SCATTERING

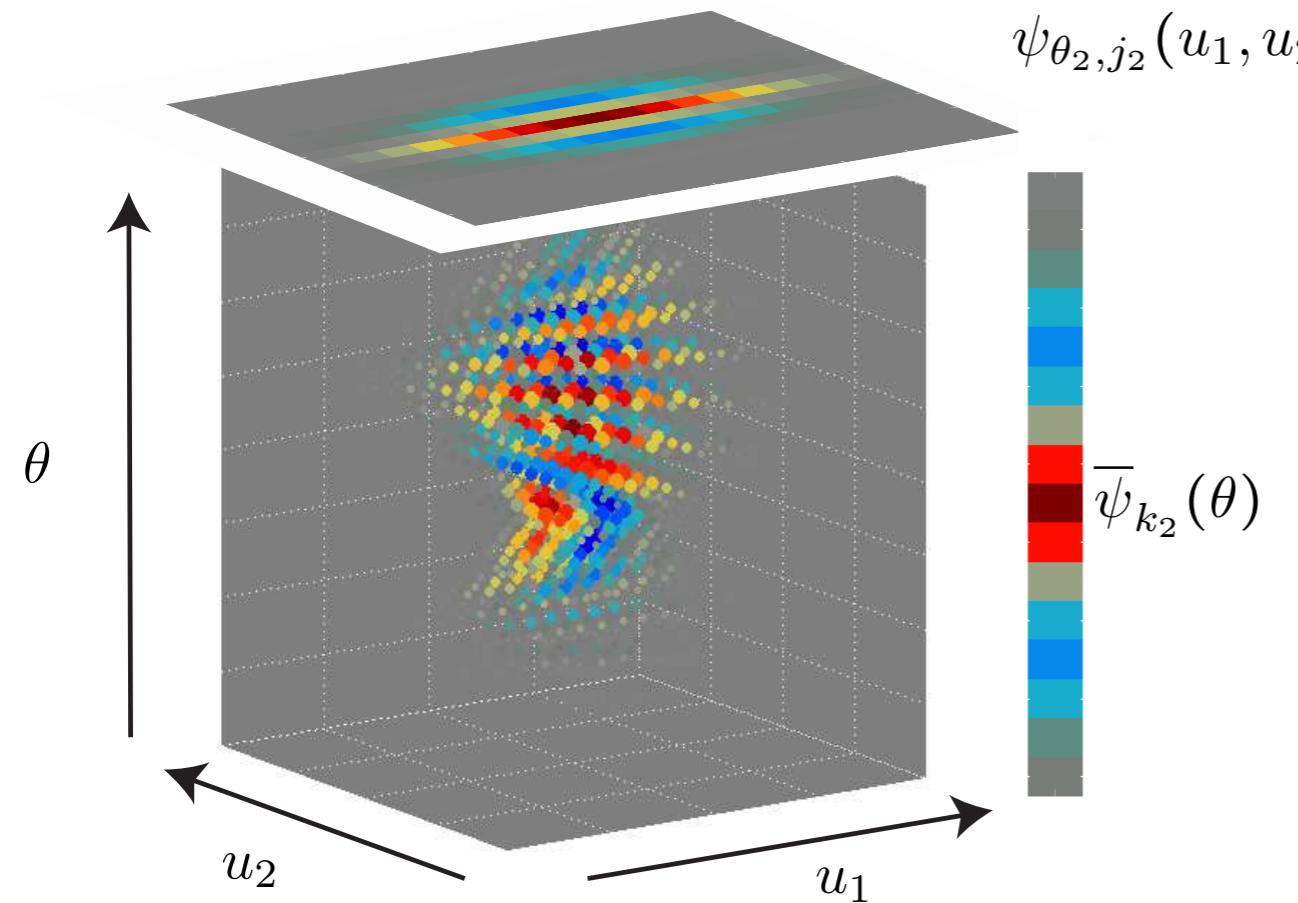
- We start by *lifting* the image with spatial wavelet convolutions: stable and covariant to roto-translations.



- We then adapt the second wavelet operator to its joint variability structure.
- More discriminability.
- Requires defining wavelets on more complicated domains:
 $\Omega = \mathbb{R}^2 \times S^2$.

EXAMPLE: ROTO-TRANSLATION SCATTERING

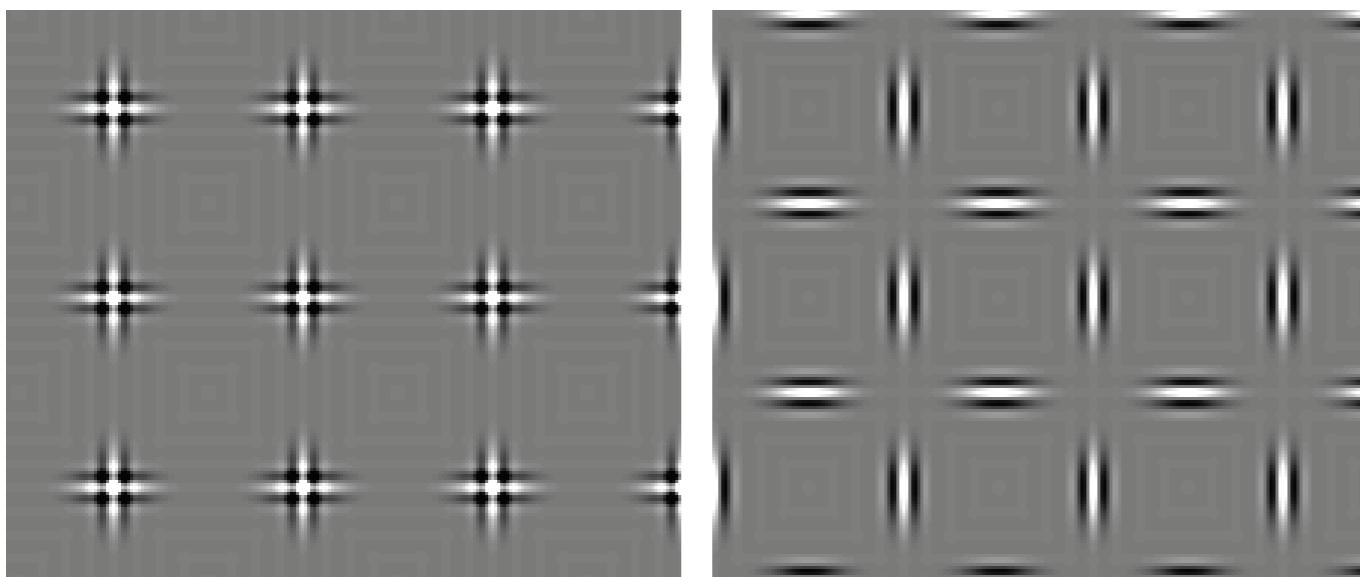
- [Sifre and Mallat'13]



$$\psi_{\theta_2, j_2}(u_1, u_2)$$

second layer wavelets constructed by a separable product on spatial and rotational wavelets:

$$\Psi_\lambda(u, \theta) = \psi_{\lambda_1}(u)\psi_{\lambda_2}(\theta)$$



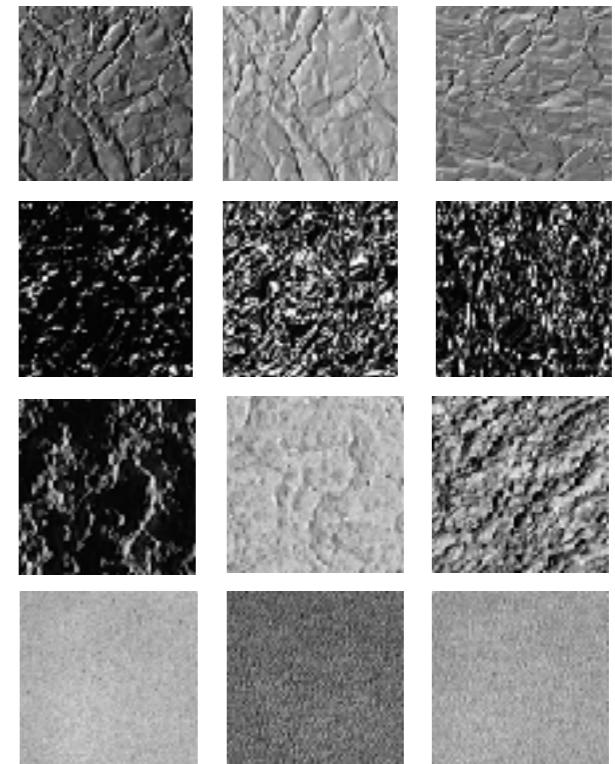
example of patterns that are discriminated by joint scattering but not with separable scattering.

CLASSIFICATION WITH SCATTERING

- State-of-the art on pattern and texture recognition using separable scattering followed by SVM:

- MNIST, USPS [Pami'13]

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

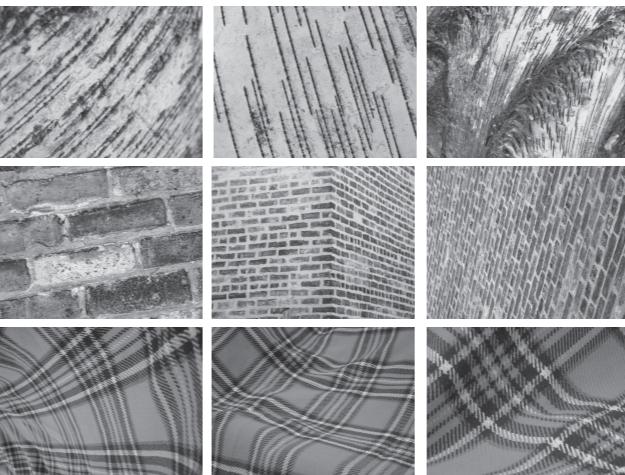


- Texture (CUREt) [Pami'13]

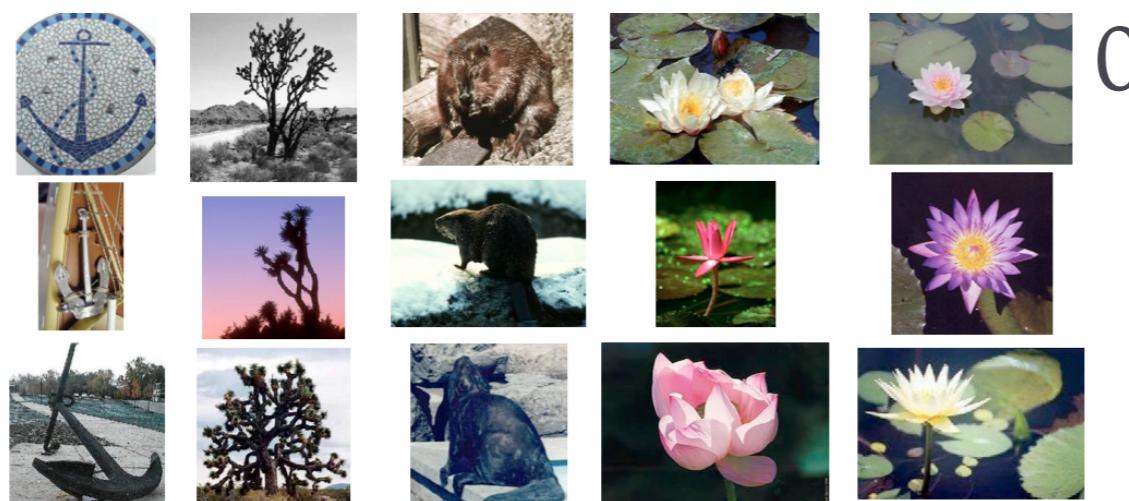
- Music Genre Classification (GTZAN) [IEEE Acoustic '13]

CLASSIFICATION WITH SCATTERING

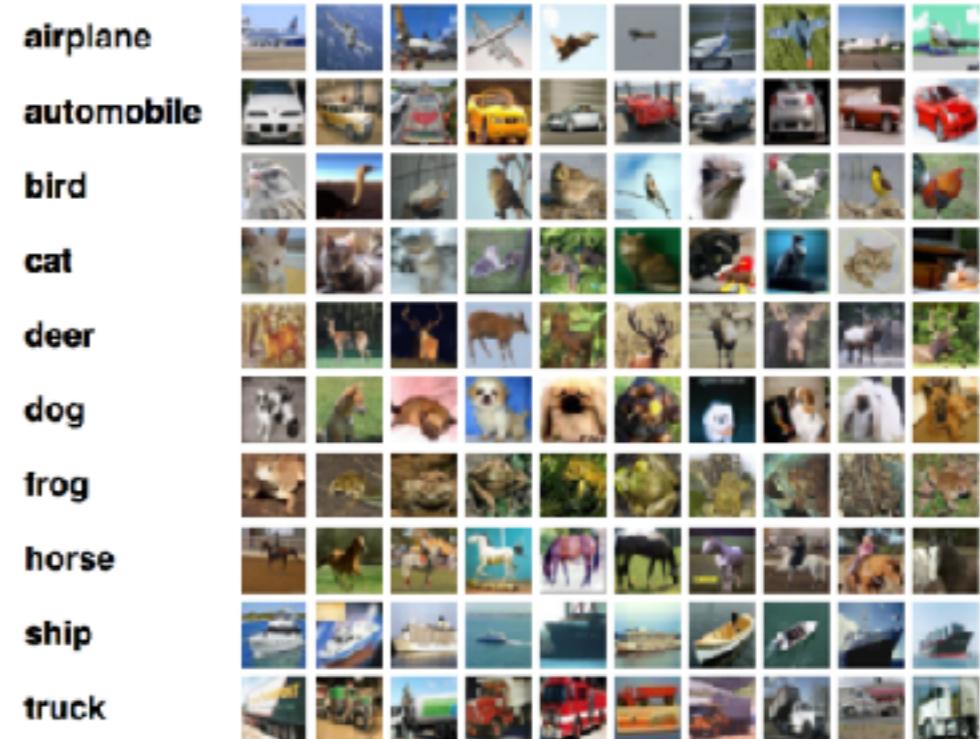
- Joint Scattering Improves Performance:
 - More complicated Texture (KTH,UIUC,UMD)
[Sifre&Mallat, CVPR'13]



- Small-mid scale Object Recognition (Caltech, CIFAR)
[Oyallon&Mallat, CVPR'15]



0



LIMITATIONS OF JOINT SCATTERING

- Variability from physical world expressed in the language of transformation groups and deformations
 - However, there are not many possible groups: essentially the affine group and its subgroups.
- As a new wavelet layer is introduced, we create new coordinates, but we do not destroy existing coordinates
 - Hard to scale: dimensionality reduction is needed.
 - Wavelet design complicated beyond roto-translation groups.
- Beyond physics, many deformations are class-specific and not small.
 - Learning filters from data rather than designing them.

FROM SCATTERING TO CNNS

- Given $x(u, \lambda)$ and a group G acting on both u and λ , we defined wavelet convolutions over G as

$$x \star_G \psi_{\lambda'}(u, \lambda) = \int_v \int_\alpha \psi_\lambda(R_{-\alpha}(u - v)) x(v, \alpha) dv d\alpha$$

- In discrete coordinates,

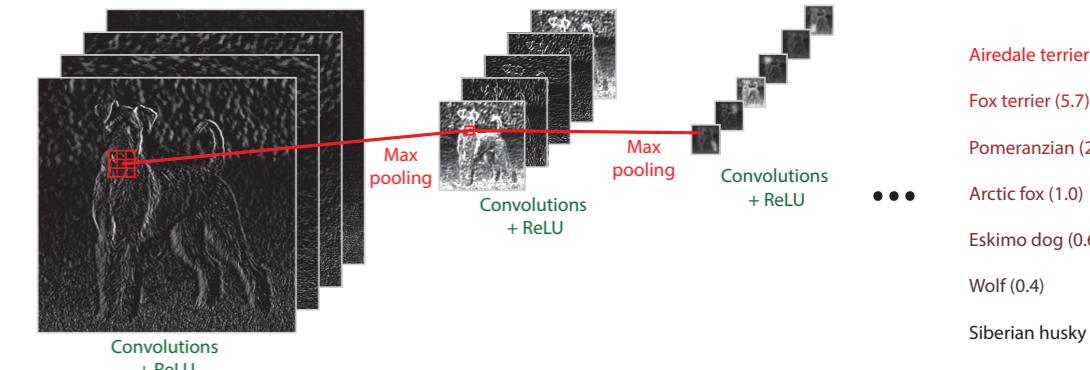
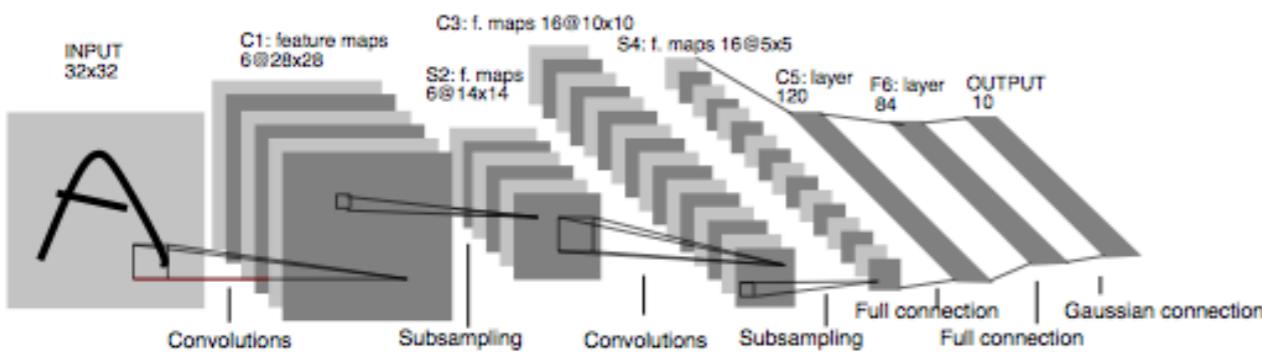
$$x \star_G \psi_{\lambda'}(u, \lambda) = \sum_v \sum_\alpha \bar{\psi}_{\lambda'}(u - v, \alpha, \lambda) x(v, \alpha)$$

- Which in general is a convolutional tensor.

CONVOLUTIONAL NEURAL NETWORKS

[LeCun, 80s,90s]

- Stack multiple layers of **localized convolutional operators** and point-wise contractive non-linearities:



Input: $x \in L^2(\Omega, \mathbb{R}^p)$.

$$\tilde{x}_{\tilde{j}}(u) = \rho \left(\sum_{j=1}^p x_j \star \theta_{j,\tilde{j}}(u) \right), \quad \tilde{j} \leq \tilde{p}.$$

Output: $\tilde{x} \in L^2(\Omega, \mathbb{R}^{\tilde{p}})$.

$\rho(z)$: point-wise nonlinearity
(e.g. $\max(0, z)$).

$\Theta = (\theta_{j,\tilde{j}})$: localized convolutional kernel.

- Down-sampling via *pooling* (can be either linear with average, or nonlinear with max) in invariant tasks:

$$\bar{x}_{\tilde{j}}(\bar{u}) = \|\tilde{x}_{\tilde{j}}(\mathcal{N}(u))\|$$

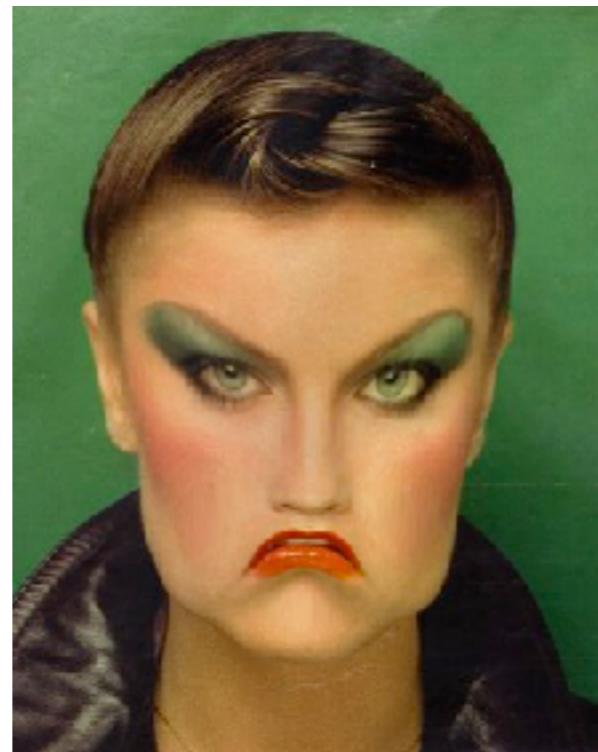
$\mathcal{N}(u)$: Neighborhood of u .

CONVOLUTIONAL NEURAL NETWORKS

- Why are CNNs geometrically stable?



$x(u)$



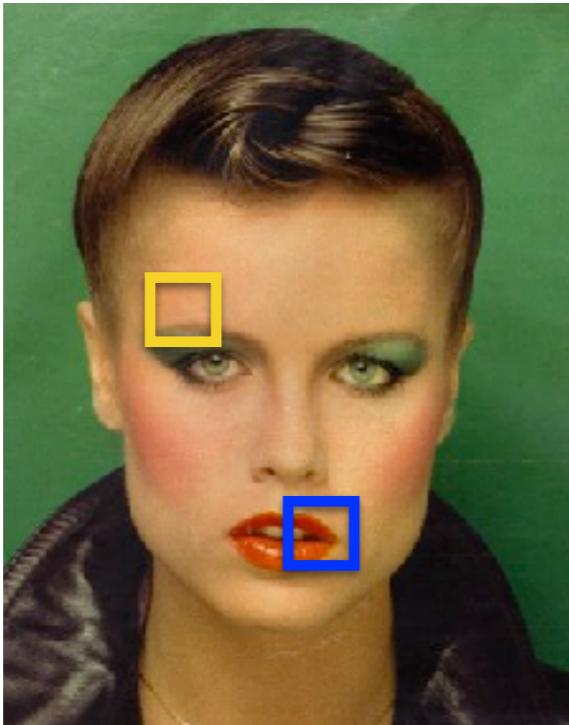
$x_\tau(u)$



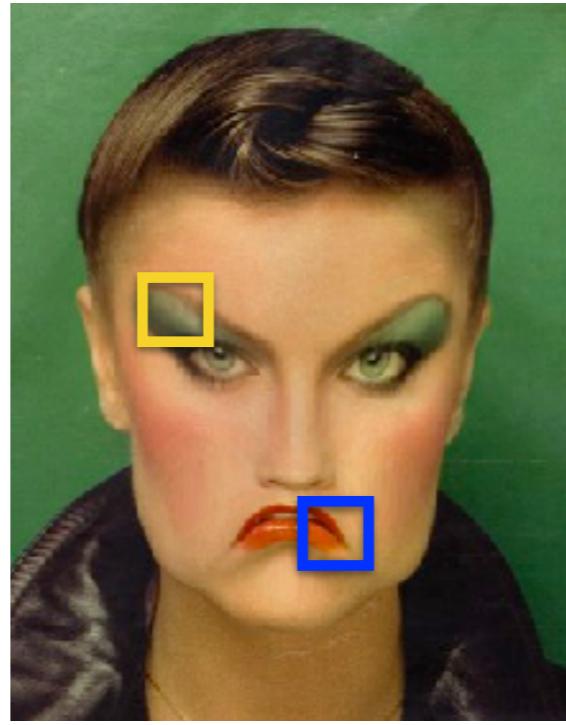
$x_{\tau'}(u)$

CONVOLUTIONAL NEURAL NETWORKS

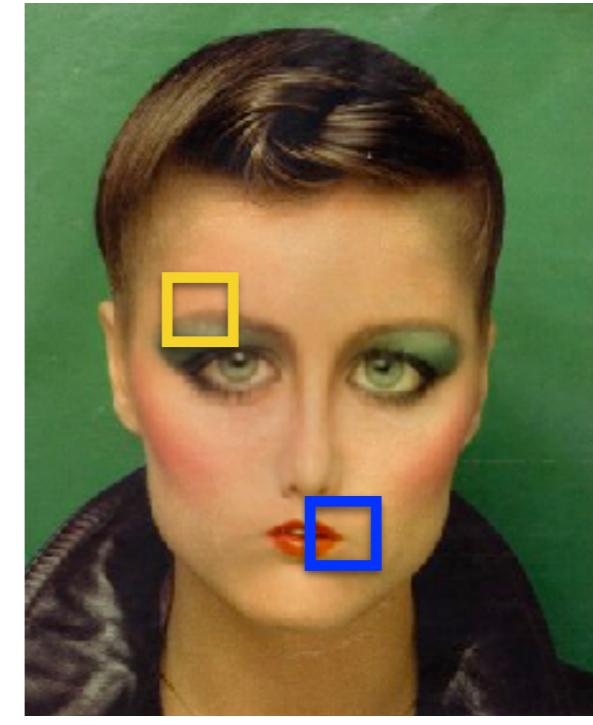
- Why are CNNs geometrically stable?



