

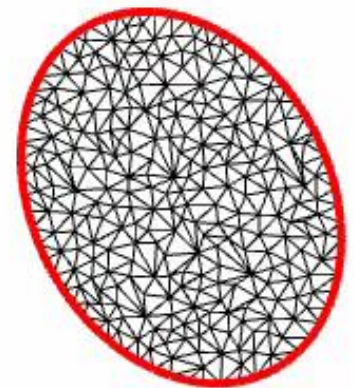
Spectral Perturbations and Generative Models in Geometric Deep Learning

Emanuele Rodolà

GLADIA lab



SAPIENZA
UNIVERSITÀ DI ROMA



European Research Council
Established by the European Commission

Overview

- ❑ Sci-fi
- ❑ Shape-from-spectrum (optimization)
- ❑ Shape-from-spectrum (learning-based)
- ❑ Spectral adversarial attacks



CAN ONE HEAR THE SHAPE OF A DRUM?

MARK KAC, The Rockefeller University, New York

To George Eugene Uhlenbeck on the occasion of his sixty-fifth birthday

“La Physique ne nous donne pas seulement l’occasion de résoudre des problèmes . . . , elle nous fait sentir la solution.” H. POINCARÉ.

Wave equation



$$\frac{\partial^2 u(x, y; t)}{\partial t^2} = \Delta u(x, y; t)$$

$$\Delta \phi_i = \lambda_i \phi_i$$

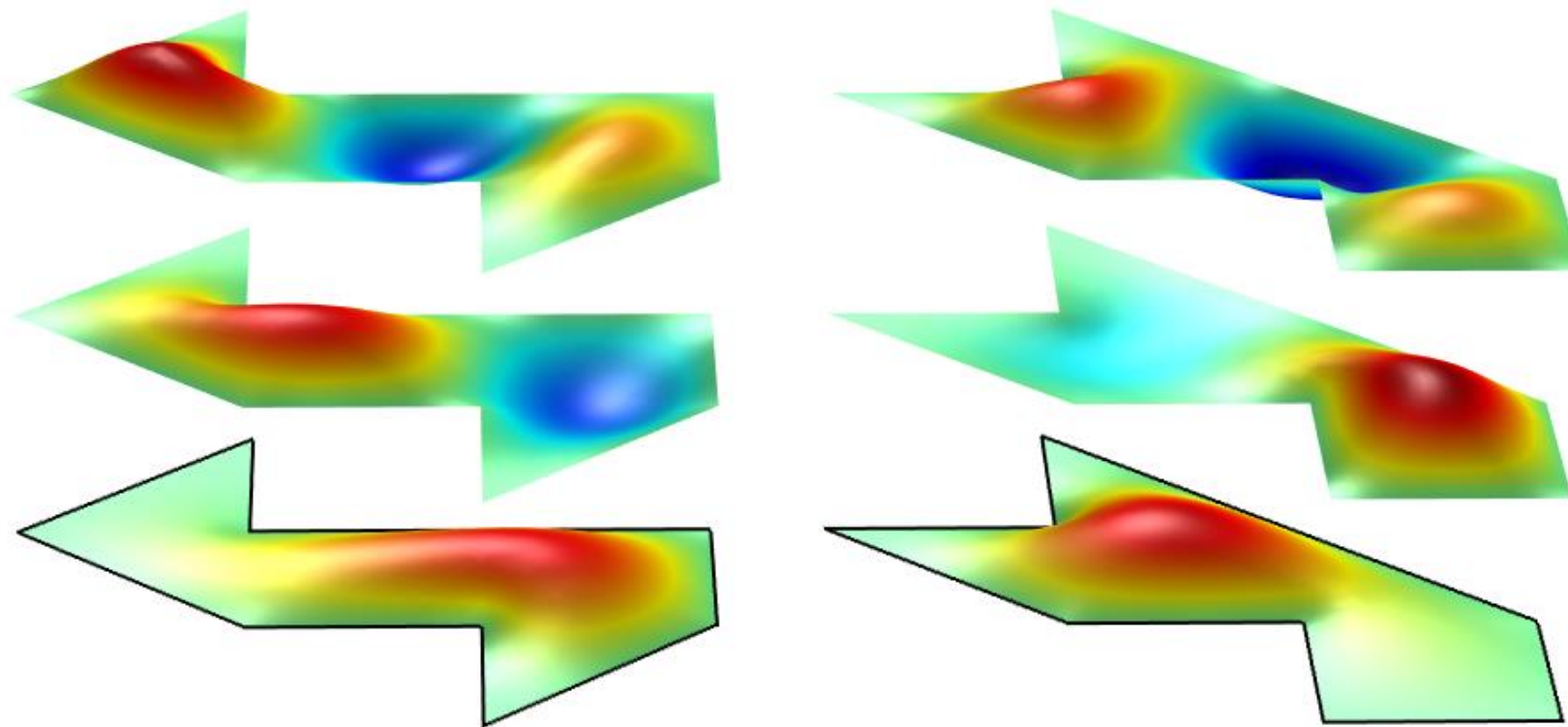
$$u(x, y; t) = \sum_j d_j \phi_j(x, y) \left(\cos \left(\sqrt{\lambda_j} t \right) + i \sin \left(\sqrt{\lambda_j} t \right) \right)$$