$x(u)$



$x_\tau(u)$



$x_{\tau'}(u)$

- A non-rigid deformation locally looks like a translation if $\|\nabla \tau\|$ small:

$$\Rightarrow x_\tau \star \theta(u) \approx [x \star \theta]_\tau(u)$$

- A point-wise nonlinearity commutes with deformations:

$$\Rightarrow \rho(x_\tau \star \theta(u)) \approx \rho([x \star \theta]_\tau(u)) = [\rho(x \star \theta)]_\tau(u)$$

- Pooling progressively creates invariance to geometric deformations:

$$\|x_\tau(\mathcal{N}(u))\| \approx \|x(\mathcal{N}(u))\| \text{ if } |\tau| \text{ small}$$

CONVOLUTIONAL NEURAL NETWORKS



- **Convolutions** to exploit translation invariance/equivariance.
- **Localized** to exploit geometric stability: leads to multi scale architecture.
- These two properties lead to models with $O(\log N)$ trainable parameters.
- Stability is only part of the story. Discriminability via learning/optimization is another major component for success.

INVARIANCE, LINEARIZATION AND GEODESICS

- We related stability with the ability to linearize deformations:

$$\begin{aligned}\tau \mapsto \Phi(\varphi_\tau x) \text{ Lipschitz} \Rightarrow \\ \Phi(\varphi_\tau x) = \Phi(x) + D(\Phi \circ \varphi_\cdot(x))\tau + O(\|\tau\|)\end{aligned}$$

INVARIANCE, LINEARIZATION AND GEODESICS

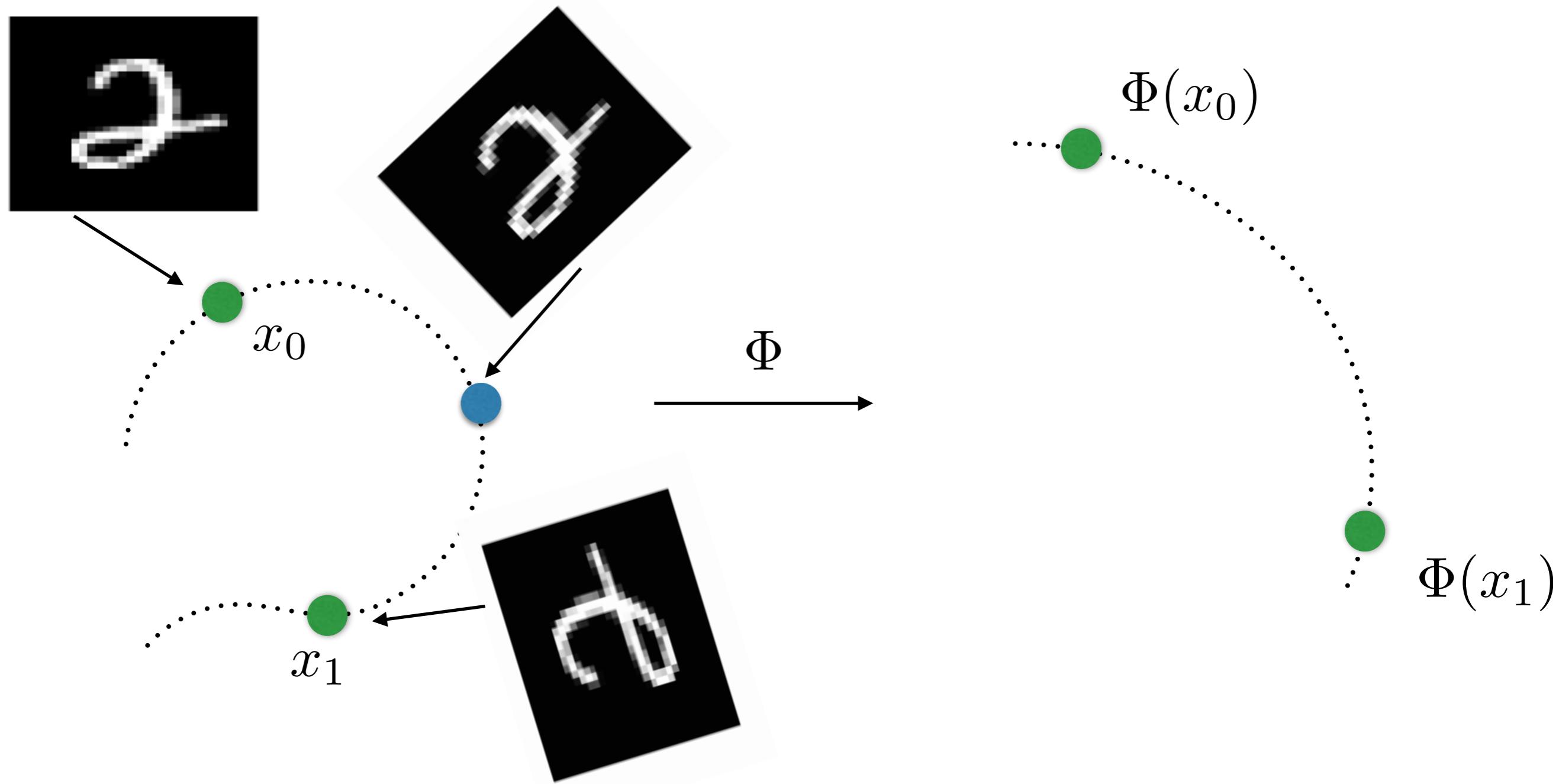
- We related stability with the ability to linearize deformations:

$$\begin{aligned}\tau \mapsto \Phi(\varphi_\tau x) \text{ Lipschitz} \Rightarrow \\ \Phi(\varphi_\tau x) = \Phi(x) + D(\Phi \circ \varphi_\cdot(x))\tau + O(\|\tau\|)\end{aligned}$$

- One can test this property over learnt representations by inspecting geodesics.
 - They become linear paths in feature space under the metric
$$d(x, x') = \|\Phi(x) - \Phi(x')\|$$
- [Bengio et al. '11], [Goroshin et al'15], [Henaff et al '16]

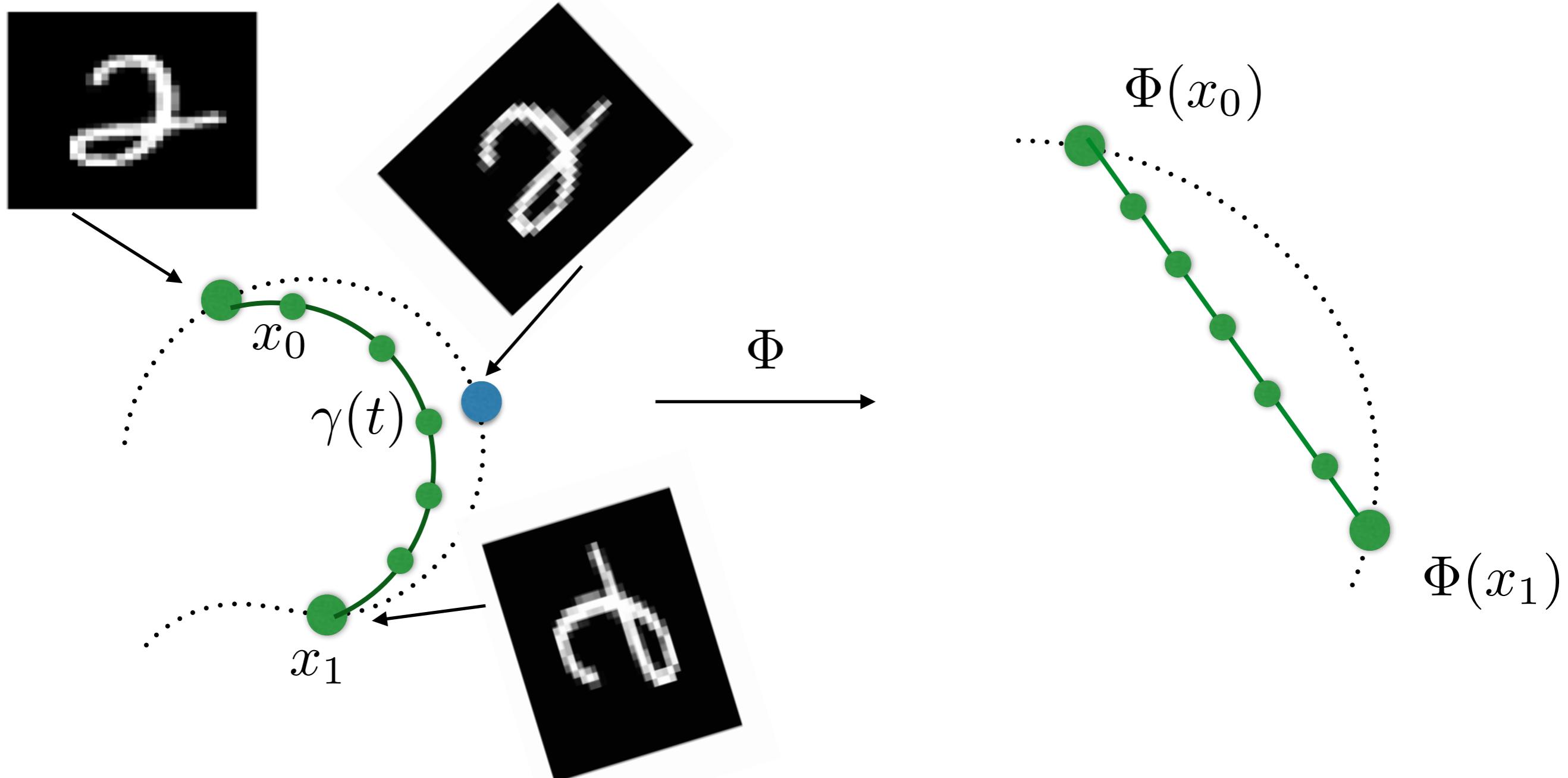
INVARIANCE, LINEARIZATION AND GEODESICS

- Algorithm from [Henaff & Simoncelli '16]:



INVARIANCE, LINEARIZATION AND GEODESICS

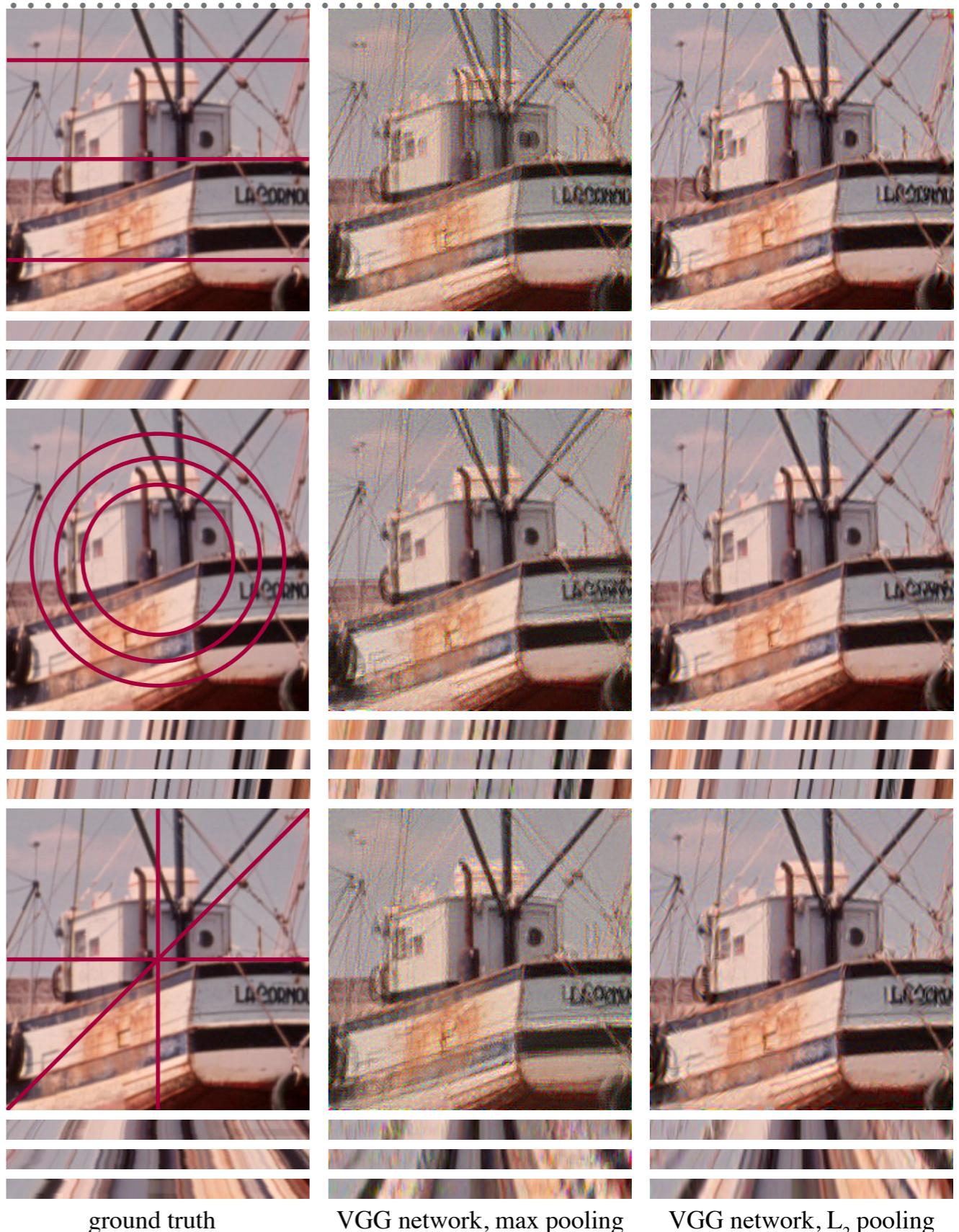
► Algorithm from [Henaff & Simoncelli '16]:



$$\min_{\gamma(0)=x_0, \gamma(1)=x_1} \int_0^1 |\dot{\gamma}(t)| dt + \int_0^1 |(\dot{\Phi}\gamma)(t)| dt$$

INVARIANCE, LINEARIZATION AND GEODESICS

- On pertained CNNs (VGG oxford net), linearization is empirically verified for various groups.
- Continuous transformation groups are better linearized with energy pooling than with max-pooling



[Henaff and Simoncelli'16]

REDUNDANCY IN CNNS

$$\Phi(x) = \rho(\dots \rho(x * \Psi_1) * \dots * \Psi_k))$$

- Large-scale networks contain > 10 layers and $> 10^6$ parameters.
- Q: Is there a smaller parametric model that contains good representations?

REDUNDANCY IN CNNS

- “Post-training” model compression:
Given parameters $\Theta = (\Theta_1, \dots, \Theta_k)$, find a
reparametrization $\tilde{\Phi}$ such that $\mathbb{E}\|\Phi(x; \Theta) - \tilde{\Phi}(x)\|$ is small.

REDUNDANCY

- “Post-training” model compression:
Given parameters $\Theta = (\Theta_1, \dots, \Theta_k)$, find a
reparametrization $\tilde{\Phi}$ such that $\mathbb{E}\|\Phi(x; \Theta) - \tilde{\Phi}(x)\|$ is small.
- Useful to accelerate evaluation of large networks ([Denton et al,’14], [Jaderberg et al’14]) (“Optimal Brain Damage” [LeCun et al,’90]) $\tilde{\Phi}(x) = \Phi(x, \tilde{\Theta})$, $\tilde{\Theta}_i = F(\beta_i)$
- Typically we restrict the new class to be $\dim(\beta_i) \ll \dim(\Theta_i)$
- Explore low-rank tensor factorizations of each convolutional tensor.

REDUNDANCY

- “Post-training” model compression:
Given parameters $\Theta = (\Theta_1, \dots, \Theta_k)$, find a reparametrization $\tilde{\Phi}$ such that $\mathbb{E}\|\Phi(x; \Theta) - \tilde{\Phi}(x)\|$ is small.
 - Useful to accelerate evaluation of large networks ([Denton et al,’14], [Jaderberg et al’14]) (“Optimal Brain Damage” [LeCun et al,’90])
 - Typically we restrict the new class to be
$$\tilde{\Phi}(x) = \Phi(x, \tilde{\Theta}) , \quad \tilde{\Theta}_i = F(\beta_i) \qquad \dim(\beta_i) \ll \dim(\Theta_i)$$
 - Explore low-rank tensor factorizations of each convolutional tensor.
- “Pre-training” model compression:
 - Train directly in the compressed domain ([“Predicting parameters in Deep Learning”, Denil et al,’13]).
 - Mild regularization effect. *Interplay between statistical performance and optimization performance.*

INVERTIBILITY: NO TRAINING AND NO STRUCTURE

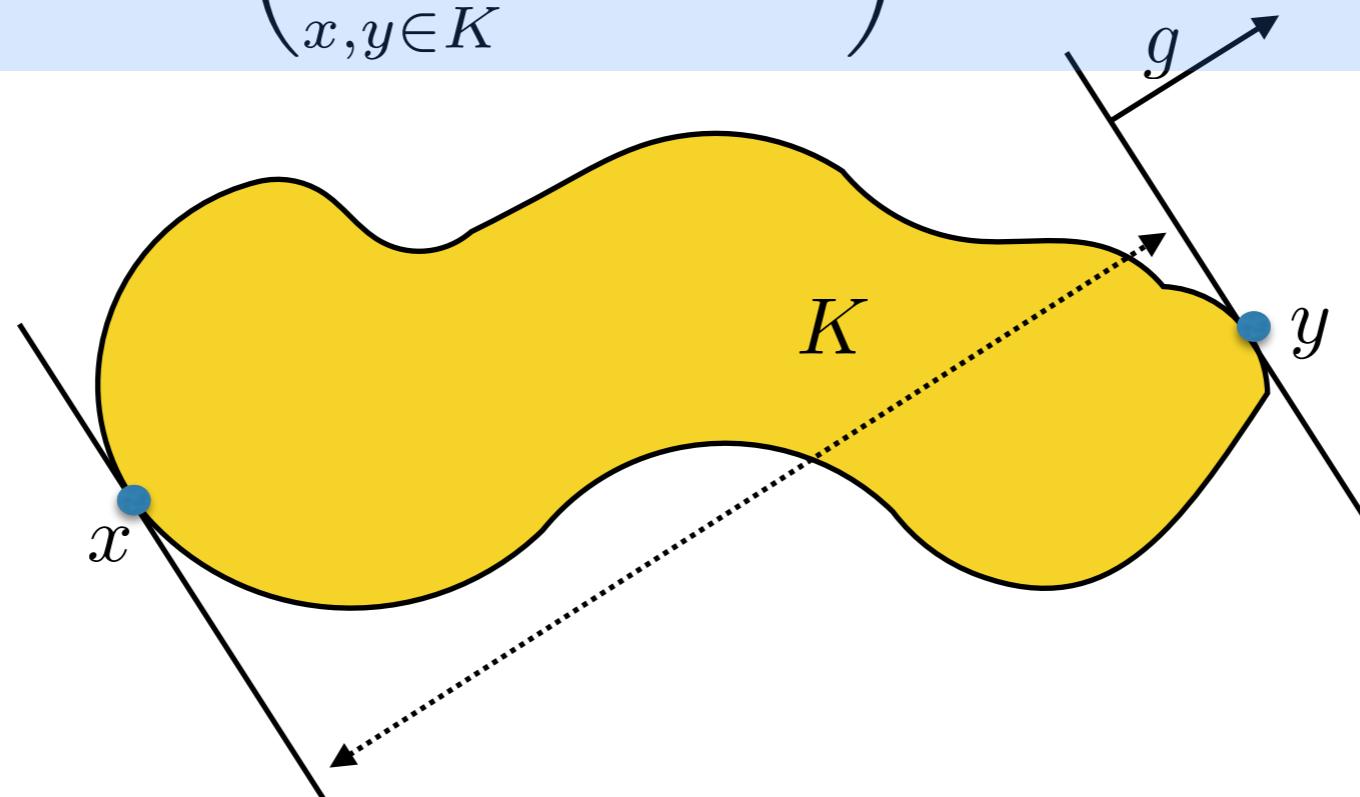
- [Giryes, Sapiro and Bronstein, '15] $\Phi = \text{Random Convnet}$

INVERTIBILITY: NO TRAINING AND NO STRUCTURE

- [Giryes, Sapiro and Bronstein, '15]
Gaussian mean width of a set K :

$\Phi = \text{Random Convnet}$

$$\omega(K) := \mathbb{E} \left(\sup_{x,y \in K} \langle g, x - y \rangle \right), \quad g \sim \mathcal{N}(0, \mathbf{I})$$

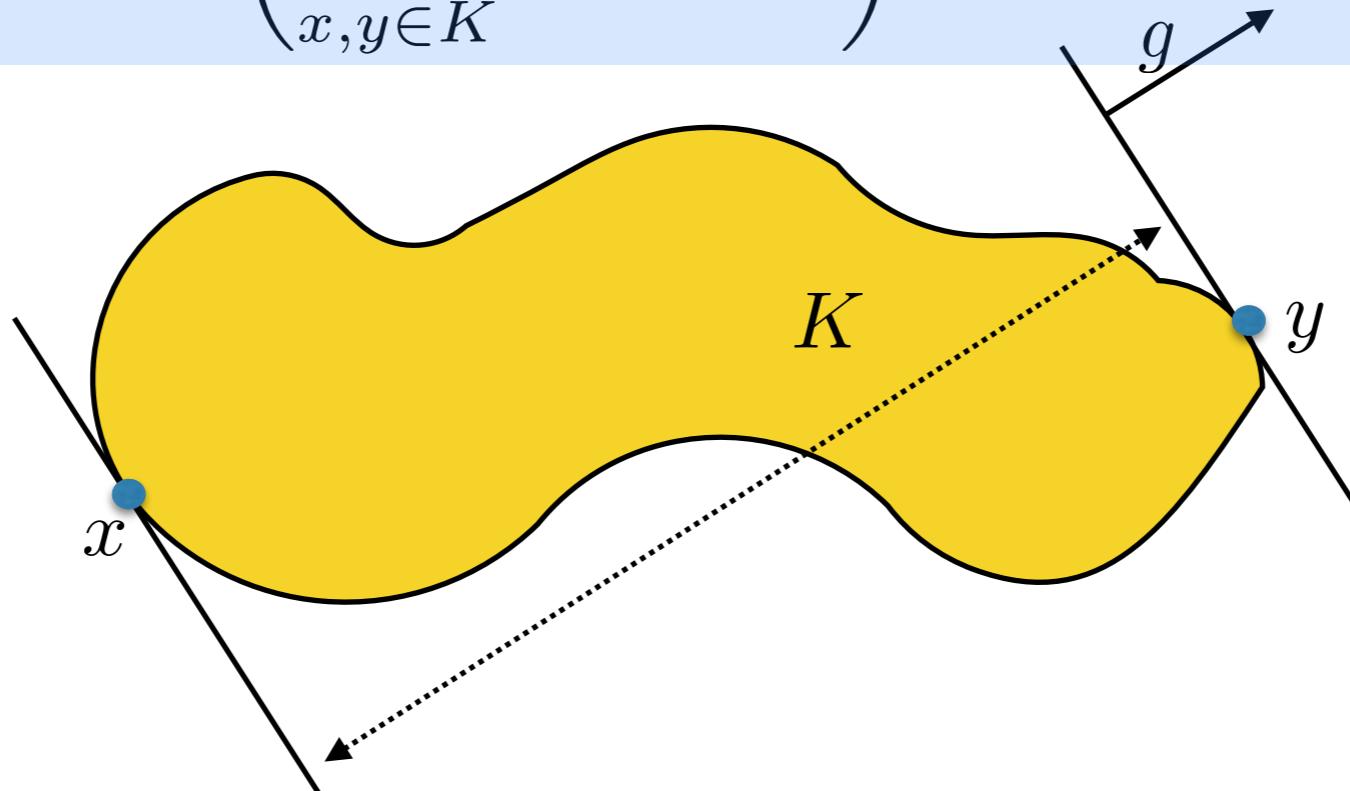


INVERTIBILITY: NO TRAINING AND NO STRUCTURE

► [Giryes, Sapiro and Bronstein, '15] Gaussian mean width of a set K :

$\Phi = \text{Random Convnet}$

$$\omega(K) := \mathbb{E} \left(\sup_{x,y \in K} \langle g, x - y \rangle \right), \quad g \sim \mathcal{N}(0, \mathbf{I})$$



► K : mixture of L gaussians of dimension k : $\omega(K) = O(\sqrt{k + \log L})$.

► Proxy for the dimensionality of a set.

► K : k -sparse signals in a dictionary of size L : $\omega(K) = O(\sqrt{k \log(L/k)})$.

INVERTIBILITY: NO TRAINING AND NO STRUCTURE

Theorem [GSB'15]: Let $\rho(\cdot)$ be the ReLU and $K \subset \mathbb{B}_1^n$ the dataset. If $\sqrt{m}W \in \mathbb{R}^{m \times n}$ is a random matrix with iid normally distributed entries and $m \geq C\delta^{-4}\omega(K)^4$ then with high probability

$$|\|\rho(Wx) - \rho(Wy)\|^2 - (0.5\|x - y\|^2 + \|x\|\|y\|\beta(x, y))| \leq \delta.$$

Moreover, if K is sufficiently away from 0, there exists $C > 0$ such that whp

$$|\cos \angle(\rho(Wx), \rho(Wy)) - \cos(\angle(x, y)) - \beta(x, y)| \leq C\delta.$$

$$\angle(x, y) = \cos^{-1} \left(\frac{x^T y}{\|x\| \|y\|} \right) \text{ angle between } x \text{ and } y$$

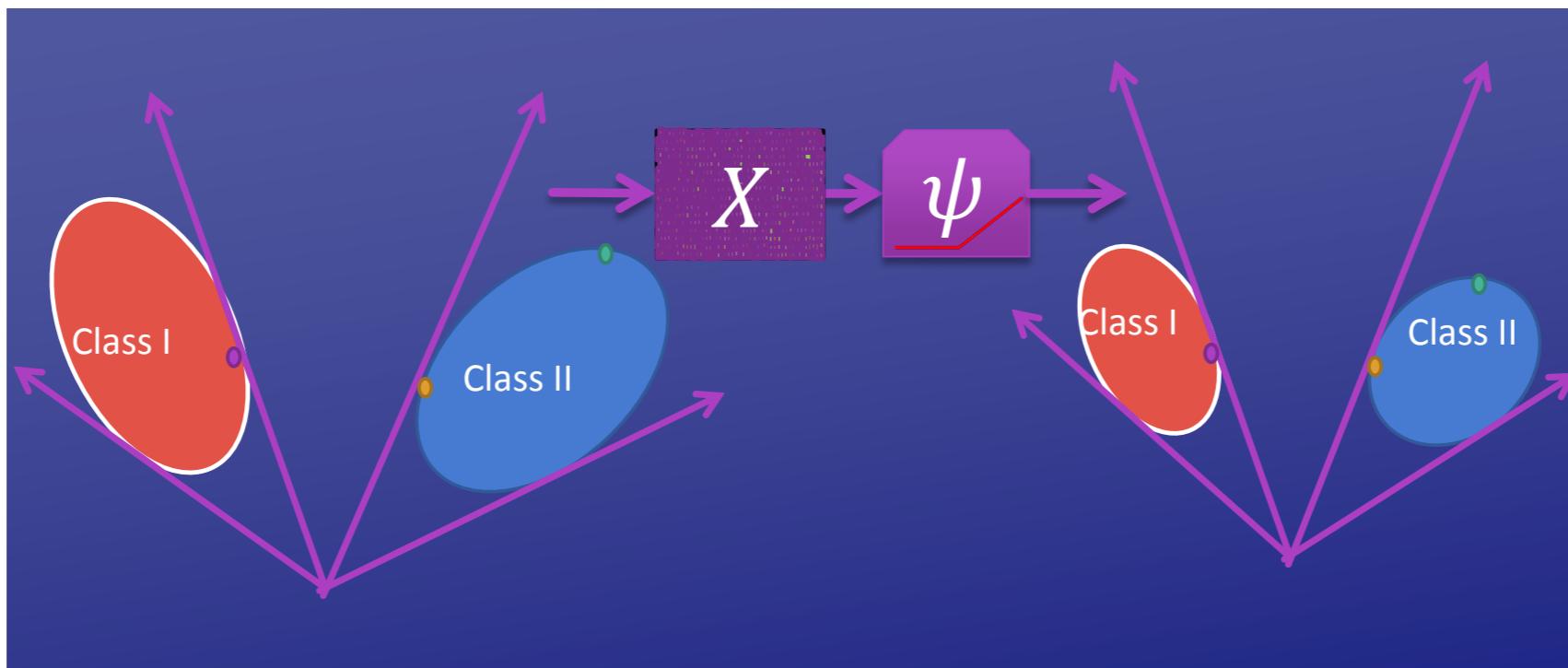
$$\beta(x, y) = \pi^{-1}(\sin(\angle(x, y)) - \angle(x, y) \cos(\angle(x, y)))$$

INTERPRETATION

- If $\angle(x, y)$ is small, then $\beta(x, y) \approx 0$:
distances are approx. shrunk by 2, angles are preserved.

INTERPRETATION

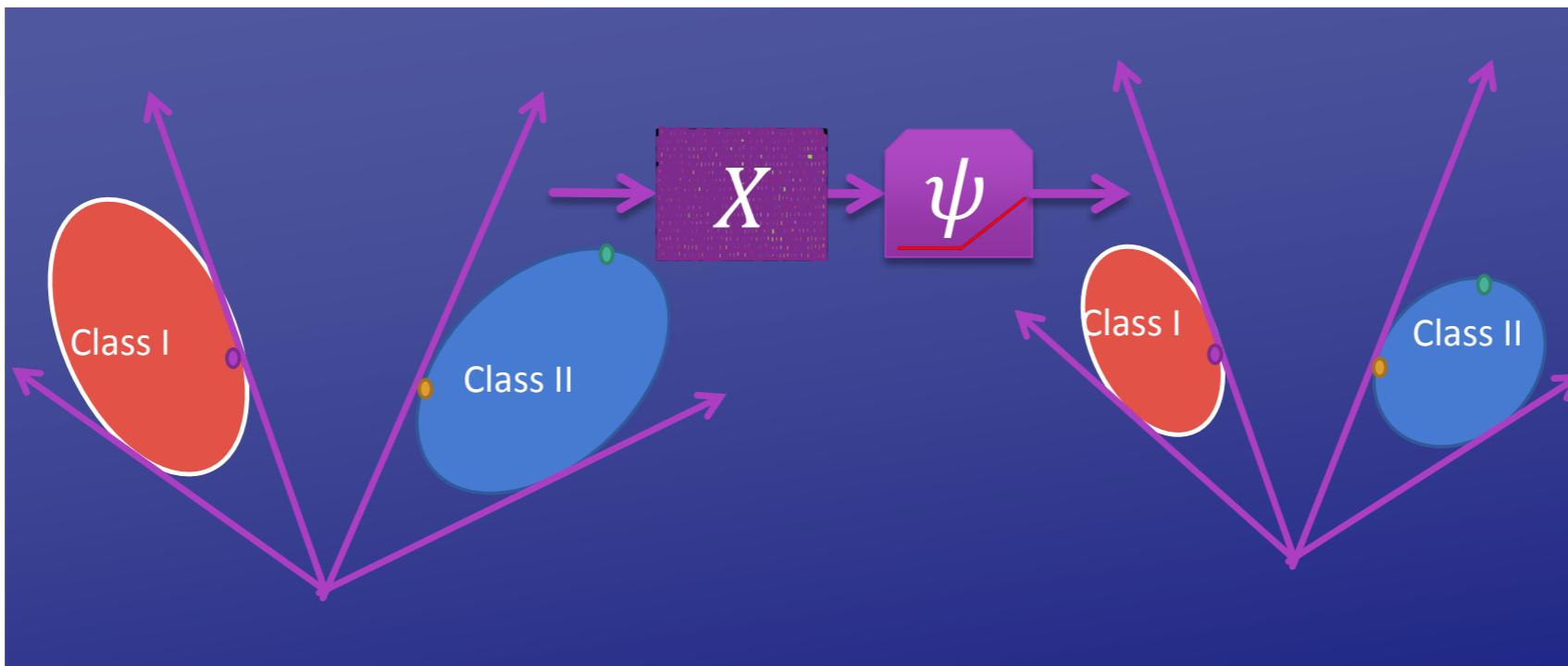
- If $\angle(x, y)$ is small, then $\beta(x, y) \approx 0$: distances are approx. shrunk by 2, angles are preserved.
- If $\angle(x, y)$ is large, then $\beta(x, y) \approx 0.5$: distances are shrunk by a smaller factor.



points with small angles between them become closer than points with larger angles between them

INTERPRETATION

- If $\angle(x, y)$ is small, then $\beta(x, y) \approx 0$: distances are approx. shrunk by 2, angles are preserved.
- If $\angle(x, y)$ is large, then $\beta(x, y) \approx 0.5$: distances are shrunk by a smaller factor.



[Raja Giryes]

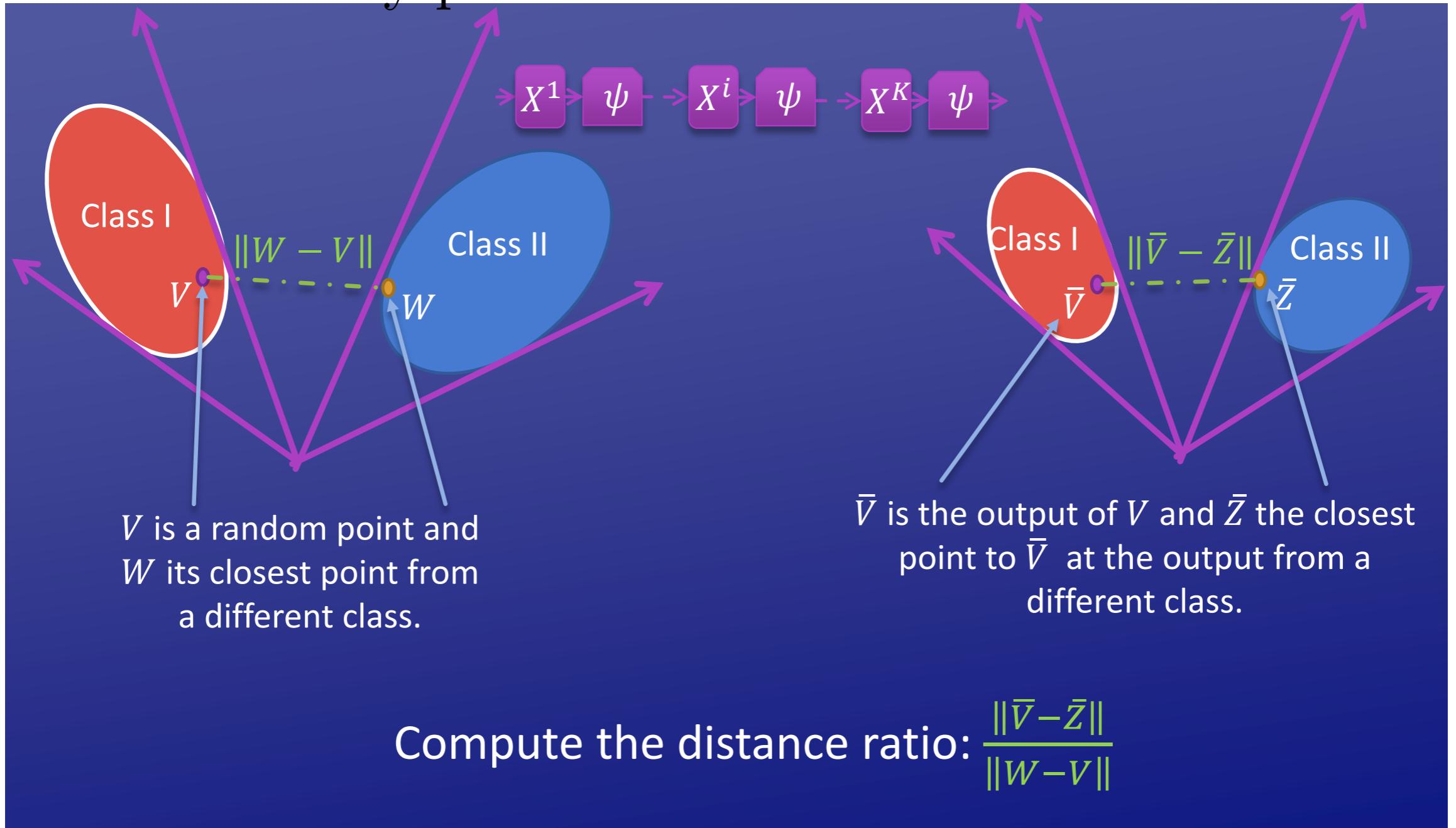
points with small angles between them become closer than points with larger angles between them

The result can be cascaded since gaussian mean width is approximately preserved by each layer.

ROLE OF TRAINING?

ROLE OF TRAINING?

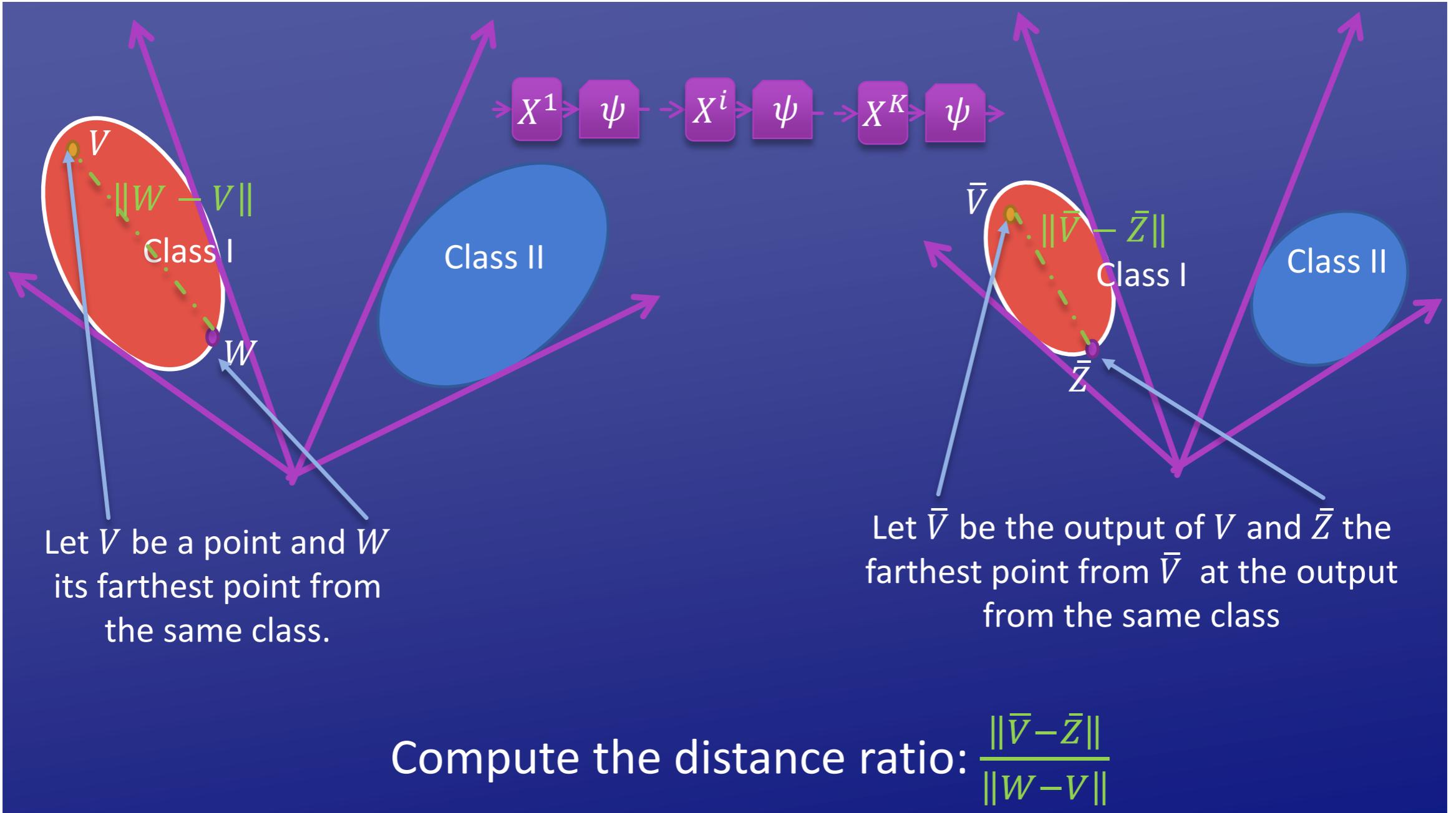
Inter Boundary points distance ratio



[Raja Giryes]

ROLE OF TRAINING?

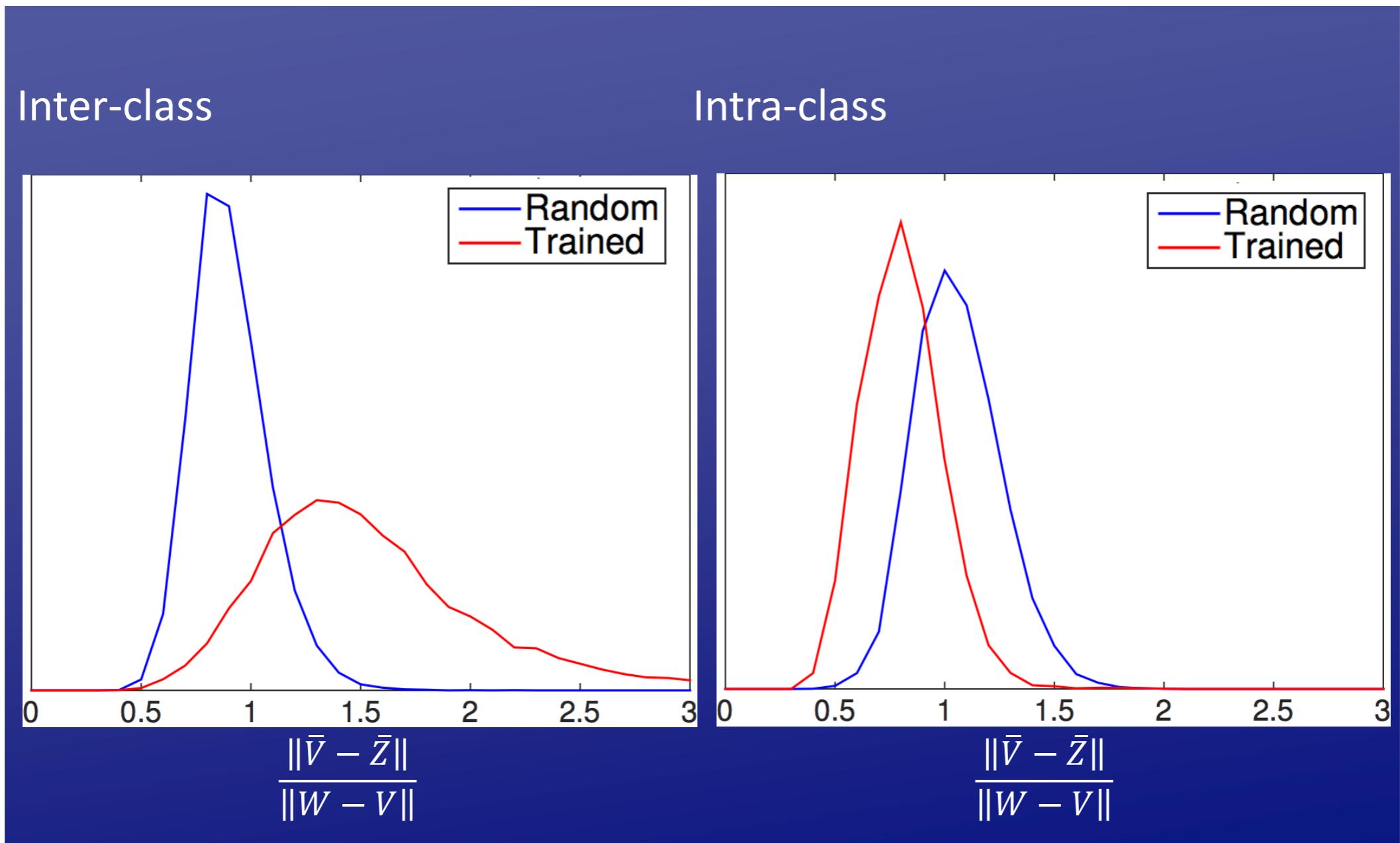
Intra Boundary points distance ratio



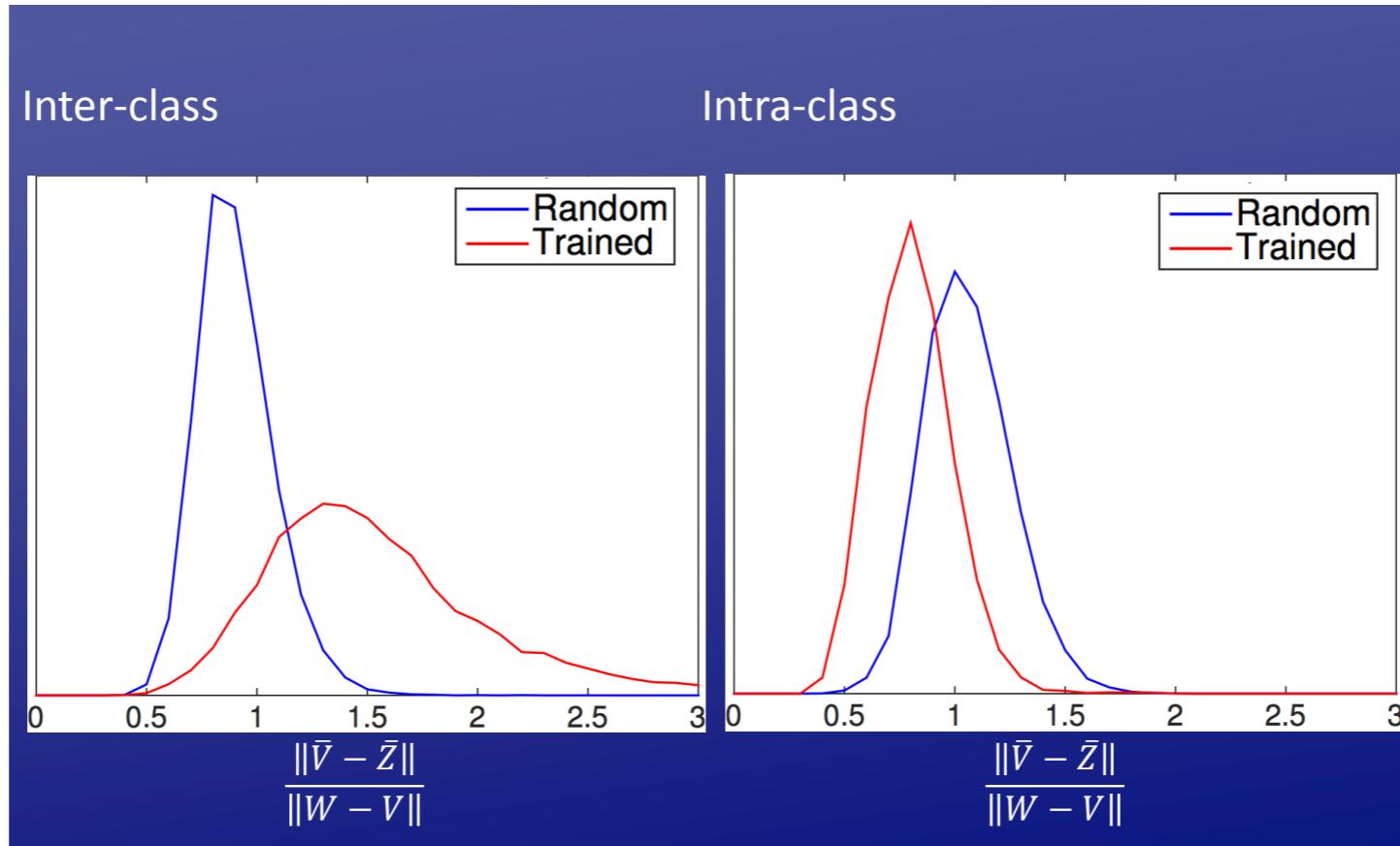
[Raja Giryes]

ROLE OF TRAINING?