Isometry invariance



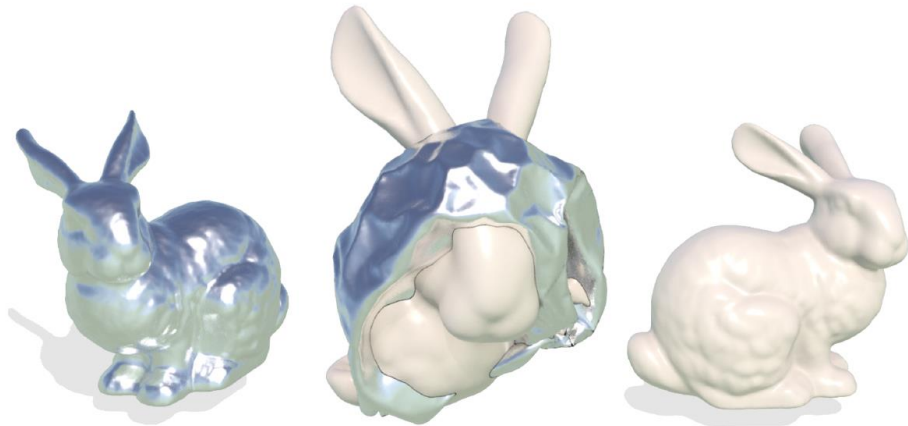
Isometric shapes have the same Laplacian eigenvalues

Isospectral \neq Isometric



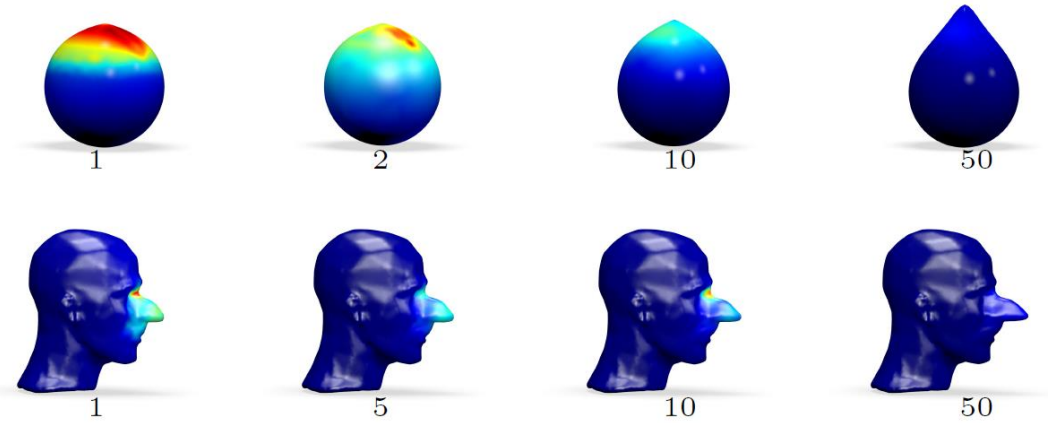
Existing approaches

Shape-from-metric



Chern et al 2018

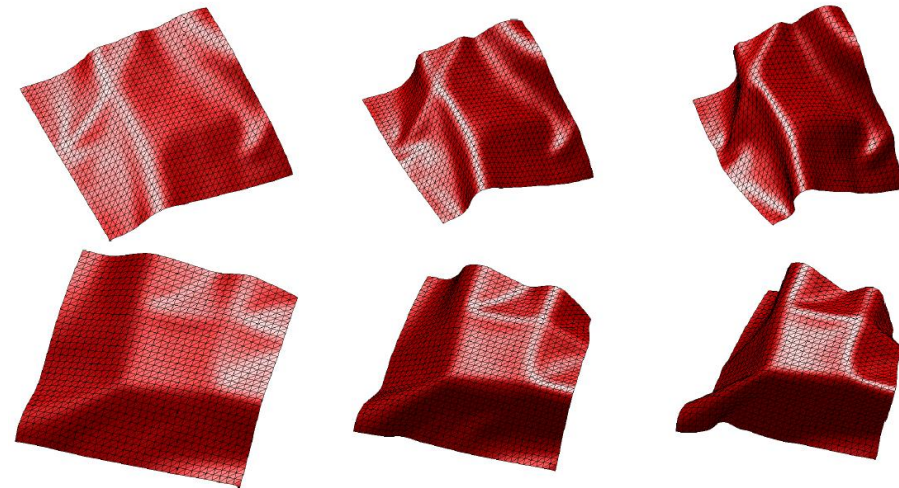
Shape-from-operator



Boscaini et al 2014



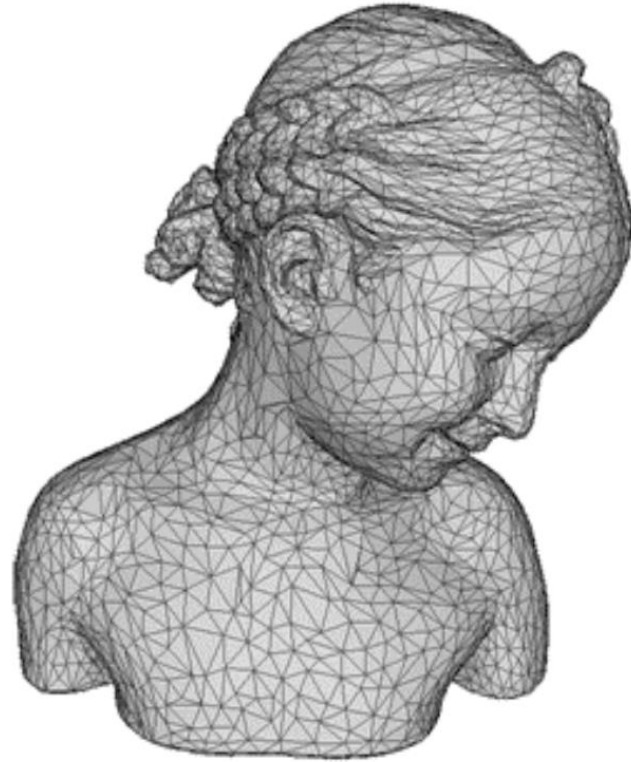
Borrelli et al 2012



Corman et al 2017

An empirical approach

Discrete setting



Isospectralization

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times d}} \|\boldsymbol{\lambda}(\Delta_X(\mathbf{X})) - \boldsymbol{\mu}\|_{\omega} + \rho_X(\mathbf{X})$$

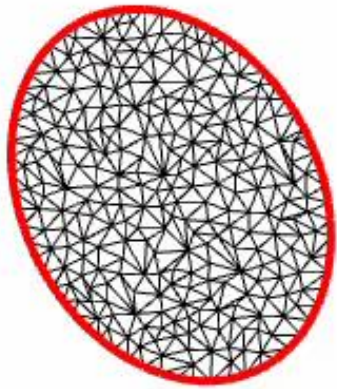
- **Data term:** weighted norm (frequency balancing)

$$\|\boldsymbol{\lambda} - \boldsymbol{\mu}\|_{\omega}^2 = \sum_{i=1}^k \frac{1}{\mu_i^2} (\lambda_i - \mu_i)^2$$

- **Regularizers** to promote smoothness / maximize volume
- **Input:** ≤ 30 eigenvalues
- **Optimization:** Nonlinear conjugate gradient with automatic differentiation

Example: Mickey-from-spectrum

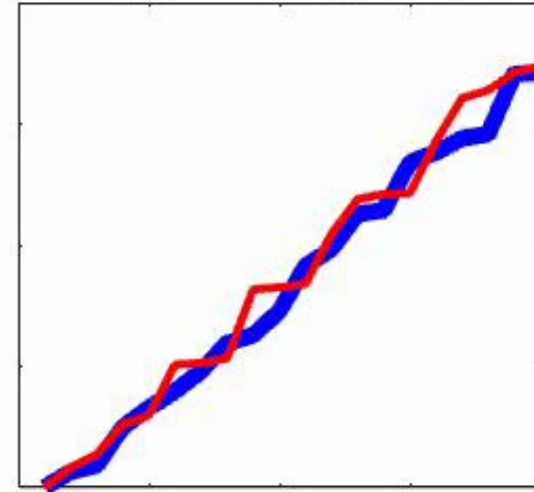
iter 1



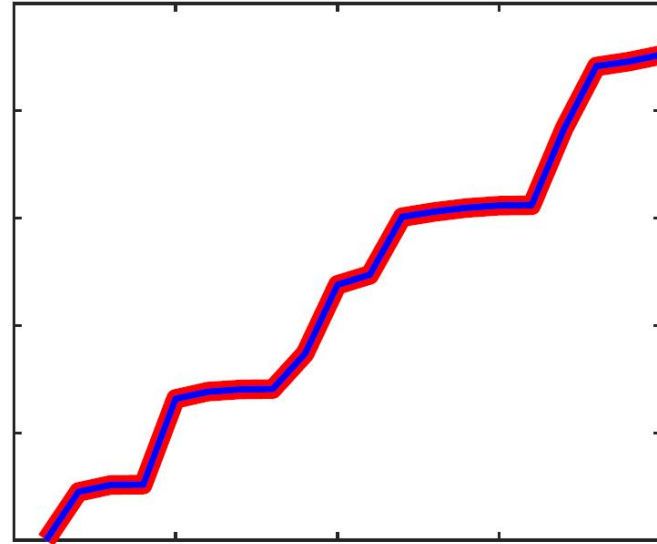
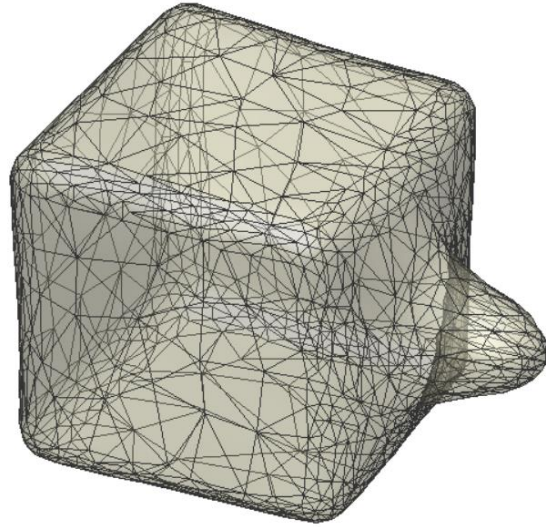
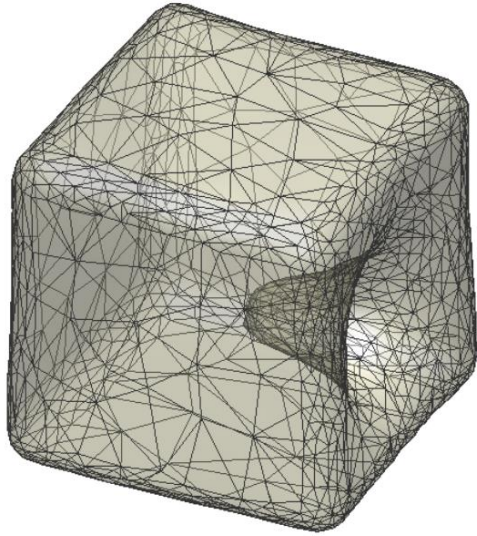
Target shape