Boundary distance ratios
measured on Imagenet using VGG oxfordnet



ROLE OF TRAINING?

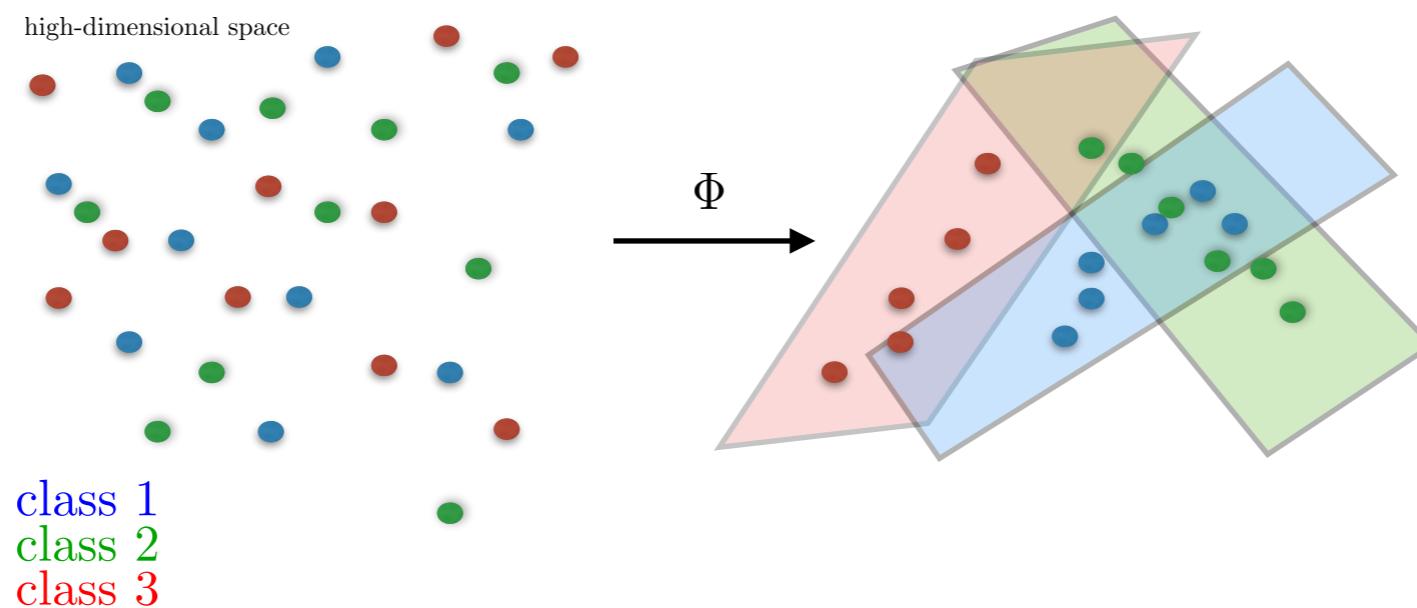


[Raja Giryes]

- Training the network does not affect the *bulk* of distances
- However, it critically changes the behavior at the boundary points:
 - Inter-class distances expand (as expected).
 - Intra-class distances shrink (as expected).

INVERTIBILITY WITH TRAINING

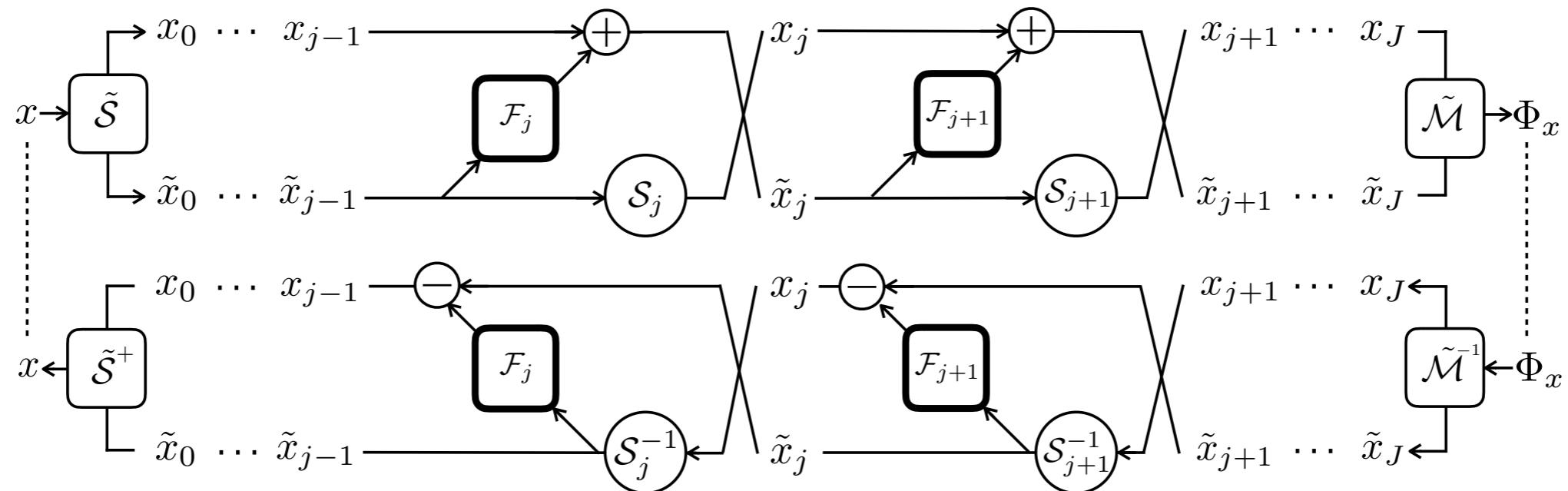
- Hypothesis: Deep CNNs progressively “linearize” intra-class variability.
$$\hat{f}(x) = \langle \Phi(x), w \rangle$$
- Hypothesis: Once it is “linearized”, it can be “killed” by linear projections that contract the space:
$$\Phi(x) = \rho W_L \rho W_{L-1} \cdots W_1 x$$



- Does this imply that the network loses information?
 - As implied by “Information Bottleneck”, Tibshy et al.

INVERTIBILITY WITH TRAINING

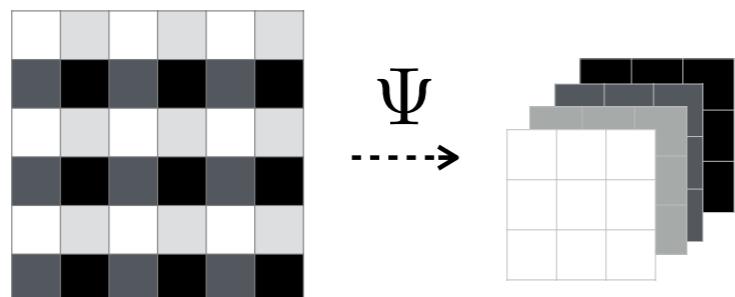
- iRevNets [Jacobsen, Smeulders, Oyallon, ICLR'18] construct provably invertible CNN representations for large-scale image classification.



- \mathcal{S}_j : linear *invertible* reshuffling operator.
- \mathcal{F}_j : non-linear convolutional net.
- Same as residual networks, but downsampling is made invertible.

INVERTIBILITY WITH TRAINING

- Invertible down-sampling operator:



$$\mathcal{S}_j x(u, \lambda) := x(\Psi(u, \lambda)) .$$

- The operator is certified invertible:

$$\begin{cases} x_{j+1} = \mathcal{S}_{j+1} \tilde{x}_j \\ \tilde{x}_{j+1} = x_j + \mathcal{F}_{j+1} \tilde{x}_j \end{cases} \iff \begin{cases} \tilde{x}_j = \mathcal{S}_{j+1}^{-1} x_{j+1} \\ x_j = \tilde{x}_{j+1} - \mathcal{F}_{j+1} \tilde{x}_j \end{cases}$$

- Related to wavelet *lifting* schemes [Sweldens'98]
- Classification performance is preserved despite invertibility:

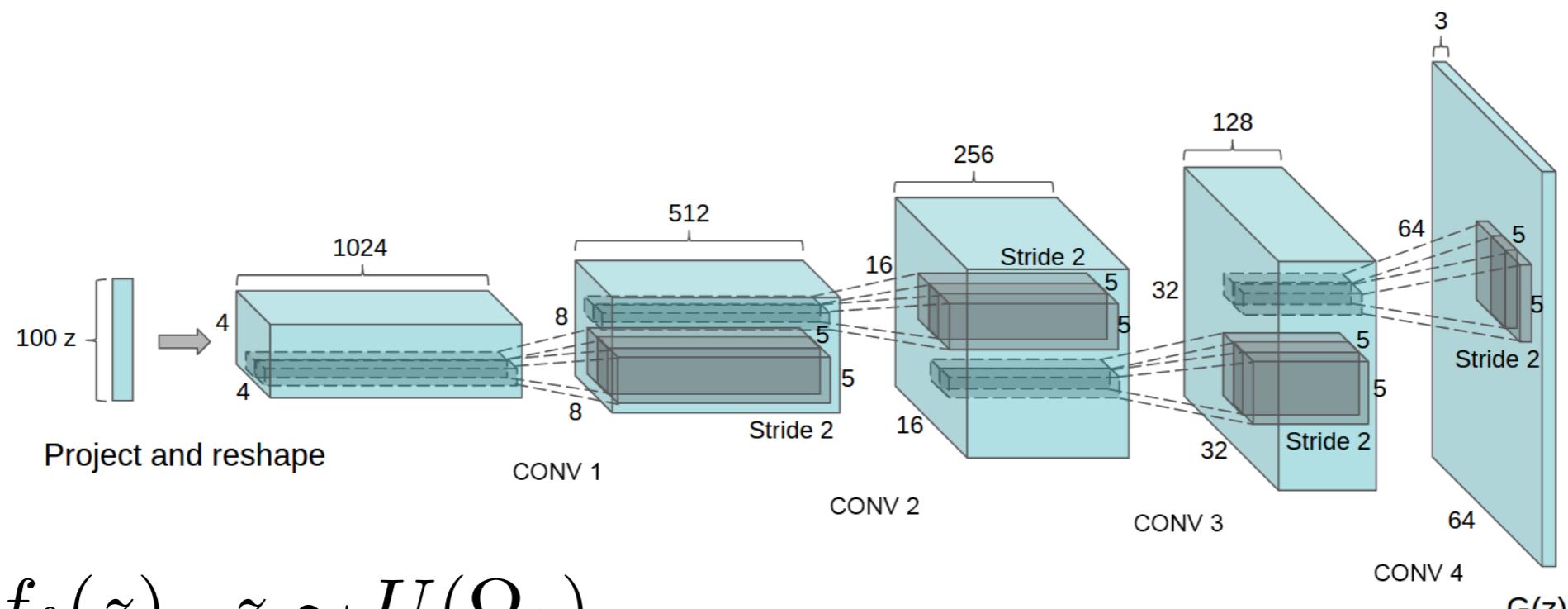
Architecture	Injective	Bijective	Top-1 error	Parameters
ResNet	-	-	24.7	26M
RevNet	-	-	25.2	28M
i -RevNet (a)	yes	-	24.7	181M
i -RevNet (b)	yes	yes	26.7	29M

DEEP IMAGE PRIOR [ULYANOV ET AL '17]

- How much prior information is captured by the choice of convolutional architecture?

DEEP IMAGE PRIOR [ULYANOV ET AL '17]

- How much prior information is captured by the choice of convolutional architecture?
- We will see next week that CNNs can be used also as *generators*:



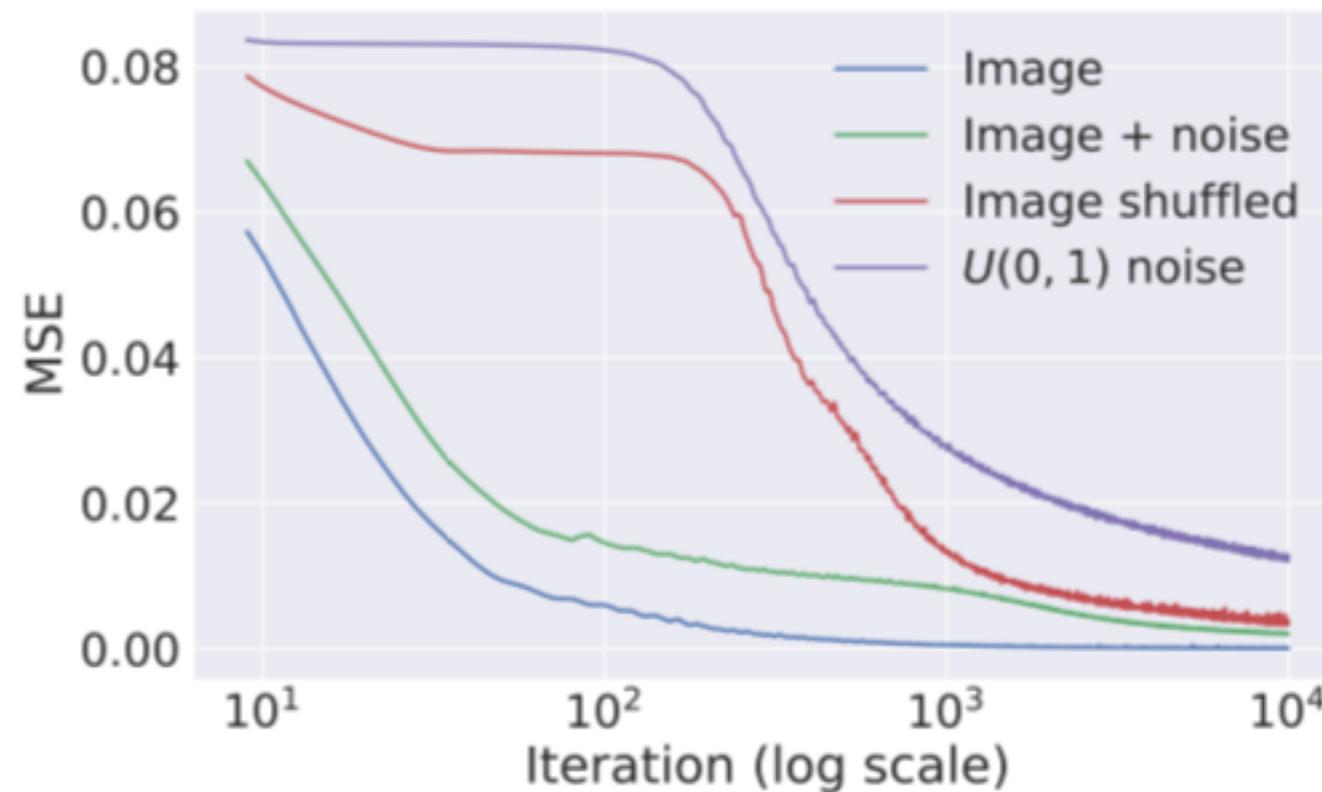
$$x = f_{\theta}(z) , z \sim U(\Omega_0) .$$

DC-Gan Architecture [Radford et al'16]

- Which $x \in L^2(\Omega)$ are well-approximated by such CNNs?

DEEP IMAGE PRIOR [ULYANOV ET AL '17]

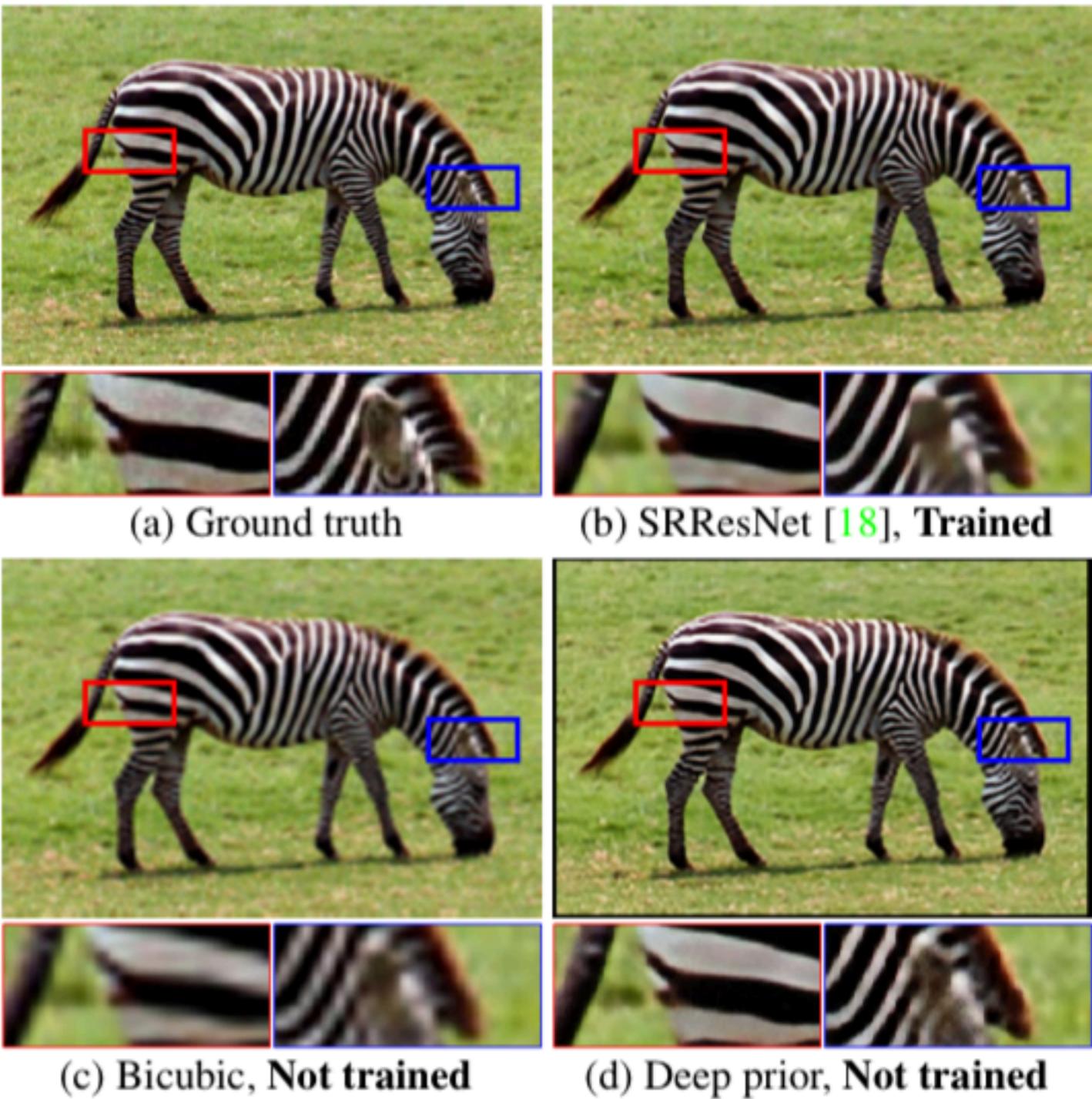
- Given $x \in L^2(\Omega)$, fix a sample $z \sim U$ and consider the optimization problem $\min_{\theta} \ell(x, f_{\theta}(z))$
- Remarkably, *natural* images are easier to approximate than “non-natural” images:



DEEP IMAGE PRIOR [ULYANOV ET AL '17]

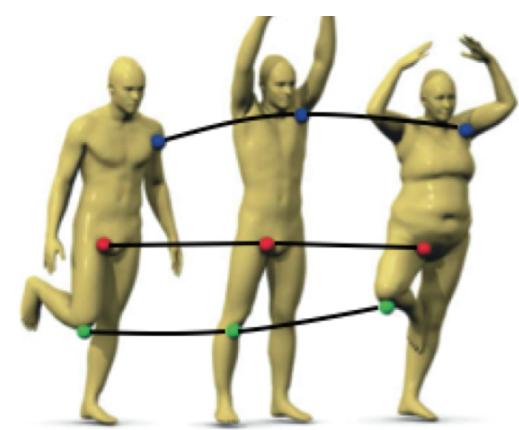
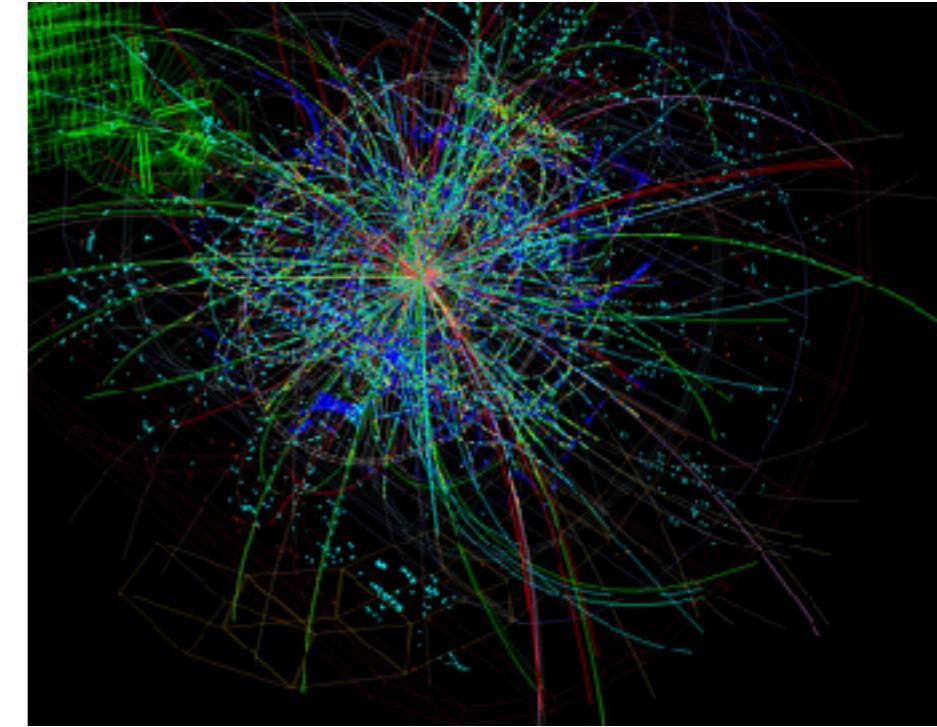
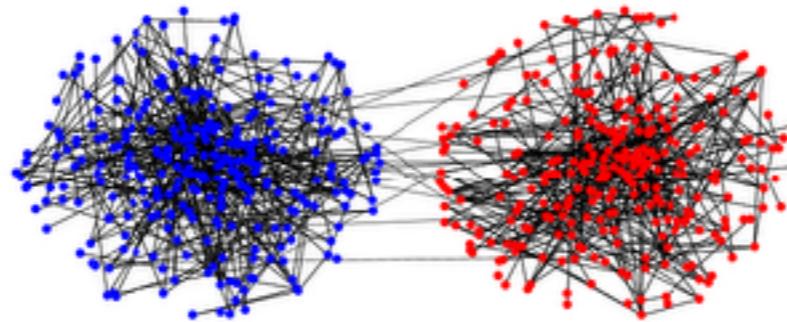
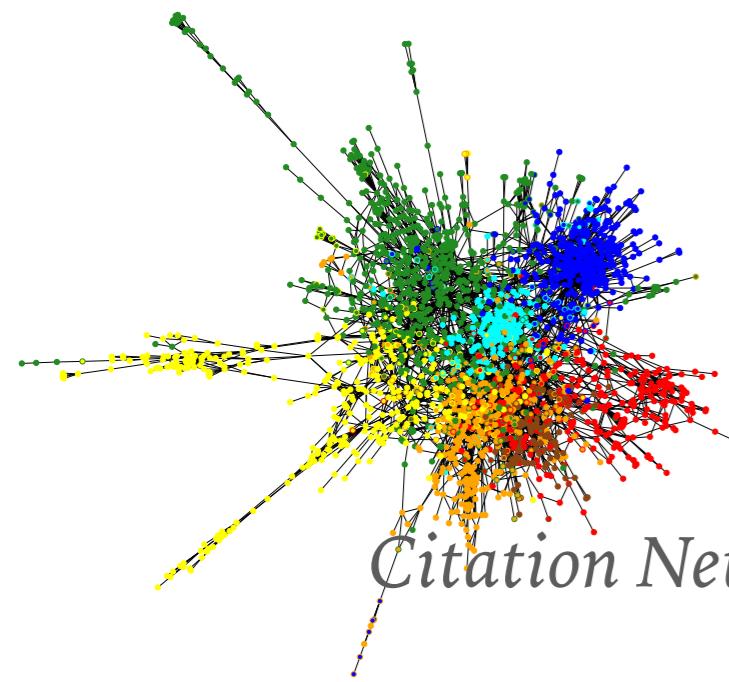
$$\hat{x} = f_{\theta^*}(z) , \theta^* \in \arg \min_{\theta} \|\Gamma f_{\theta}(z) - x_0\|^2$$

- Applications to imaging inverse problems *without learning.*
- More on that next week.

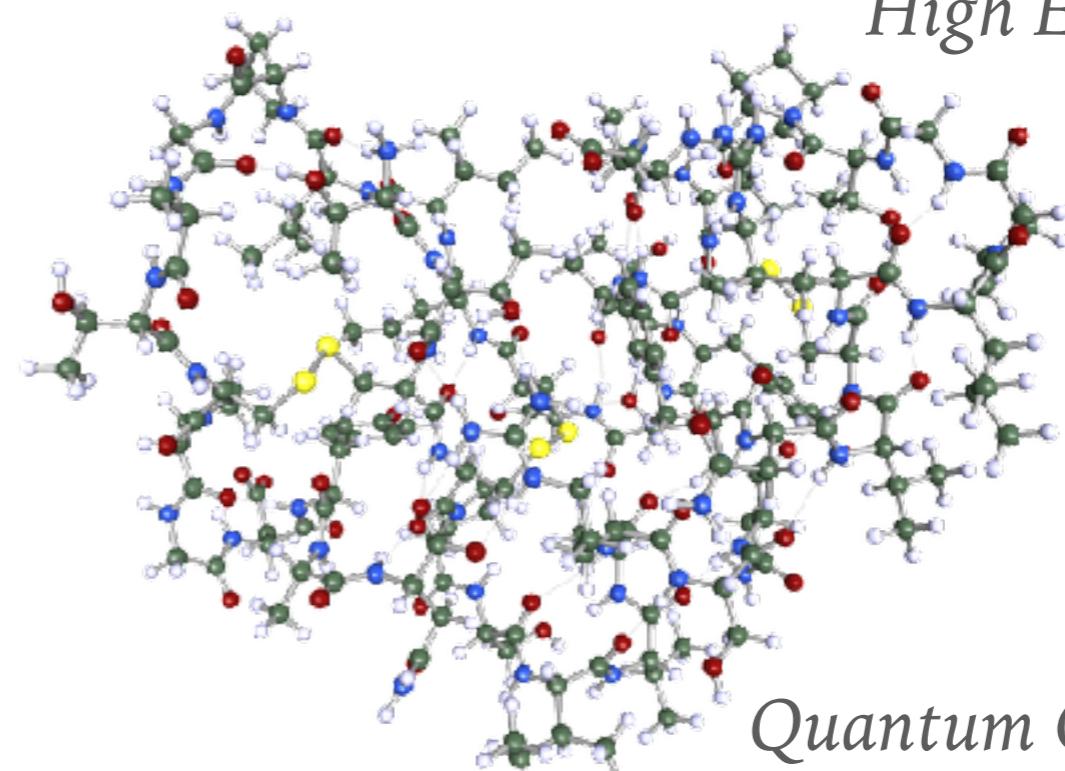


TOWARDS NON-EUCLIDEAN GEOMETRIES

- How about problems/tasks defined over more general domains?



Graphics

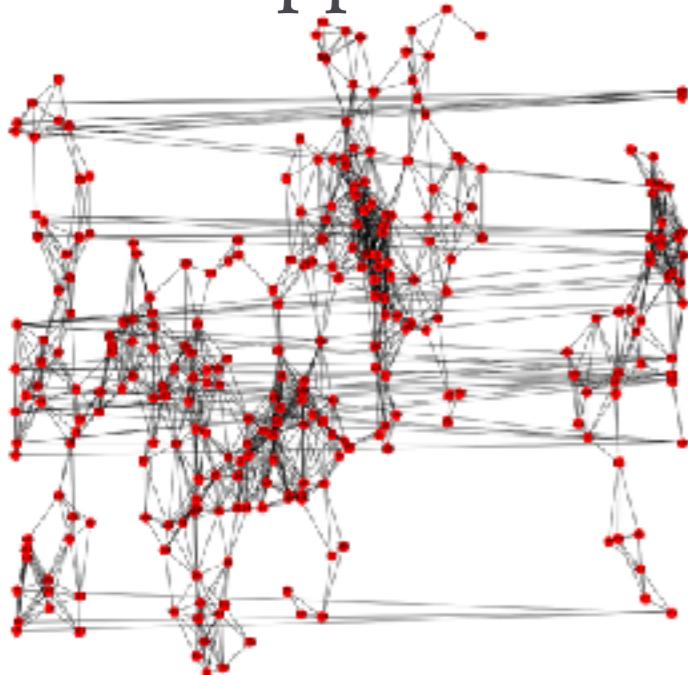


NON-EUCLIDEAN GEOMETRIC STABILITY

- We replace the Euclidean domain Ω by a general graph $G = (V, E)$.
 $x(u) \in L^2(\Omega) \rightarrow x(u) \in L^2(\textcolor{blue}{G})$, $G = (V, E)$.
 - In some applications, the input is the graph itself: $x \leftrightarrow G$
- We focus on undirected, possibly weighted graphs:
 $W \in \mathbb{R}^{|V| \times |V|}$: similarity matrix

NON-EUCLIDEAN GEOMETRIC STABILITY

- We replace the Euclidean domain Ω by a general graph $G = (V, E)$.
 $x(u) \in L^2(\Omega) \rightarrow x(u) \in L^2(G)$, $G = (V, E)$.
- In some applications, the input is the graph itself: $x \leftrightarrow G$
- We focus on undirected, possibly weighted graphs:
 $W \in \mathbb{R}^{|V| \times |V|}$: similarity matrix
- Suppose first that G admits a low-dimensional embedding, ie,



$$w_{i,j} = \varphi(x_i, x_j) , \quad x_i \in \Omega \subset \mathbb{R}^d , i, j \leq |V|.$$

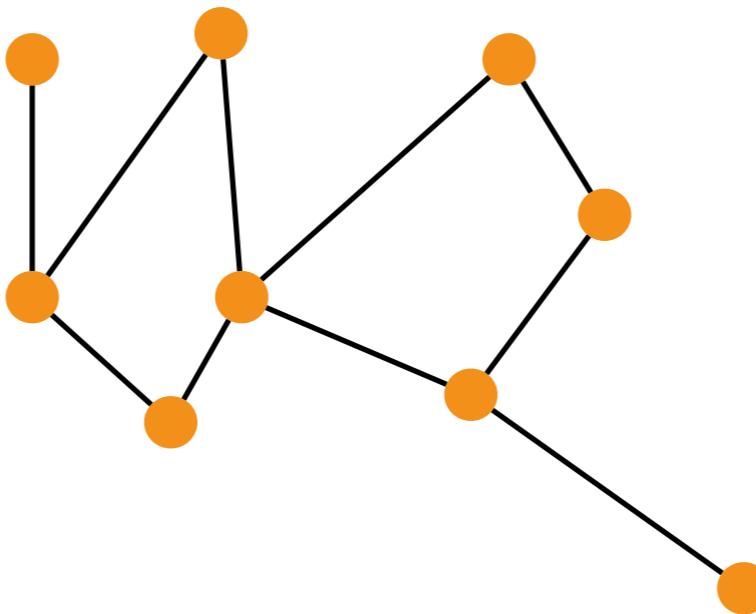
$\varphi(\cdot, \cdot)$: psd kernel (e.g. RBF, dot-product).

NON-EUCLIDEAN EXTRINSIC GEOMETRIC STABILITY

- A deformation field τ in Ω induces a deformation on G :

$$W_\tau = (w_\tau)_{i,j} , \quad (w_\tau)_{i,j} = \varphi(\tau(x_i), \tau(x_j)) .$$

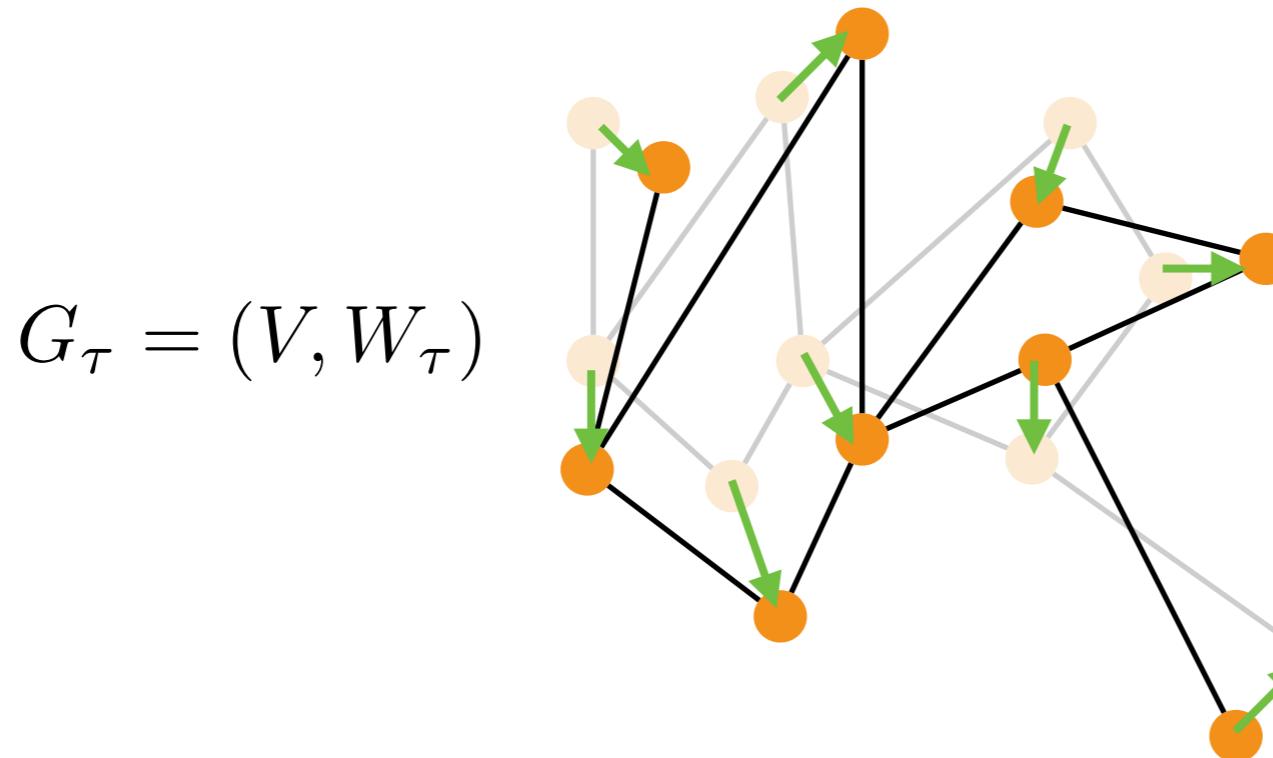
$$G = (V, W)$$



NON-EUCLIDEAN EXTRINSIC GEOMETRIC STABILITY

- A deformation field τ in Ω induces a deformation on G :

$$W_\tau = (w_\tau)_{i,j} , \quad (w_\tau)_{i,j} = \varphi(\tau(x_i), \tau(x_j)) .$$



- Similarly as before, many tasks satisfy geometric stability:
 - particle physics / chemistry.
 - 3D surfaces.
- Can we define geometric deformation/stability intrinsically?

$$f(G) \approx f(G_\tau) \text{ if } \|\nabla \tau\| \text{ small.}$$

DEFORMATIONS AND METRICS

[with F. Gama and A. Ribeiro (U Penn)]

- A deformation in an Euclidean domain Ω induces a change of metric in Ω :

$$\begin{aligned}\langle x_\tau, x'_\tau \rangle_{L^2} &= \int x(u - \tau(u))x'(u - \tau(u))du = \int x(v)x'(v)|\mathbf{1} - \nabla\tau(v)^{-1}|dv \\ &= \int x(v)x'(v)dg(v) = \langle x, x' \rangle_\tau\end{aligned}$$

- A small deformation cost corresponds to a small change of the metric.

$$(1 - o(\|\tau\|))dv \leq dg(v) \leq (1 + o(\|\tau\|))dv$$

DEFORMATIONS AND METRICS

[with F. Gama and A. Ribeiro (U Penn)]

- A deformation in an Euclidean domain Ω induces a change of metric in Ω :

$$\begin{aligned}\langle x_\tau, x'_\tau \rangle_{L^2} &= \int x(u - \tau(u))x'(u - \tau(u))du = \int x(v)x'(v)|\mathbf{1} - \nabla\tau(v)^{-1}|dv \\ &= \int x(v)x'(v)dg(v) = \langle x, x' \rangle_\tau\end{aligned}$$

- A small deformation cost corresponds to a small change of the metric.

$$(1 - o(\|\tau\|))dv \leq dg(v) \leq (1 + o(\|\tau\|))dv$$

- Can we generalize this notion of distance between metric spaces? ie on metrics associated with an arbitrary graph?

GROMOV-HAUSDORFF DISTANCE

[with F. Gama and A. Ribeiro (U Penn)]

- An undirected graph $G = (V, E; W)$ generates a metric given by *shortest-paths*:

$$d_G(i, j) = \text{shortest path between nodes i and j.}$$

GROMOV-HAUSDORFF DISTANCE

[with F. Gama and A. Ribeiro (U Penn)]

- An undirected graph $G = (V, E; W)$ generates a metric given by *shortest-paths*:

$$d_G(i, j) = \text{shortest path between nodes } i \text{ and } j.$$

- One can measure similarity between metric spaces using e.g. Gromov-Hausdorff distance:

$$d_{\text{GH}}(\mathcal{M}, \mathcal{Q}) = \frac{1}{2} \inf_{\begin{array}{l} \varphi : \mathcal{M} \mapsto \mathcal{Q} \\ \psi : \mathcal{Q} \mapsto \mathcal{M} \end{array}} \max\{\|\varphi\|, \|\psi\|, \|(\varphi, \psi)\|\}.$$

$$\|(\varphi, \psi)\| = \sup_{m \in \mathcal{M}, q \in \mathcal{Q}} |d_{\mathcal{M}}(m, \psi(q)) - d_{\mathcal{Q}}(q, \varphi(m))|, \|\varphi\| = \sup_{m, m' \in \mathcal{M}} |d_{\mathcal{M}}(m, m') - d_{\mathcal{Q}}(\varphi(m), \varphi(m'))|.$$

- Introduced on surfaces/point-clouds in [Memoli & Sapiro'05], [Bronstein et al'06].

- Corresponds to a permutation distance when $|V| = |V'| :$

$$d_{\text{P}}(G, G') = \frac{1}{2} \min_{\pi \in \Pi_n} \max_{i, j} |d_G(i, j) - d_{G'}(\pi(i), \pi(j))|.$$

INTRINSIC GEOMETRIC STABILITY PRIORS

[with F. Gama and A. Ribeiro (U Penn)]

- Many inference problems on graphs are stable to intrinsic geometric deformations, in the sense that

$$|f(G) - f(G')| \lesssim d(G, G')$$

- Community Detection.
- Planning, Routing.
- How to leverage geometric stability on graphs?

LINEAR STABLE GENERATORS

[with F. Gama and A. Ribeiro (U Penn)]

- In Euclidean domains Ω , we have seen that *localized*, multiscale filters provide the key to geometric stability.
- These can be expressed as linear operators A of $L^2(\Omega)$ that nearly commute with deformations T_τ :

$$\|AT_\tau - T_\tau A\| \sim \|\nabla \tau\|$$

$$T_\tau x(u) = x(u - \tau(u))$$

$$\begin{array}{|c|c|c|}\hline & -1 & \\ \hline & 1 & \\ \hline & & \\ \hline \end{array}$$

A_1

$$\begin{array}{|c|c|c|}\hline & -1 & \\ \hline & 1 & \\ \hline & & \\ \hline \end{array}$$

A_2

$$\begin{array}{|c|c|c|}\hline & & \\ \hline & -1 & 1 \\ \hline & & \\ \hline \end{array}$$

A_3

$$\begin{array}{|c|c|c|}\hline & & \\ \hline & 1 & \\ \hline & -1 & \\ \hline \end{array}$$

A_4

• • •

LINEAR STABLE GENERATORS

[with F. Gama and A. Ribeiro (U Penn)]

- In Euclidean domains Ω , we have seen that *localized*, multiscale filters provide the key to geometric stability.
- These can be expressed as linear operators A of $L^2(\Omega)$ that nearly commute with deformations T_τ :

$$\|AT_\tau - T_\tau A\| \sim \|\nabla \tau\|$$

$$T_\tau x(u) = x(u - \tau(u))$$

$$\begin{array}{|c|c|c|}\hline & -1 & \\ \hline & 1 & \\ \hline & & \\ \hline \end{array}$$

A_1

$$\begin{array}{|c|c|c|}\hline & -1 & \\ \hline & & 1 \\ \hline & & \\ \hline \end{array}$$

A_2

$$\begin{array}{|c|c|c|}\hline & & \\ \hline & -1 & 1 \\ \hline & & \\ \hline \end{array}$$

A_3

$$\begin{array}{|c|c|c|}\hline & & \\ \hline & & 1 \\ \hline & -1 & \\ \hline \end{array}$$

A_4

• • •

- We can write a CNN layer as a linear combination of such operators:

$$\tilde{x} = \rho \left(\sum_k (A_k x) \theta_k \right) . \quad \theta_1, \dots, \theta_k, \in \mathbb{R}^{p \times \tilde{p}} .$$

- What about general graphs?

LINEAR STABLE GENERATORS

[with F. Gama and A. Ribeiro (U Penn)]

- Linear diffusion on graphs is given by its adjacency matrix $A(G)$

$A(G)_{i,j} = 1$ iff $(i, j) \in E$. $W_{i,j}$ in weighted graphs.

- By definition, this is a localized operator. Local smoothing.

LINEAR STABLE GENERATORS

[with F. Gama and A. Ribeiro (U Penn)]

- Linear diffusion on graphs is given by its adjacency matrix $A(G)$

$A(G)_{i,j} = 1$ iff $(i, j) \in E$. $W_{i,j}$ in weighted graphs.

- By definition, this is a localized operator. Local smoothing.
- Q: Stable to deformations? By definition,

$$\inf_{P \in \Pi_n} \|W - PW'P^\top\| = d_{\text{P}}(G, G') \lesssim d_{\text{GH}}(G, G')$$

- Up to rigid (isometric) transformation, local diffusion is continuous with respect to metric deformations.

LINEAR STABLE GENERATORS

[with F. Gama and A. Ribeiro (U Penn)]

- Linear diffusion on graphs is given by its adjacency matrix $A(G)$

$A(G)_{i,j} = 1$ iff $(i, j) \in E$. $W_{i,j}$ in weighted graphs.

- By definition, this is a localized operator. Local smoothing.
- Q: Stable to deformations? By definition,
$$\inf_{P \in \Pi_n} \|W - PW'P^\top\| = d_P(G, G') \lesssim d_{\text{GH}}(G, G')$$
- Up to rigid (isometric) transformation, local diffusion is continuous with respect to metric deformations.
- Together with the degree matrix $D = \text{diag}(W\mathbf{1})$, it defines a high-pass filter, the *Graph Laplacian*: $\Delta = D - W$.
- It is also localized and stable to deformations in the sense of GH.

GRAPH NEURAL NETWORKS

[Scarselli et al.,'09], [Gori et al. '05]

- Given a signal $x \in \mathbb{R}^{V \times p}$, a Graph Neural Network (GNN) layer considers generators $[D, W]$ and trainable coefficients $\Theta = (\theta_1, \theta_2)$:

$$\tilde{x} = \rho(Dx\theta_1 + Wx\theta_2) . \quad \theta_1, \theta_2 \in \mathbb{R}^{p \times \tilde{p}} .$$

- Flexible model: does not require fixed input graphs.
- Initial version was inspired from the Message-Passing algorithm.
Fixed point of a trainable, non-linear diffusion.
- Modernized in [Li et al.'15], [Duvenaud et al.'15], [Subkhaatar et al.'16],
- Authors also explored other forms of nonlinearity, e.g. *gating*.
- Similarly as in CNNs, we can also consider pooling layers, (*provided we have a graph coarsening scheme*).

LAPLACIAN INTERPRETATION

- Since we are learning a linear combination $A(G)$ and $D(G)$, we can reparametrize the generator in terms of the *Graph Laplacian*:

$$\Delta(G) = D(G) - A(G)$$

- If we consider generators of the form $[1, \Delta, \Delta^2, \dots]$, the resulting GNN layer is expressed as a polynomial $\hat{\Delta}$:

$$\tilde{x} = \rho(\theta(\Delta)x), \quad \theta(\Delta) = \sum_{s=0}^S \theta_s \Delta^s.$$

- In Spectral Networks [B. et al'14], we train directly on the spectrum of the Laplacian:

$$\tilde{x} = \rho(V^T \text{diag}(\alpha)Vx), \quad \alpha = \mathcal{K}(\theta), \quad \mathcal{K} : \text{spline kernel}$$

- Computationally expensive and unstable to deformations for varying graph.

EFFICIENT SPECTRAL NETWORKS

- These issues were addressed in [Defferrard et al.'16] by sticking with the polynomial representation $\theta(\Delta)$ instead.
 - Since the spectrum of Δ in finite dimensions is bounded, the optimal polynomial basis in the interval is given by Chebyshev polynomials:

GNN with generators $[F_j(\tilde{\Delta})]_{j \leq J}$, $\tilde{\Delta} = \frac{2}{\|\Delta\|}\Delta - 1$.

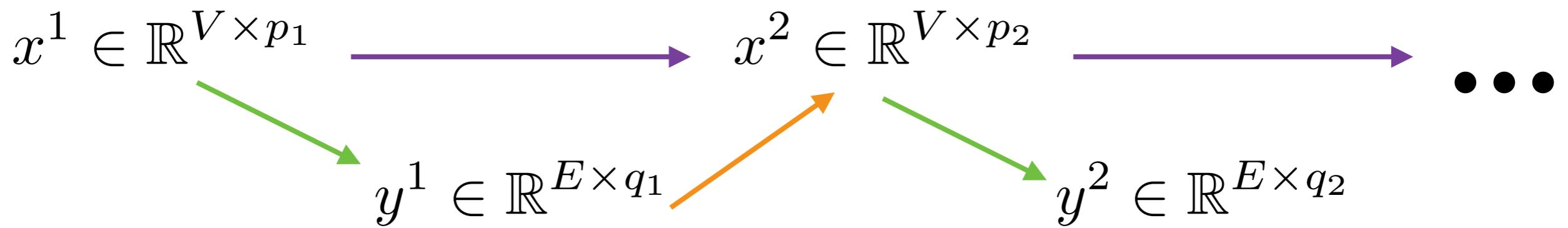
- [Kipf & Welling, '17] Further simplified the construction using only $J = 1$ and the normalized Laplacian
$$\bar{\Delta} = \mathbf{1} - D^{-1/2}WD^{-1/2}.$$
- Better frequency localization can be obtained by replacing (real) Chebyshev polynomial with (complex) Cayley filters [Monti et al.'17].