Eigenvalues alignment



Geometric priors



Example: **Volume** regularizer to avoid isometric ambiguities

Geometric priors are not enough



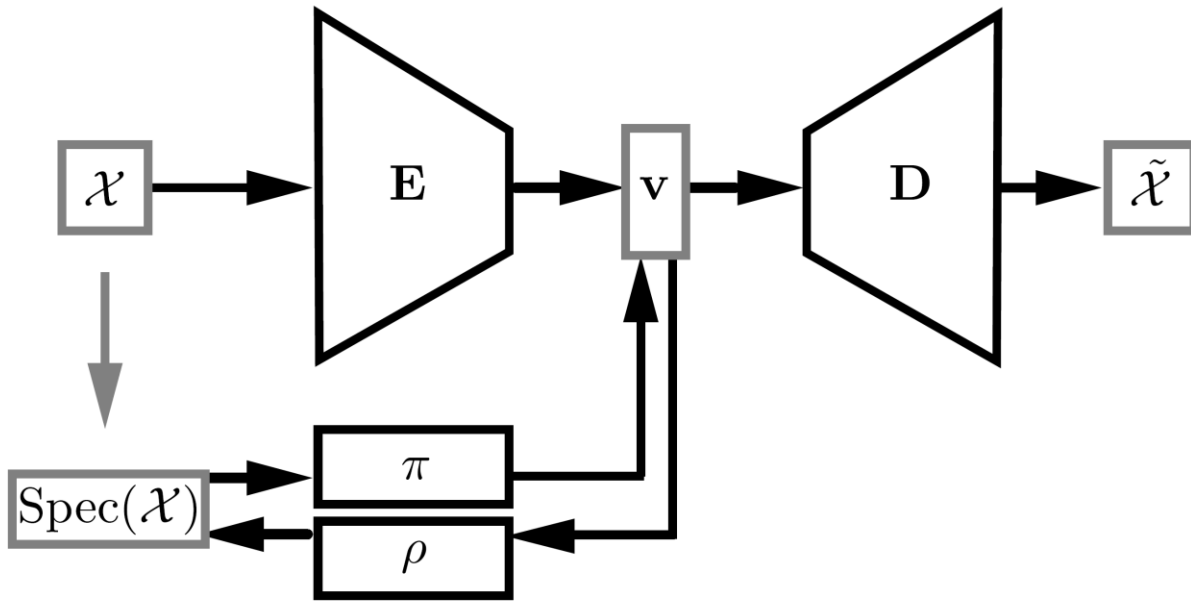
Do we really have to design regularizers?

Learn from data what is hard to
model axiomatically

Data-driven formulation

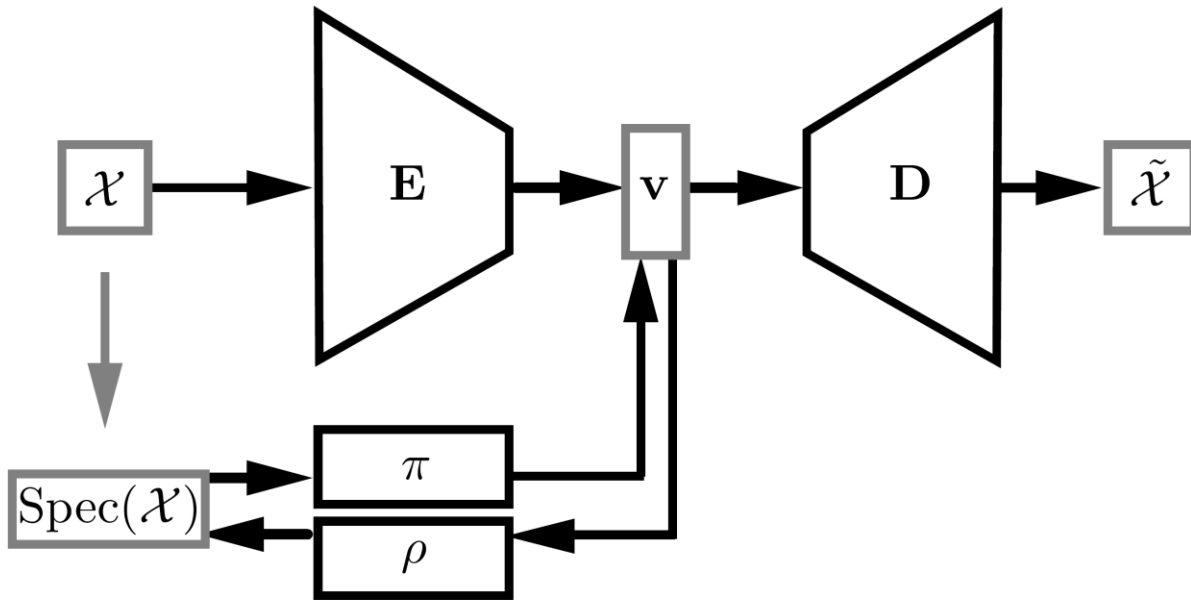
Latent space connections

AE-based learning model:



Latent space connections

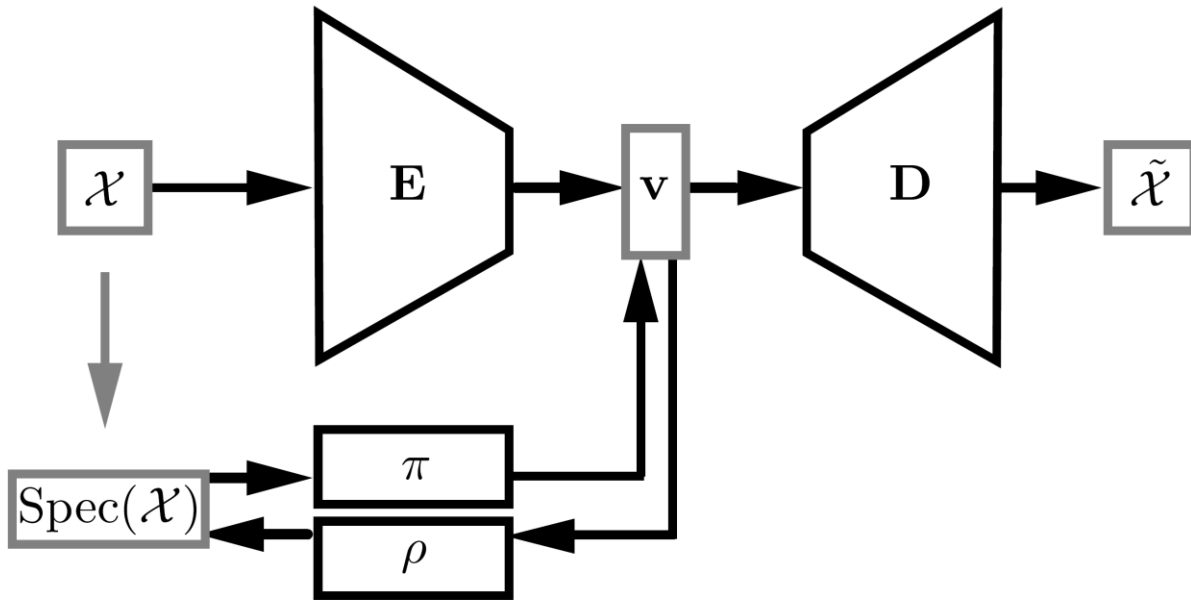
AE-based learning model:



$$l = l_{\mathcal{X}} + \alpha l_{\lambda}, \quad \text{with}$$
$$l_{\mathcal{X}} = \frac{1}{n} \|D(E(\mathbf{X})) - \mathbf{X}\|_F^2$$
$$l_{\lambda} = \frac{1}{k} (\|\pi(\boldsymbol{\lambda}) - E(\mathbf{X})\|_2^2 + \|\rho(E(\mathbf{X})) - \boldsymbol{\lambda}\|_2^2)$$

Latent space connections

AE-based learning model:



$$l = l_{\mathcal{X}} + \alpha l_{\lambda}, \quad \text{with}$$

$$l_{\mathcal{X}} = \frac{1}{n} \|D(E(\mathbf{X})) - \mathbf{X}\|_F^2$$

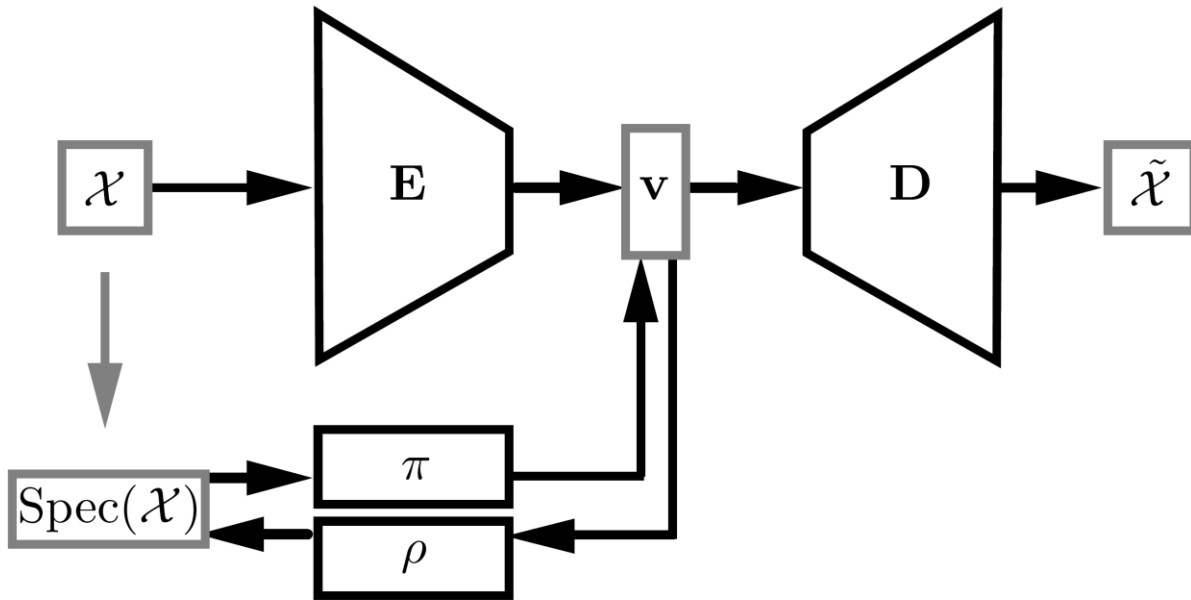
$$l_{\lambda} = \frac{1}{k} (\|\pi(\boldsymbol{\lambda}) - E(\mathbf{X})\|_2^2 + \|\rho(E(\mathbf{X})) - \boldsymbol{\lambda}\|_2^2)$$