EXTENSIONS/LIMITATIONS

- As opposed to Euclidean domains, in general graphs we only have an isotropic high-pass filter (Δ), but no oriented filters.

EXTENSIONS/LIMITATIONS

- As opposed to Euclidean domains, in general graphs we only have an isotropic high-pass filter (Δ), but no oriented filters.
- Inspired by Message-Passing algorithms, we can generalize GNNs to alternate between vertex and edge representations:



$$y^1(e) = \psi_\theta(x^1(i), x^1(j)) , \quad e = (i, j) \in E .$$

$$x^2(i) = \phi_\theta(\{y^1(e)\}_{e \in i}) , \quad i \in V .$$

- Used for example in [Battaglia et al.'16] for N-body prediction dynamics and [Gilmer et al.'17] for quantum chemistry.

SURFACE REPRESENTATIONS

[joint work with I. Kostrikov, D.Panozzo, D.Zorin (NYU)]

- In the particular case where G represents a 3D surface, we have a *mesh* representation:

$$M = (V, E, F) , \quad F = \{(i, j, k)\} \text{ triangulation}$$



credit: jonathanpuckey

SURFACE REPRESENTATIONS

[joint work with I. Kostrikov, D.Panozzo, D.Zorin (NYU)]

- In the particular case where \mathcal{G} represents a 3D surface, we have a *mesh* representation:
 $M = (V, E, F)$, $F = \{(i, j, k)\}$ triangulation

- In that case, we can compute a “proper” square root of the Laplacian, the *Dirac* operator:

$$\Delta = D^* D , D \in \mathbb{H}^{V \times F}$$

credit: jonathanpuckey

- Defined over quaternion space.
- Captures principal curvature directions (ie orientation).



SURFACE REPRESENTATIONS

- In the particular case where G represents a 3D surface, we have a *mesh* representation:

$$M = (V, E, F) , \quad F = \{(i, j, k)\} \text{ triangulation}$$

- In that case, we can compute a “proper” square root of the Laplacian, the *Dirac* operator:

$$\Delta = D^* D , \quad D \in \mathbb{H}^{V \times F}$$



- Defined over quaternion space.

credit: jonathanpuckey

GT



MLP



AvgPool



Laplace



Dirac



LAPLACE NETWORK STABILITY

- Stable graph generators result in stable GNN representations:

Theorem: [B, K, P, Z'17] Let $G = (V, E)$ and suppose $V \subset \Omega \in \mathbb{R}^d$. Let $\Phi(x; \Delta)$ be a R -layer Laplace GNN with generators $\{I, \Delta\}$, and τ a deformation field on Ω . Then

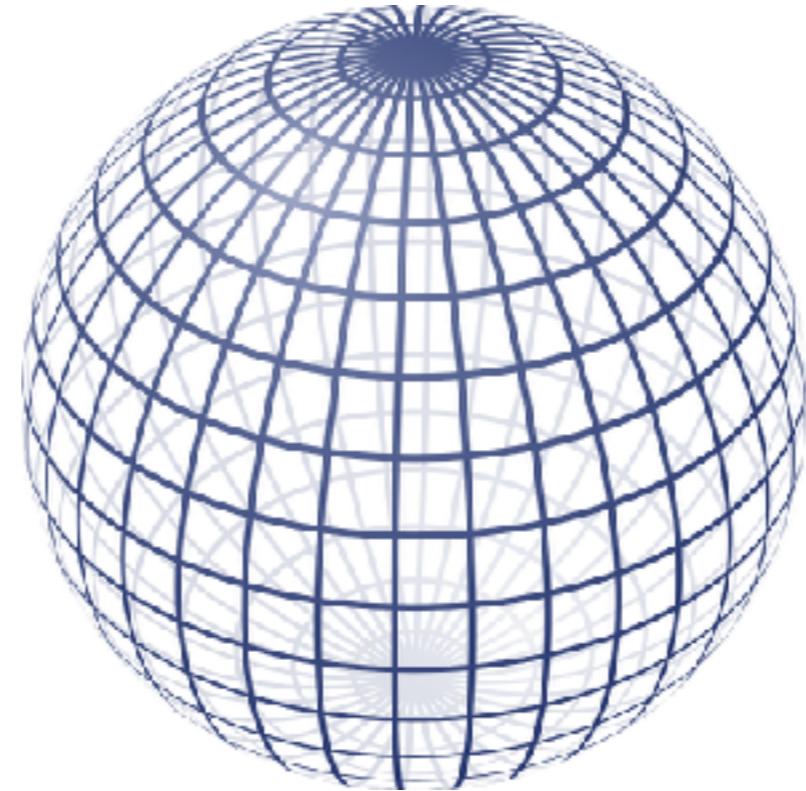
1. $\|\Phi(x; \Delta) - \Phi(x'; \Delta)\| \leq C(\Theta) \|x - x'\|^{h(\beta)} ,$
2. $\|\Phi(x; \Delta) - \Phi(x; \tau(\Delta))\| \leq C'(\Theta) \|\nabla \tau\|^{h(\beta)} ,$

where $h(\beta) = \prod_{r \leq R} \frac{\beta_r - 1}{\beta_r - 1/2}$ measures smoothness (Sobolev) of feature maps.

- In Euclidean graphs, the Laplacian is geometrically stable.
- *Caveat:* We currently require explicit smoothness decay of feature maps.
- *Future work:* Extension to intrinsic deformations.

NON-EUCLIDEAN BUT REGULAR DOMAINS

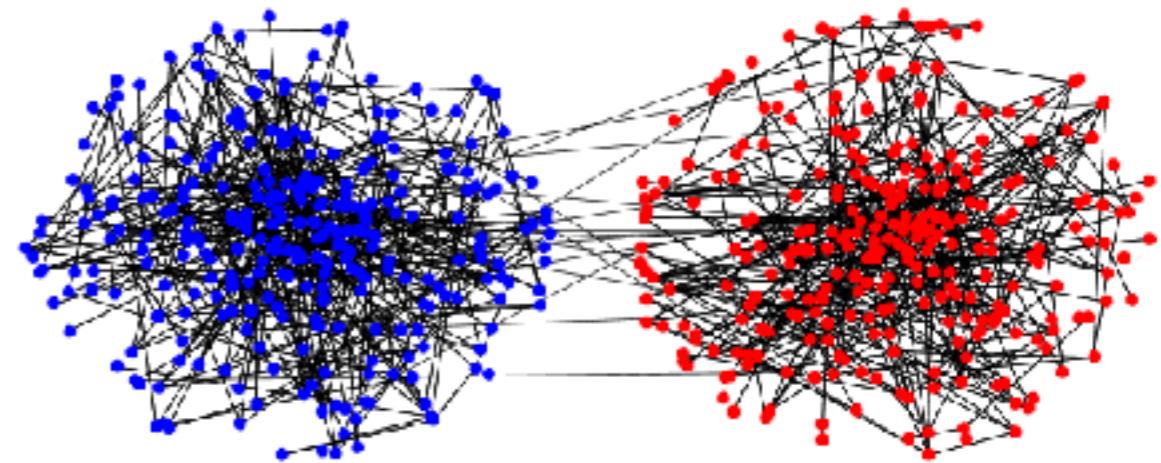
- One important particular case is when the domain Ω is regular, e.g the sphere:



- The set of transformations that leave the sphere S^2 invariant forms a Lie Group, $\text{SO}(3)$.
- We can extend CNNs to such settings by replacing planar convolutions with group convolutions [Cohen & Welling]

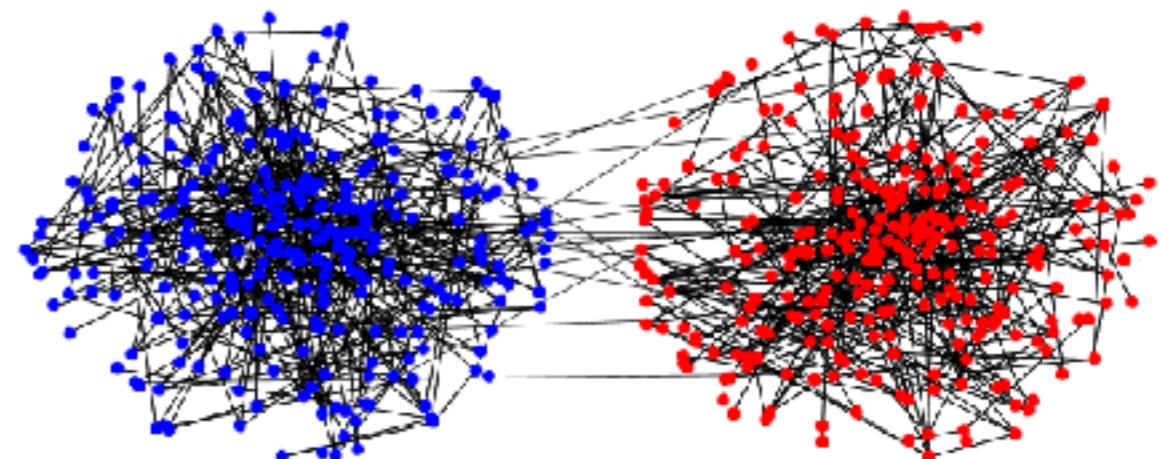
INVERSE PROBLEMS ON GRAPHS

- Consider the problem of inferring communities within a network:

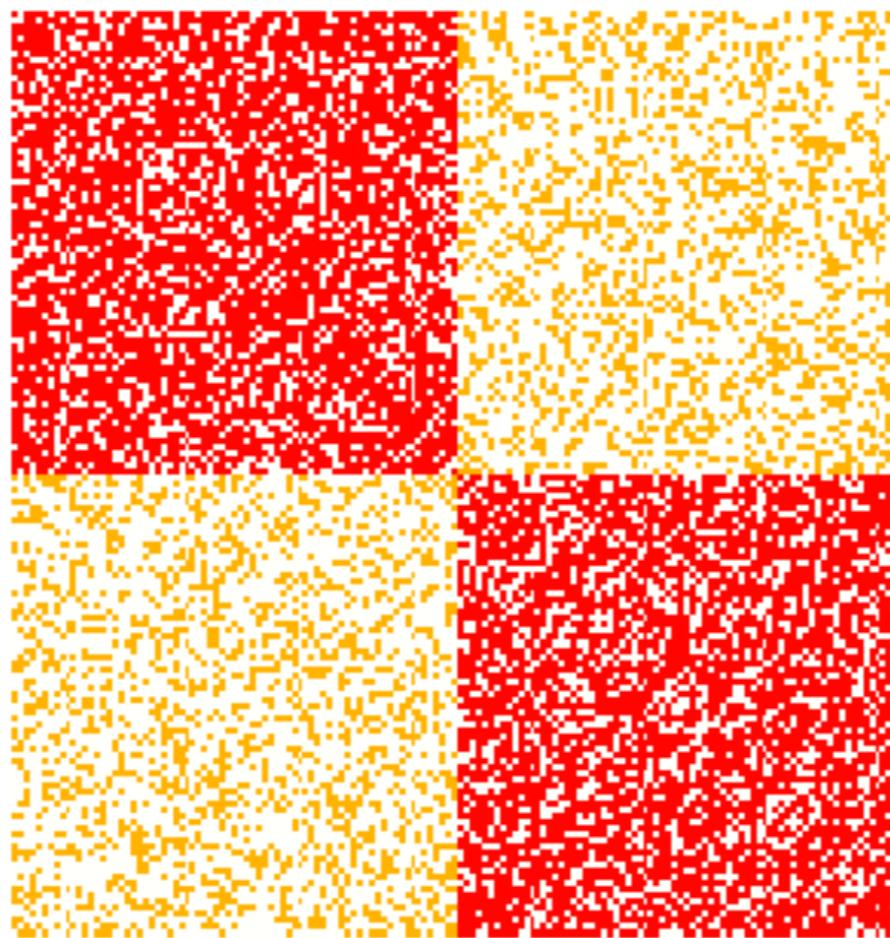


INVERSE PROBLEMS ON GRAPHS

- Consider the problem of inferring communities within a network:

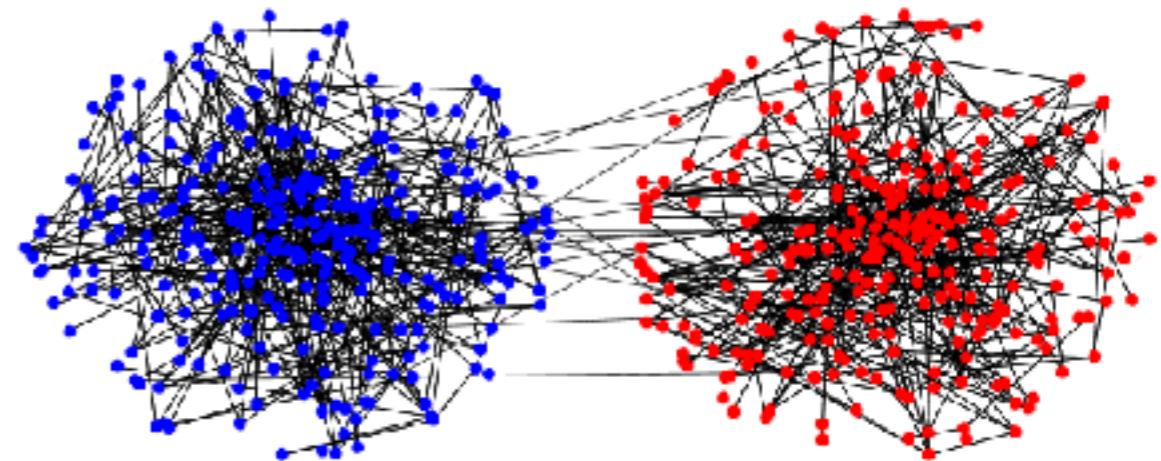


Adjacency matrix associative 2 communities

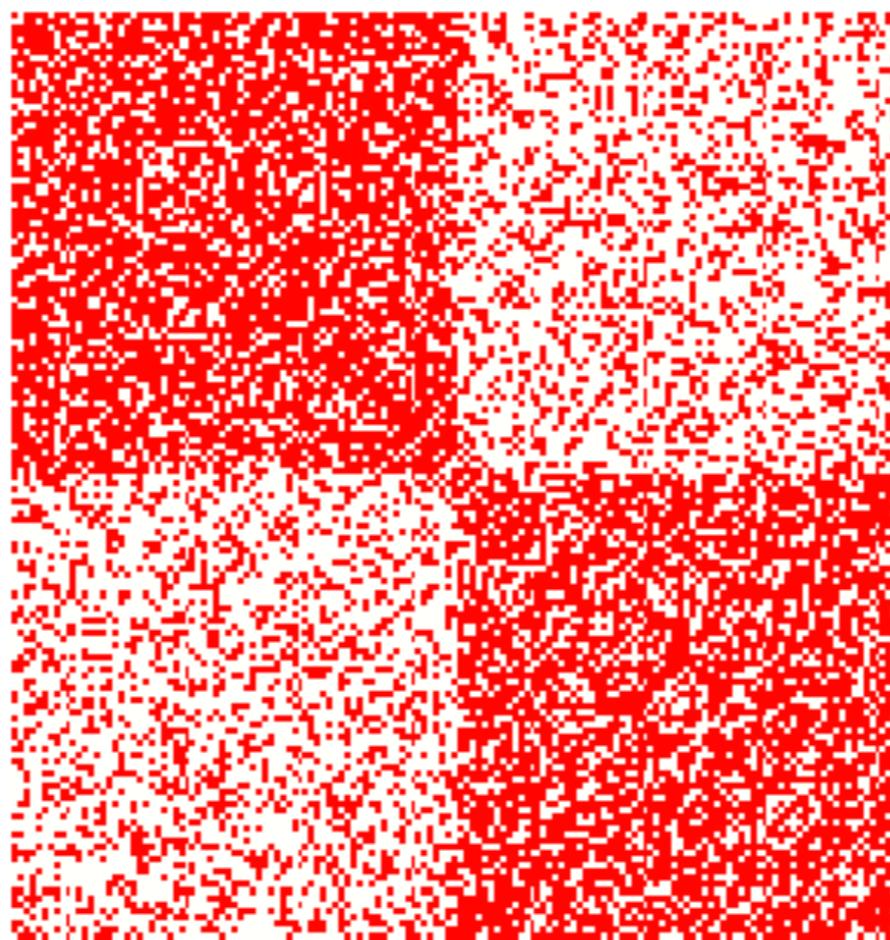


INVERSE PROBLEMS ON GRAPHS

- Consider the problem of inferring communities within a network:

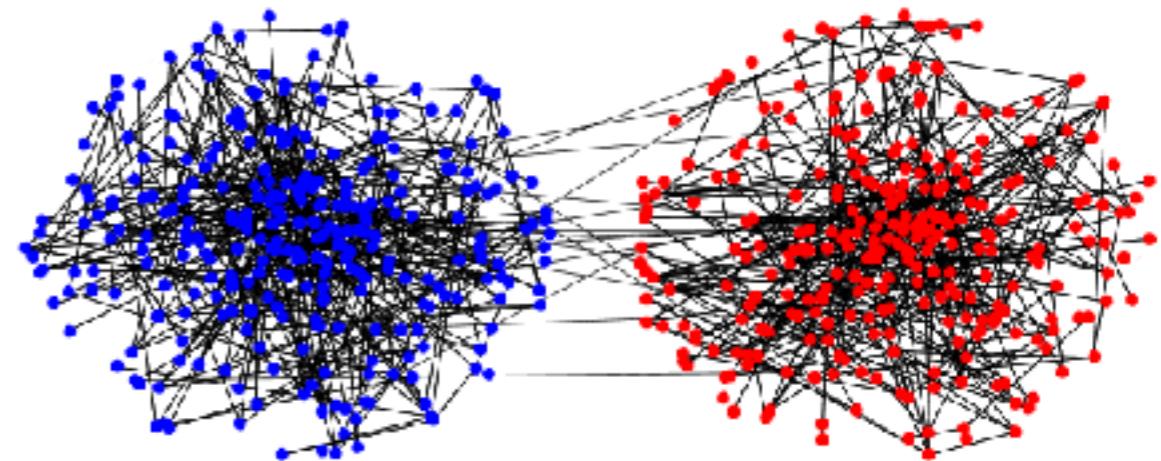


Adjacency matrix associative 2 communities

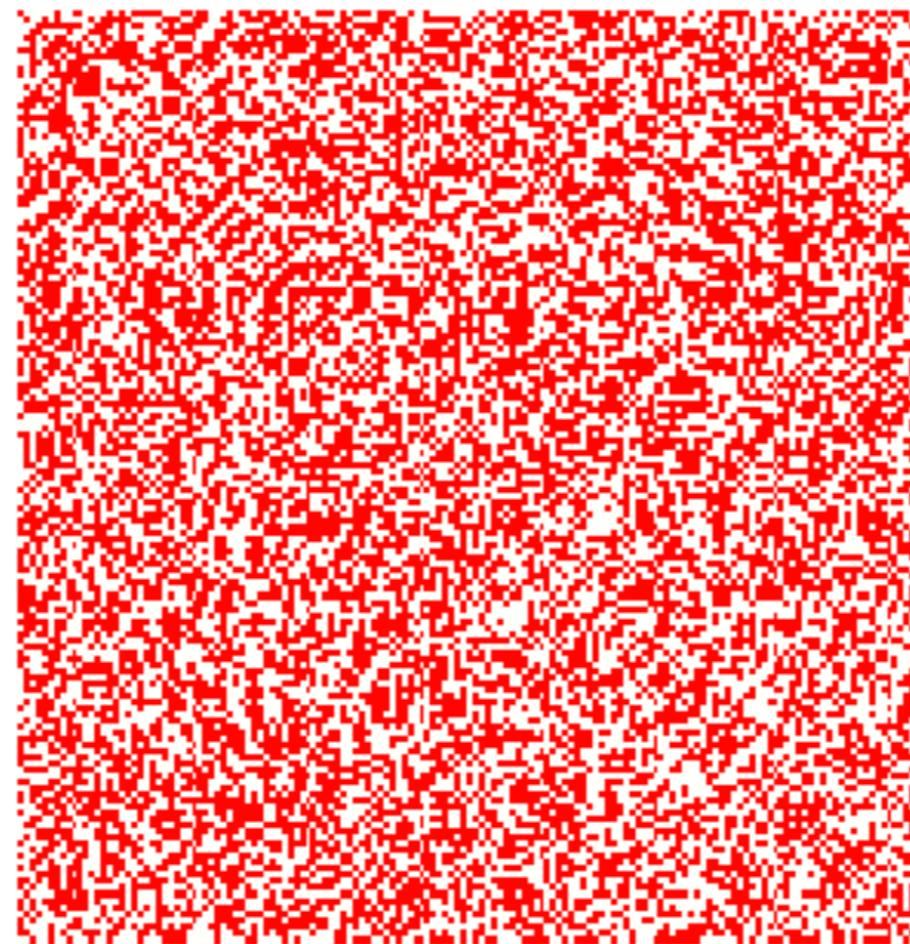


INVERSE PROBLEMS ON GRAPHS

- Consider the problem of inferring communities within a network:



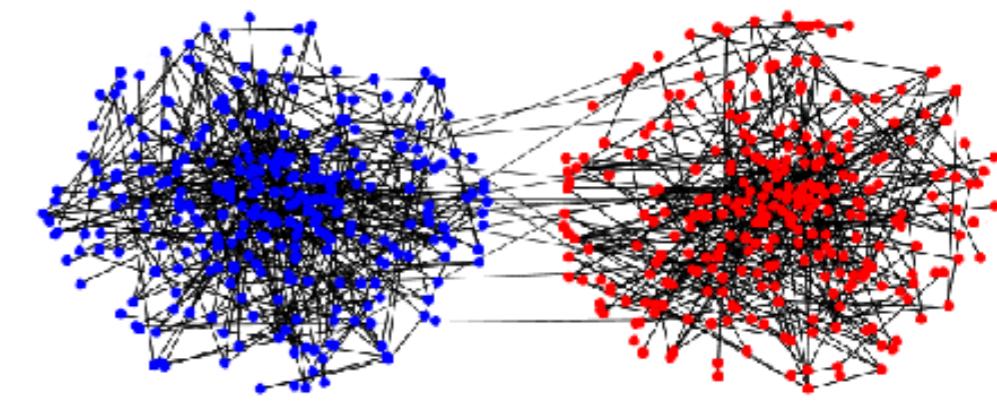
Adjacency matrix associative 2 communities



INVERSE PROBLEMS ON GRAPHS

- Community Detection in graphs.
 - Studied in the Stochastic Block Model.
 - Hardness of estimation is controlled by a Signal-to-Noise Ratio:

$$\text{SNR} = \frac{(a - b)^2}{k(a + (k - 1)b)}$$

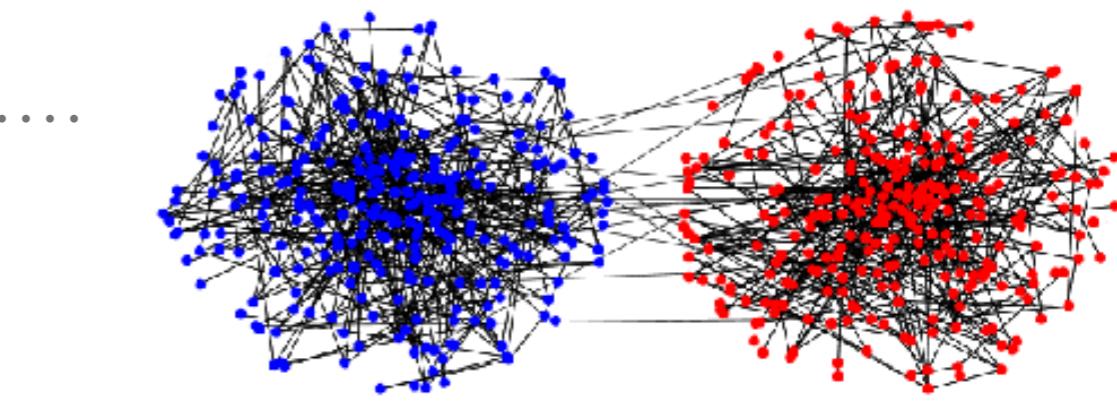


a: inner connection probability.
b: outer connection probability.

INVERSE PROBLEMS ON GRAPHS

- Community Detection in graphs.
 - Studied in the Stochastic Block Model.
 - Hardness of estimation is controlled by a Signal-to-Noise Ratio:

$$\text{SNR} = \frac{(a - b)^2}{k(a + (k - 1)b)}$$



a : inner connection probability.
 b : outer connection probability.

- Two major algorithmic frameworks:
 - Graph conductance/min-cut approach, leading to spectral clustering algorithms.

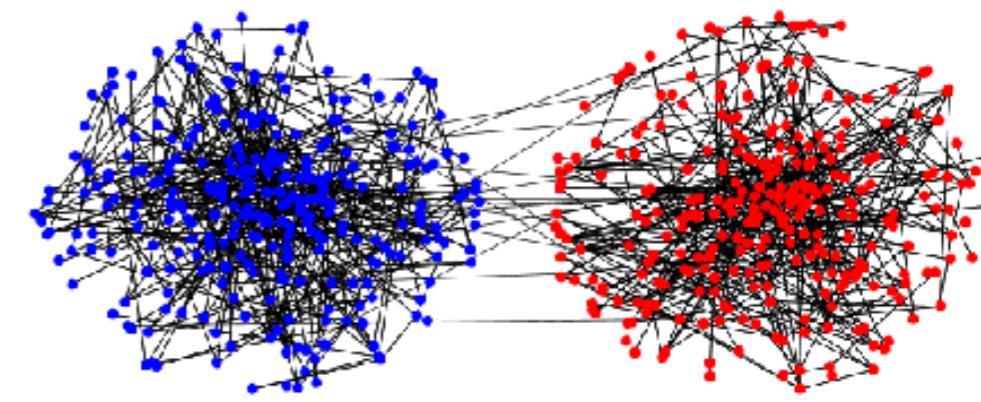
$$\min_{y_i = \pm 1; \bar{y} = 0} y^T \mathcal{A}(G) y .$$

- Probabilistic Graphical Models, leading to Belief Propagation.

$$p(y|G) \propto \prod_{(i,j) \in E} \varphi_{(i,j)}(y_i, y_j) \prod_{v \in V} \psi_i(y_i)$$

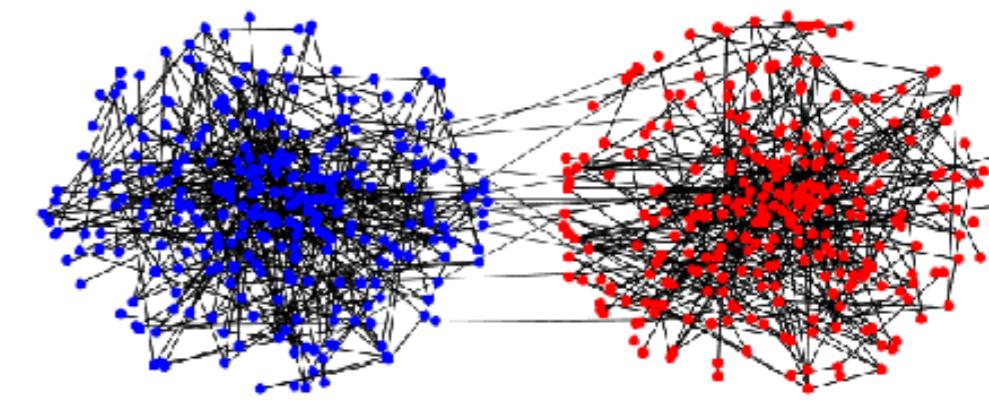
INVERSE PROBLEMS ON GRAPHS

- Community Detection in graphs.
 - Studied in the Stochastic Block Model.
 - Hardness of estimation is controlled by a Signal-to-Noise Ratio.
- Recent research program has unified both approaches using tools from statistical physics, and identified computational and information theoretic thresholds:
 - When is the detection statistically possible?
 - When is the detection feasible with polynomial-time algorithms?



INVERSE PROBLEMS ON GRAPHS

- Community Detection in graphs.
 - Studied in the Stochastic Block Model.
 - Hardness of estimation is controlled by a Signal-to-Noise Ratio.
- Recent research program has unified both approaches using tools from statistical physics, and identified computational and information theoretic thresholds:
 - When is the detection statistically possible?
 - When is the detection feasible with polynomial-time algorithms?
- Q: Can we *learn* those algorithms from the data using graph neural networks? reaching detection thresholds?



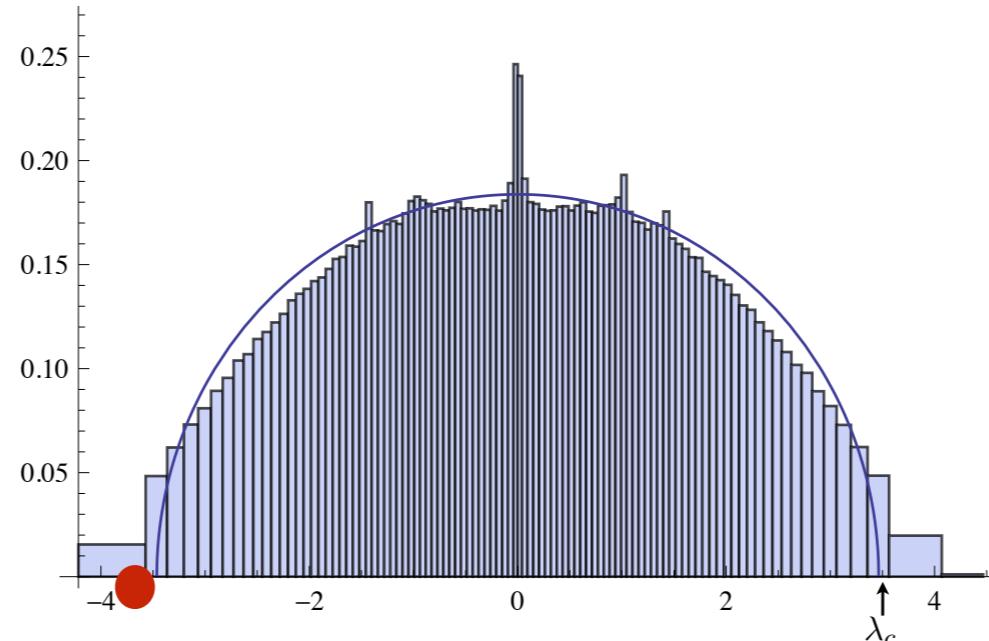
DATA-DRIVEN COMMUNITY DETECTION

[joint work with Li Shuai Li (UC Berkeley)]

- $\mathcal{A}(G)$: linear operator defined on G , eg Laplacian $\Delta = D - A$.
- Spectral Clustering estimators (2-community case):

$$\hat{y} = \text{sign}(\text{Fiedler}(\mathcal{A}(G))) ,$$

Fiedler(M): eigenvector corresponding to 2nd smallest eigenvalue



- Iterative algorithm: projected power iterations on shifted $\mathcal{A}(G)$:

$$M = \|\mathcal{A}(G)\| \mathbf{1} - \mathcal{A}(G)$$

DATA-DRIVEN COMMUNITY DETECTION

[joint work with Lisha Li (UC Berkeley)]

- We consider a GNN generated by operators $\{\mathbf{1}, A, D\}$:

$$\tilde{x} = \rho (\theta_1 x + \theta_2 Dx + \theta_3 Ax) .$$

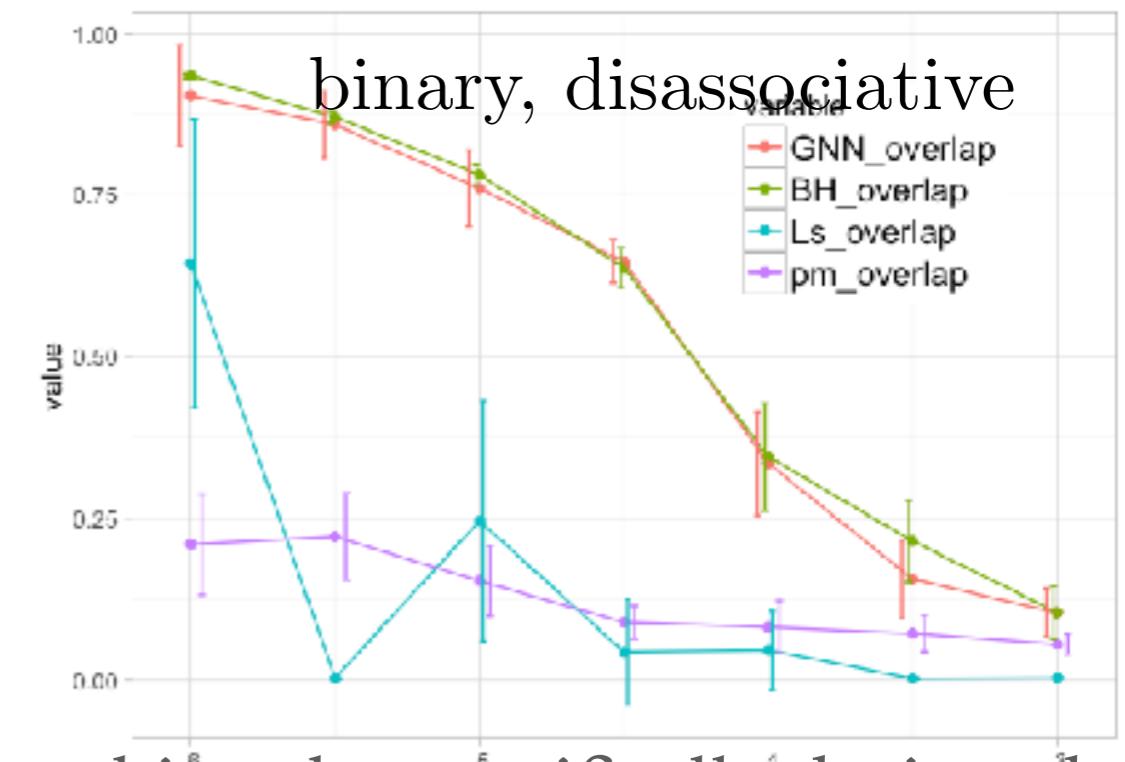
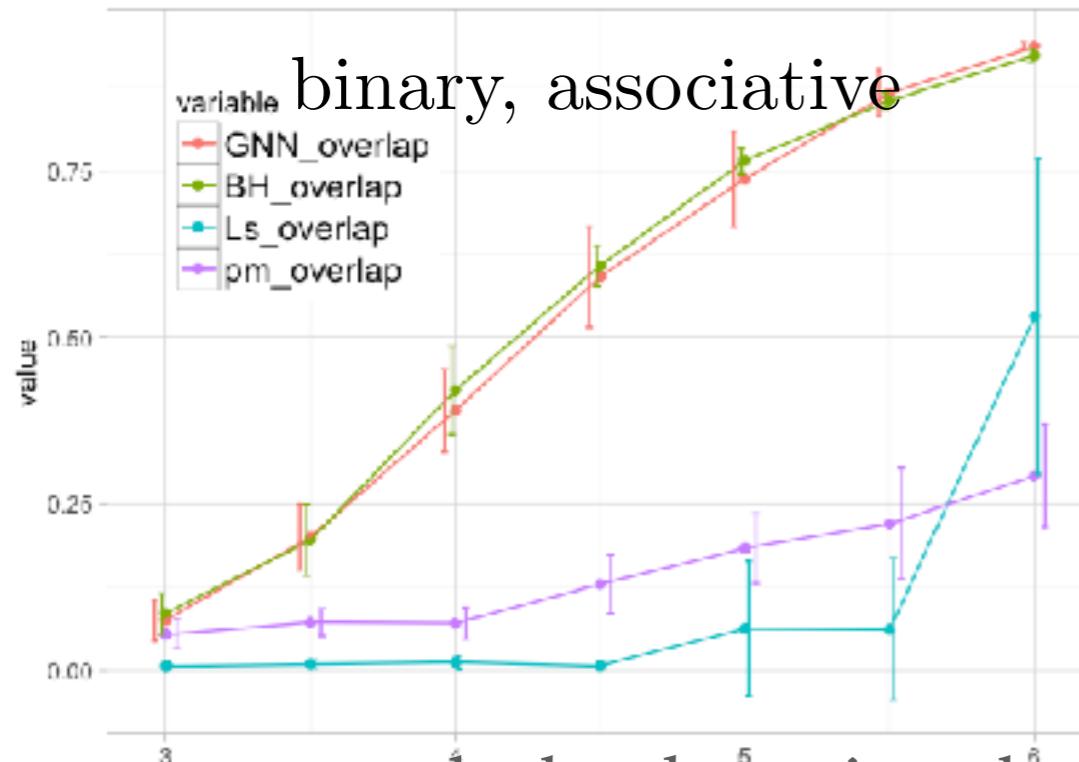
- They generate the so-called *Bethe Hessian*:

$$BH(r) = (r^2 - 1)\mathbf{1} - rA + D$$

- ❖ Second-order approximation of Bethe Free energy at critical points of BP.
- ❖ In that case, Laplacian generator does not work: its spectrum is dominated by few nodes with dominant degree.
- We train it by back propagation using a loss that is globally invariant to label permutations.

REACHING DETECTION THRESHOLD ON SBM

- Stochastic Block Model Results: [joint work with Lisha Li (UC Berkeley)]



- we reach the detection threshold, matching the specifically designed spectral method.
- Real-world community detection results on SNAP data
[Leskovec et al]

Table 1: Snap Dataset Performance Comparison between GNN and AGM

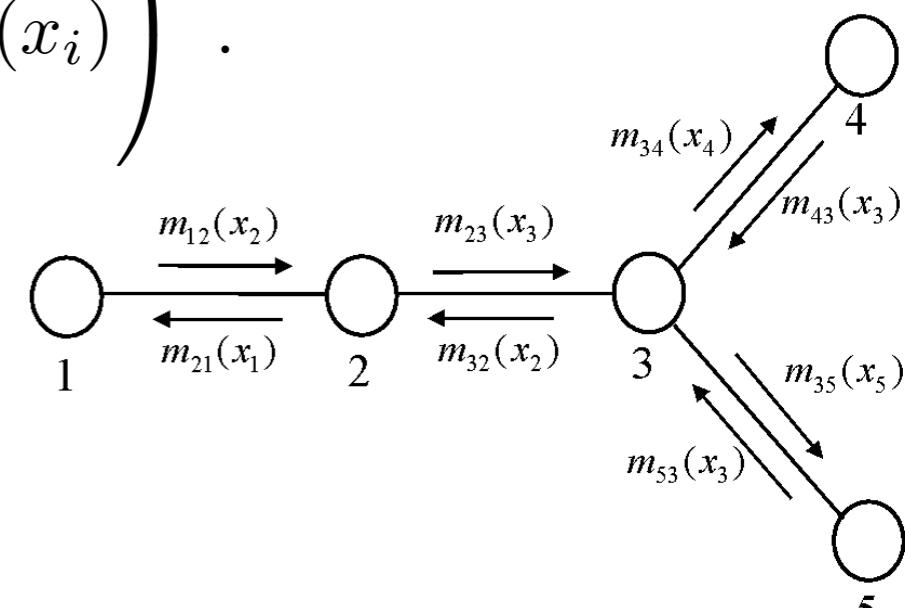
Dataset	(train/test)	Subgraph Instances		Overlap Comparison		
		Avg Vertices	Avg Edges	GNN	AGMFit	
Amazon	315 / 35	60	346	0.74 ± 0.13	0.76 ± 0.08	
DBLP	2831 / 510	26	164	0.78 ± 0.03	0.64 ± 0.01	
Youtube	48402 / 7794	61	274	0.9 ± 0.02	0.57 ± 0.01	

BELIEF PROPAGATION

- For small number of communities, the IT detection threshold is provably matched by (loopy) BP.
- BP performs message-passing updates of the form

$$m_{ij}(x_j) \leftarrow \sum_{x_i} \left(\tilde{\phi}_i(x_i; y) \psi_{ij}(x_i, x_j) \prod_{k \in N(i) \setminus j} m_{ki}(x_i) \right).$$

$$b_j(x_j) = \frac{1}{Z_j} \tilde{\phi}_j(x_j; y) \prod_{i \in N(j)} m_{ij}(x_j).$$



- exact inference on simply connected graphs.
- on general graphs, fixed points of BP correspond to critical points of the Bethe Free Energy.
- Messages are defined and propagated over edges of G , and account for non-backtracking paths.

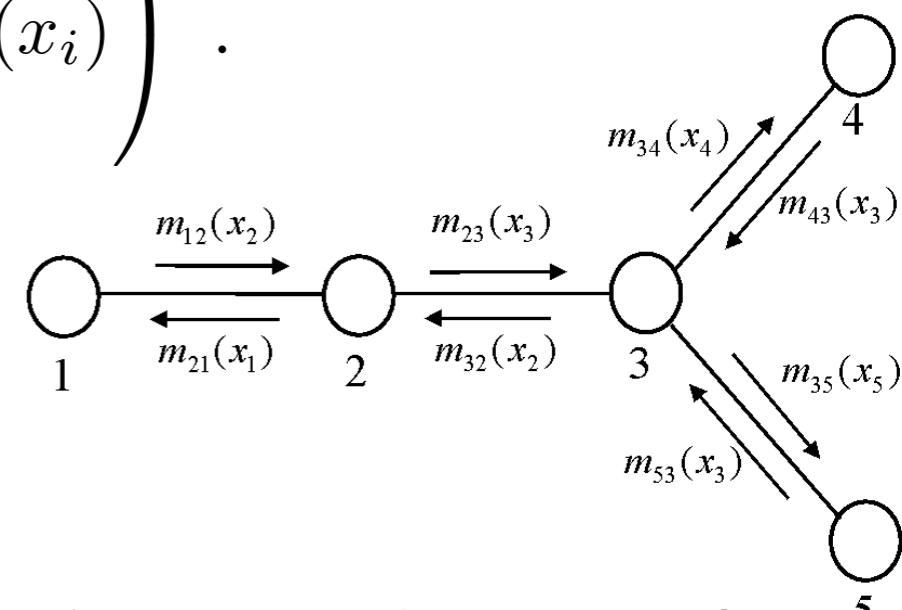
BELIEF PROPAGATION

- For small number of communities, the IT detection threshold is provably matched by (loopy) BP.
- BP performs message-passing updates of the form

$$m_{ij}(x_j) \leftarrow \sum_{x_i} \left(\tilde{\phi}_i(x_i; y) \psi_{ij}(x_i, x_j) \prod_{k \in N(i) \setminus j} m_{ki}(x_i) \right).$$

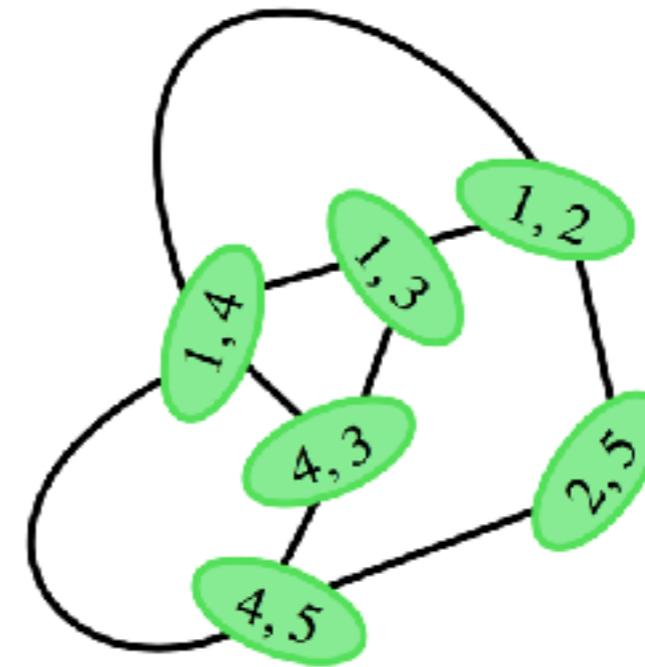
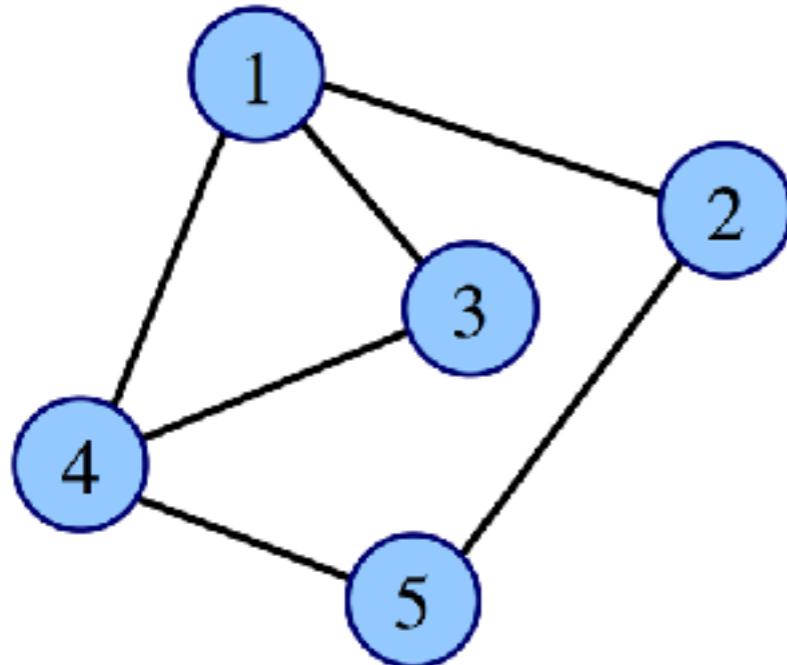
$$b_j(x_j) = \frac{1}{Z_j} \tilde{\phi}_j(x_j; y) \prod_{i \in N(j)} m_{ij}(x_j).$$

- exact inference on simply connected graphs.
- on general graphs, fixed points of BP correspond to critical points of the Bethe Free Energy.
- Messages are defined and propagated over edges of G , and account for non-backtracking paths. GNN version?



GRAPH NEURAL NETWORKS ON GRAPH HIERARCHIES

- The *line graph* of $G = (V, E)$ is a new graph $L(G) = (V', E')$ that models the adjacency of the edges: $V' \cong E$



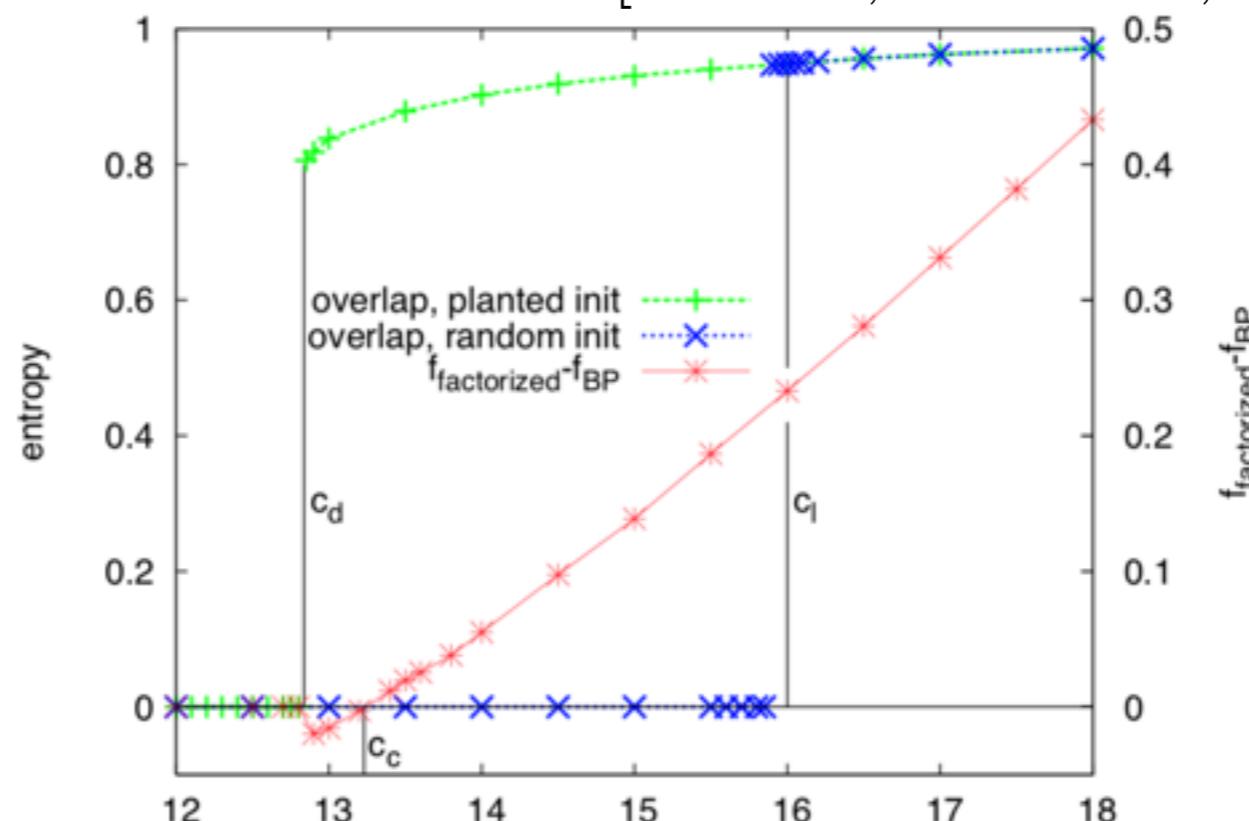
(source:wikipedia)

- We *augment* the GNN with analogous operations on $L(G)$.
- Related to Covariant Compositional Networks [Kondor et al'18], and neural message-passing [Gilmer et al.'17].

COMPUTATIONAL-TO-STATISTICAL GAPS

- When # of communities > 4 , there is a gap between the information theoretical threshold and current known polynomial-time algorithms:

[Decelle, Krzakala, Moore, Zdeborova, '13]

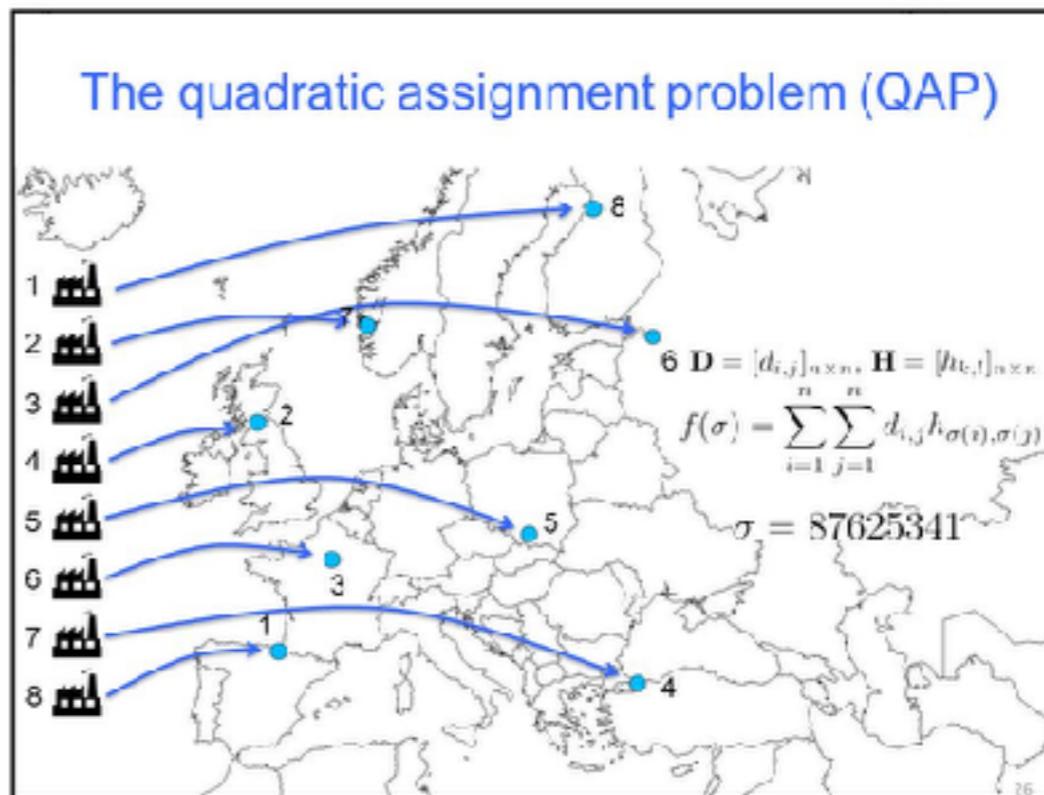


- Preliminary results for 5 communities, $\overline{\deg}^c = 14.5$:

$N = 10^3$	G	$\{G, L(G)\}$	BP
Overlap	29.5 ± 0.5	30.1 ± 0.5	30.4 ± 3

QUADRATIC ASSIGNMENT PROBLEM

- Find an assignment that optimizes the transportation cost between two graphs:



$$\min_{X \in \Pi_n} \text{Tr}(A_1 X A_2 X^T).$$

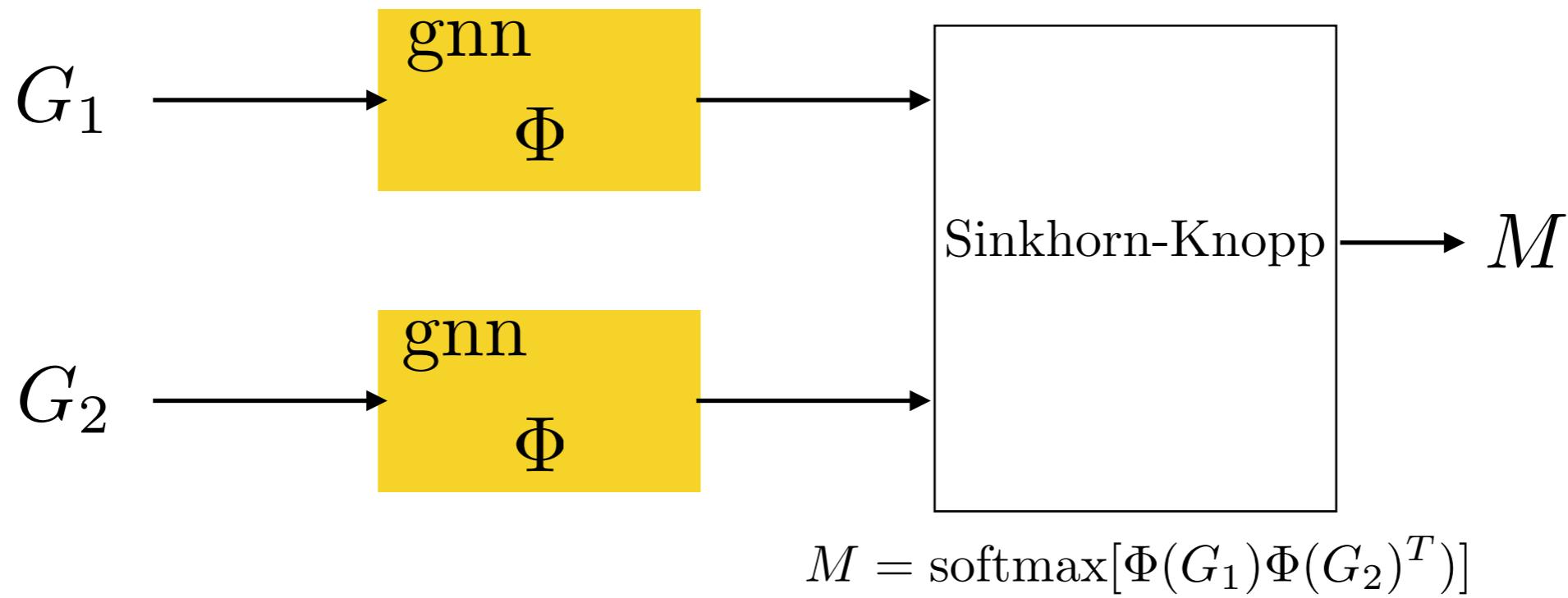
Π_n : space of $n \times n$ permutation matrices.

- NP-hard
- Contains the TSP as a particular instance.
- Relaxations using SDP and Spectral Approaches.

[with S. Villar (NYU), A. Nowak (NYU) and A. Bandeira (NYU)]

QUADRATIC ASSIGNMENT PROBLEM

- We learn approximate solutions using siamese graph neural networks:



- We train the model to predict the correct permutation matrix on a dataset of *planted* solutions:

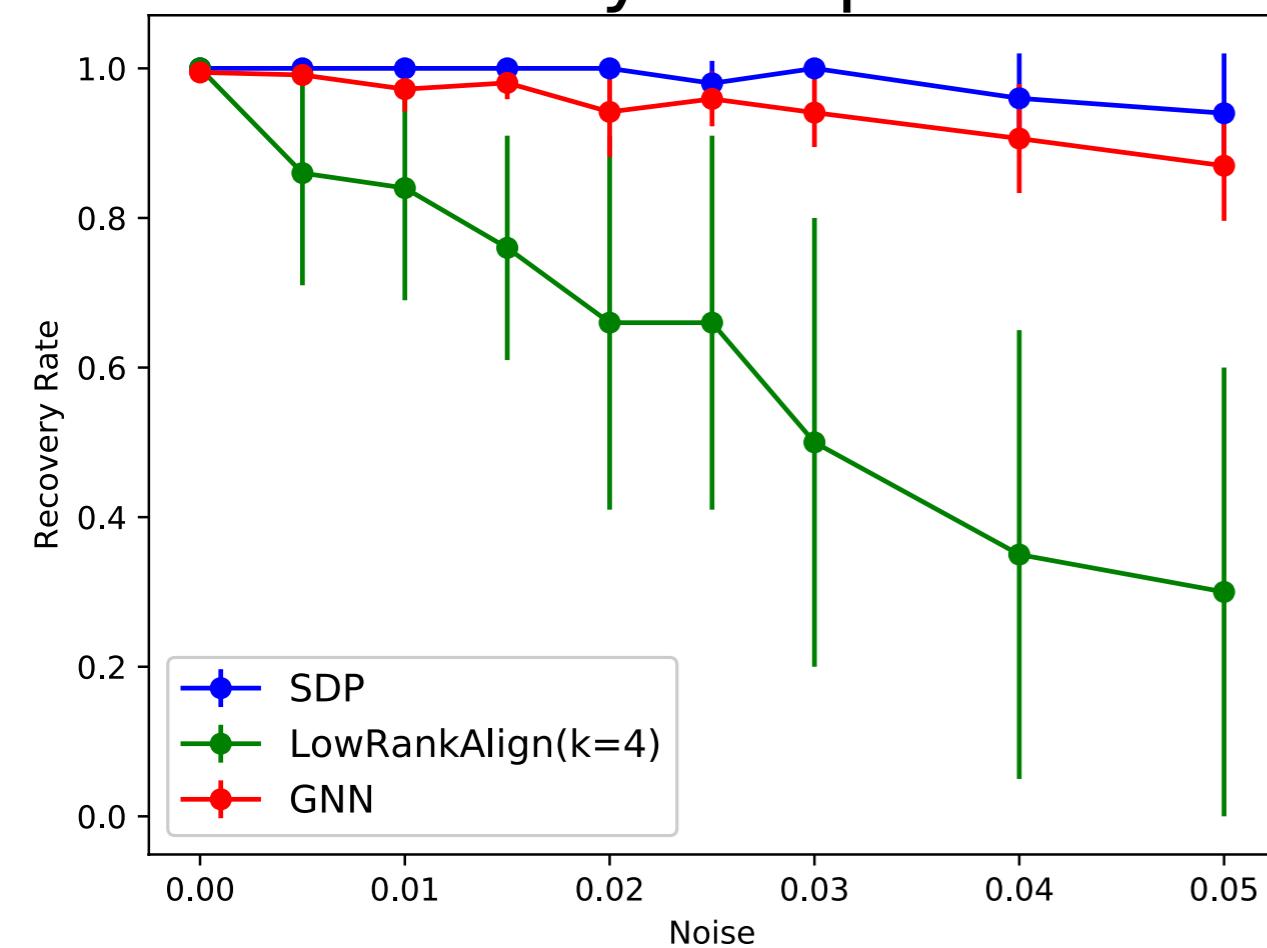
$$G_1 = PG_2 + N \quad N \sim \text{Erdos-Renyi}$$

$G_2 \sim \text{Erdos-Renyi}$

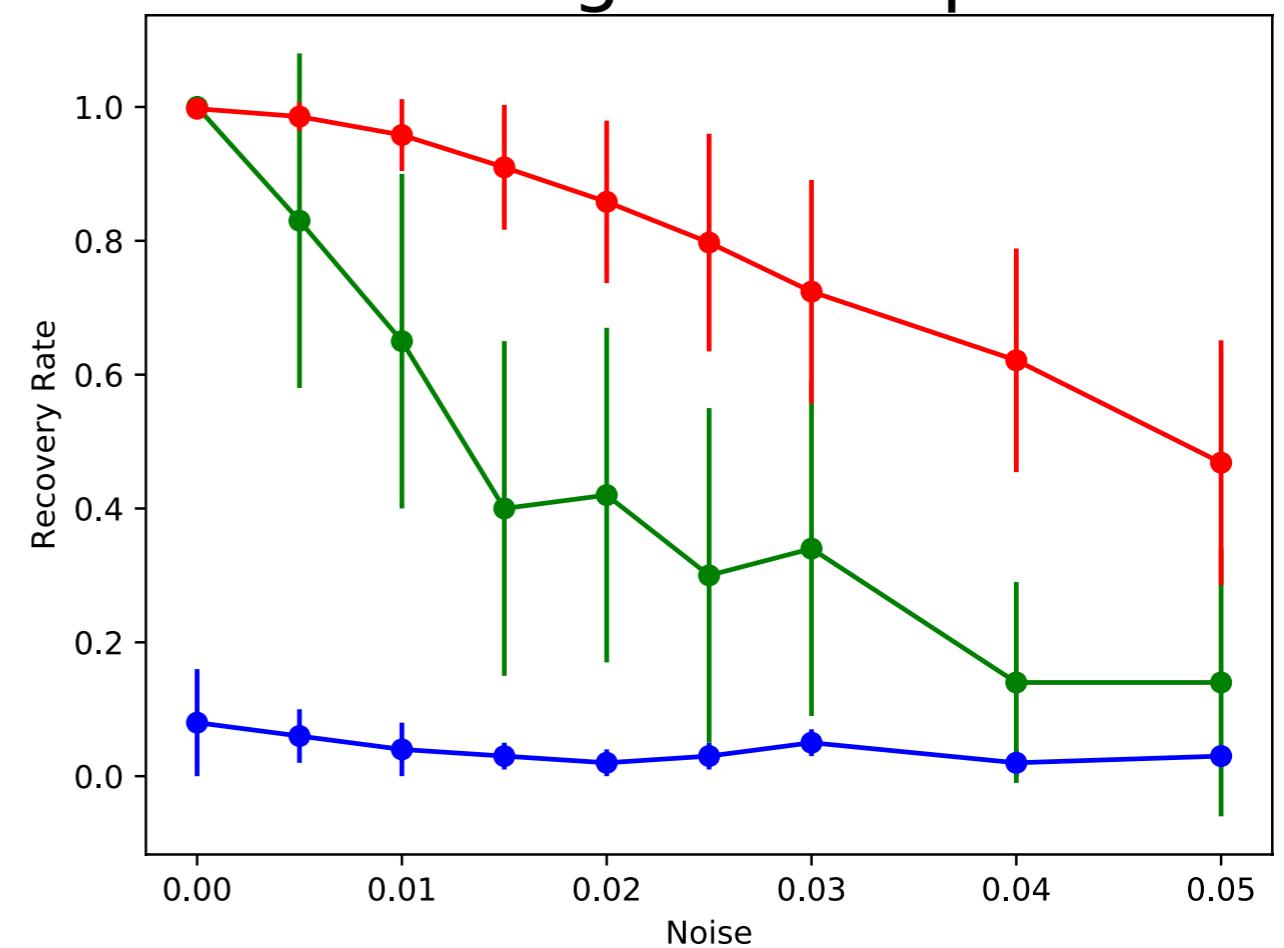
$G_2 \sim \text{Random Regular}$

QUADRATIC ASSIGNMENT PROBLEM

ErdosRenyi Graph Model



Random Regular Graph Model



- Our model runs in $o(n^2)$, LowRankAlign is $o(n^3)$ SDP in $o(n^4)$
- *Current:* What is the model learning? Link to friendly Graphs.
- *Current:* Applications to Shape Correspondence, Unaligned language translation

GIVENS FACTORIZATION OF UNITARY MATRICES

[with D. Folque (NYU)]

- Suppose we have a unitary matrix U (e.g. an eigenbasis) that we want to use extensively.
- Complexity of Matrix-vector multiplication: $\Theta(n^2)$ [Winograd]
- But structure on U can yield massive gains: FFT $\Theta(n \log n)$ [Tukey]

GIVENS FACTORIZATION OF UNITARY MATRICES

[with D. Folque (NYU)]

- Suppose we have a unitary matrix U (e.g. an eigenbasis) that we want to use extensively.
- Complexity of Matrix-vector multiplication: $\Theta(n^2)$ [Winograd]
- But structure on U can yield massive gains: FFT $\Theta(n \log n)$ [Tukey]
- General case?

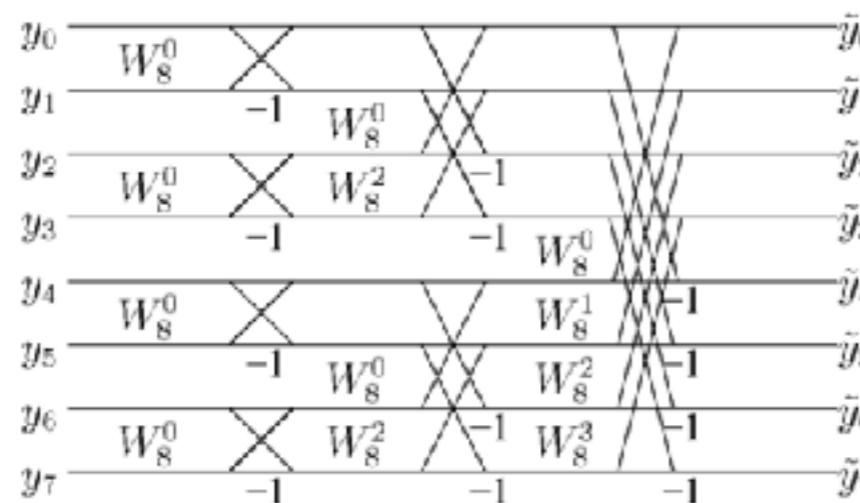
GIVENS FACTORIZATION OF UNITARY MATRICES

[with D. Folque (NYU)]

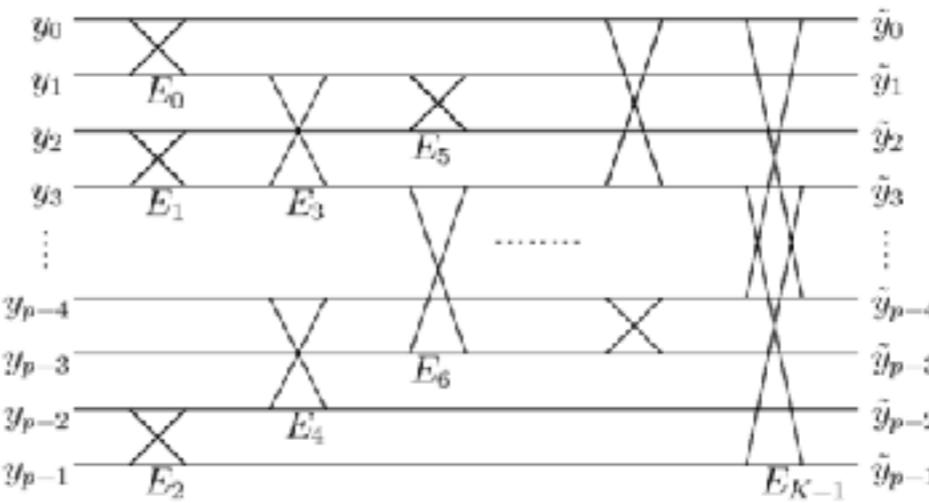
- We consider *Sparse Matrix Transforms* given by Givens plane rotations:

$$U = \prod_{k=1}^K \mathcal{O}(i_k, j_k, \alpha_k)$$

- $K = \frac{n(n-1)}{2}$ sufficient for any U FFT : $K = n \log n$
- NP-complete, non-commutative manifold-optimization problem.



(a) FFT



(b) SMT

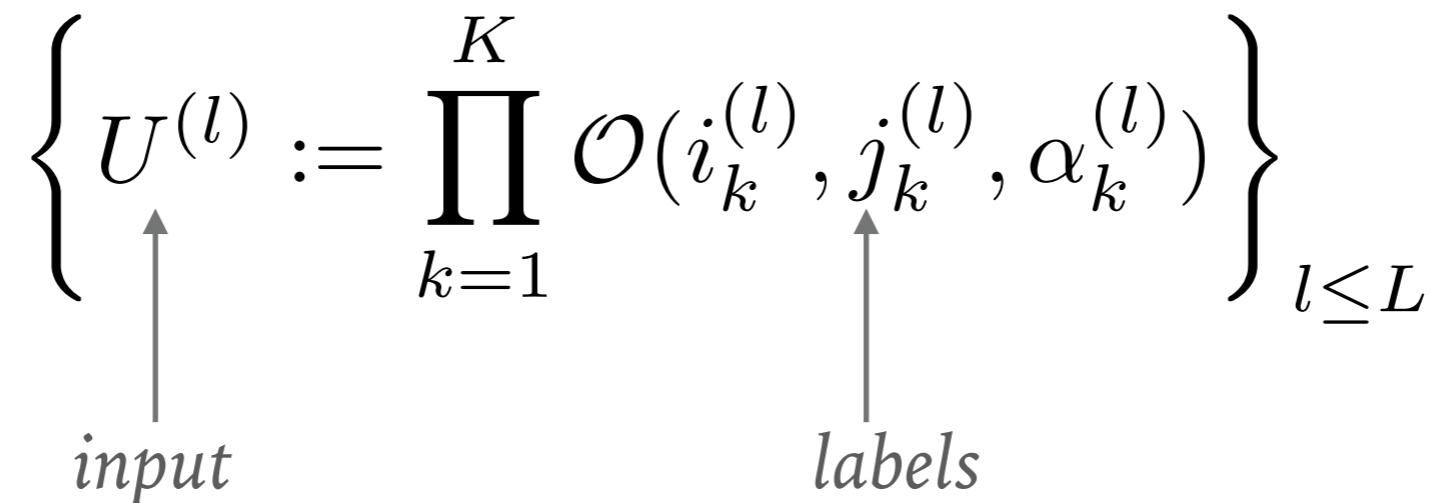
- Use GNN model to *learn* such factors.

GIVENS FACTORIZATION OF UNITARY MATRICES

[with D. Folque (NYU)]

- We consider a simple inverse-problem setup with planted solution:

$$\left\{ U^{(l)} := \prod_{k=1}^K \mathcal{O}(i_k^{(l)}, j_k^{(l)}, \alpha_k^{(l)}) \right\}_{l \leq L}$$



- Model: a GNN on the fully-connected graph, where we learn both edge and node features.
 - Uses multiset loss to account for partial permutation invariance.
 - Preliminary results: matching greedy algorithm.