The spectral loss enforces:

$$\rho \approx \pi^{-1}$$

Latent space connections

AE-based learning model:

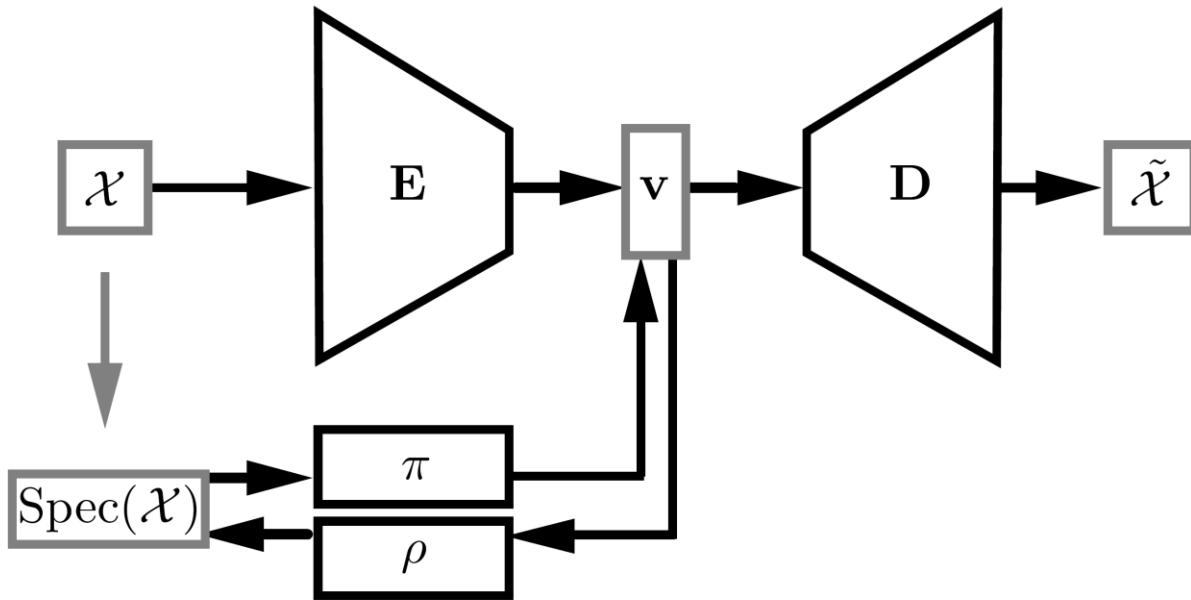


Remarks:

- **No back-propagation** through the eigen-decomposition

Latent space connections

AE-based learning model:

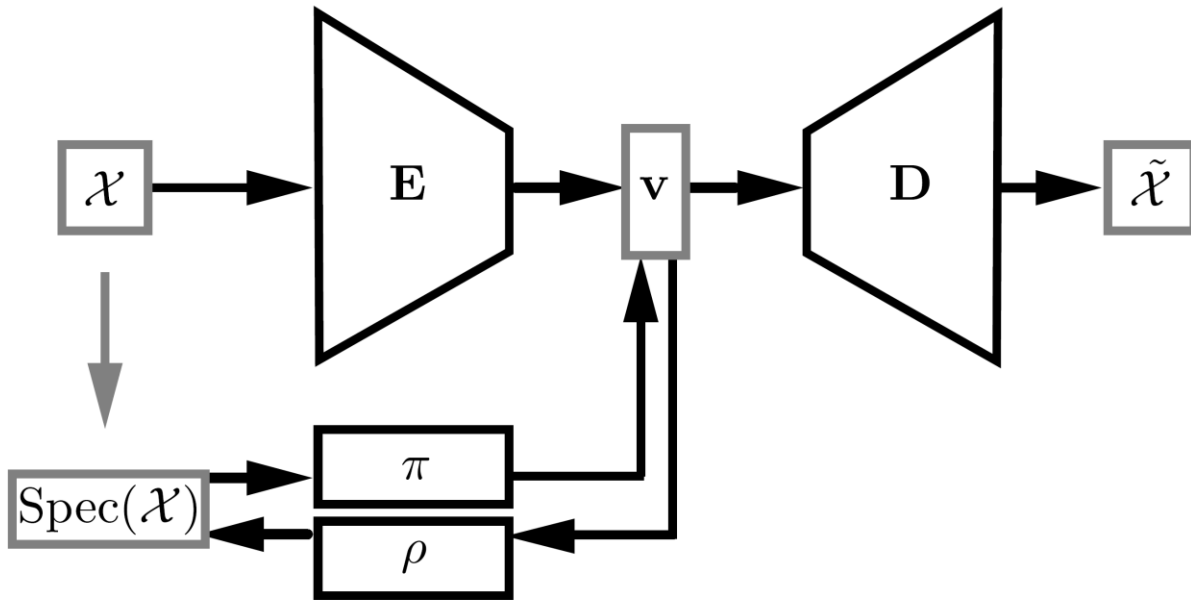


Remarks:

- **No back-propagation** through the eigen-decomposition
- The input spectrum can be **arbitrarily accurate**

Latent space connections

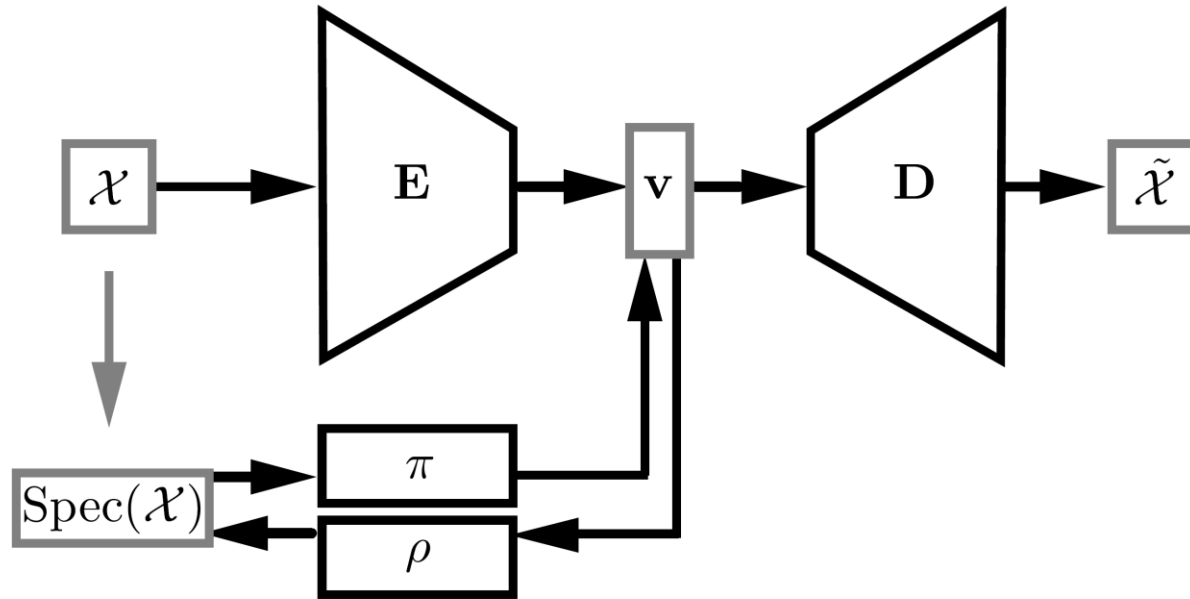
AE-based learning model:



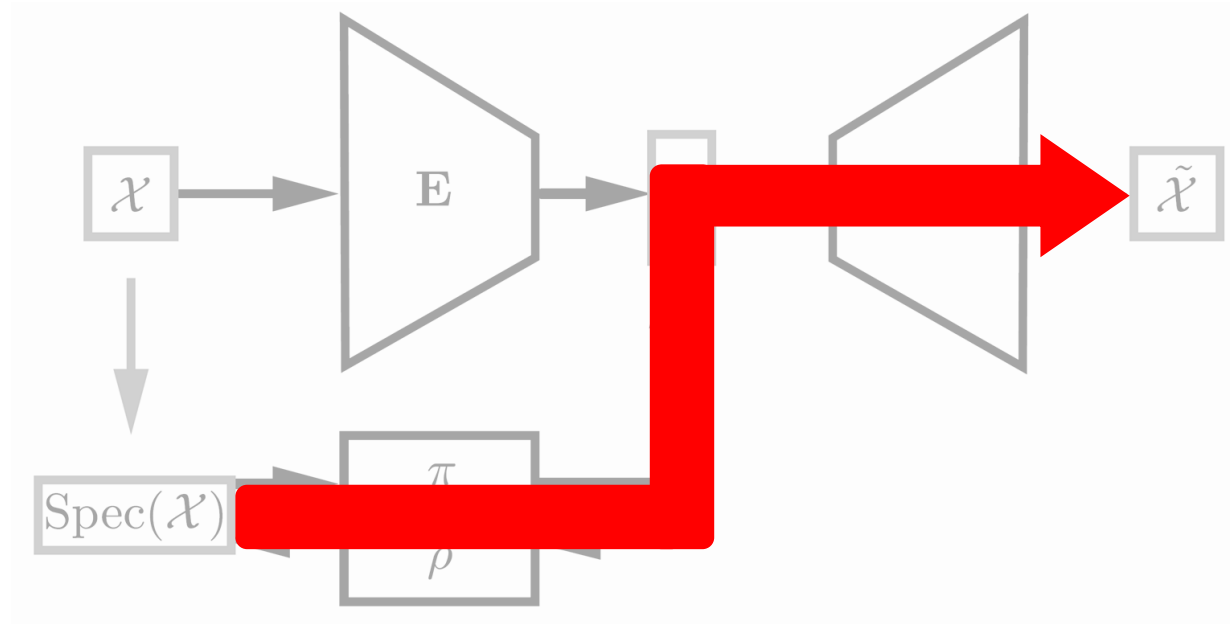
Remarks:

- **No back-propagation** through the eigen-decomposition
- The input spectrum can be **arbitrarily accurate**
- Admits **any AE** model (e.g. for point clouds, meshes, etc.)

Shape-from-spectrum reconstruction

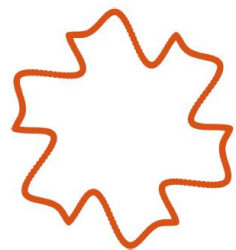


Shape-from-spectrum reconstruction

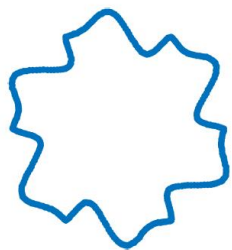


Examples

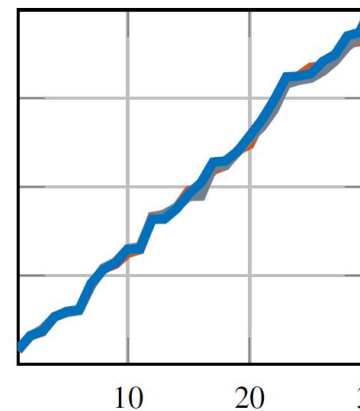
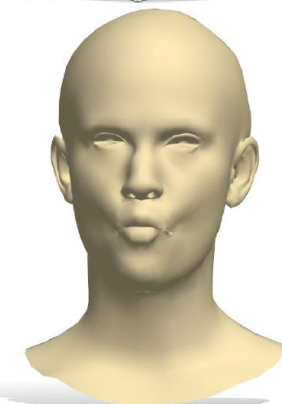
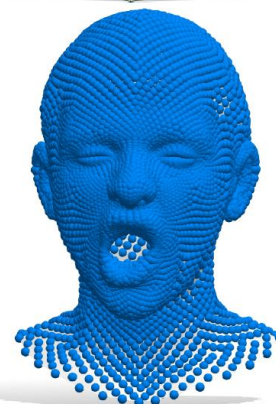
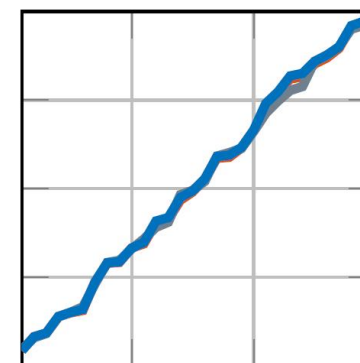
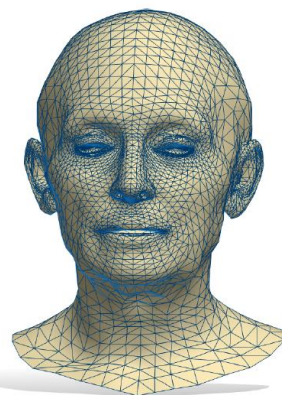
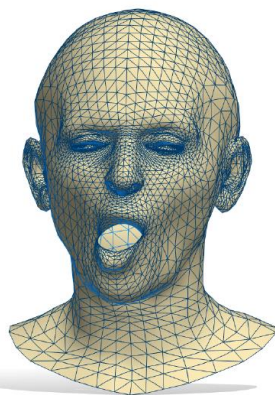
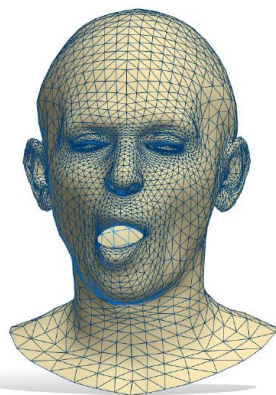
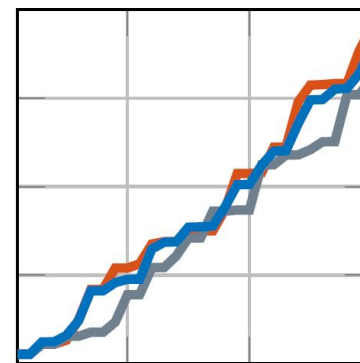
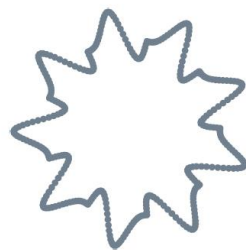
Target



Ours



NN



Application: Style transfer

$$\min_{\mathbf{v}} \|\text{Spec}(\mathcal{X}_{\text{style}}) - \rho(\mathbf{v})\|_2^2 + w \|\mathbf{v} - E(\mathcal{X}_{\text{pose}})\|_2^2$$

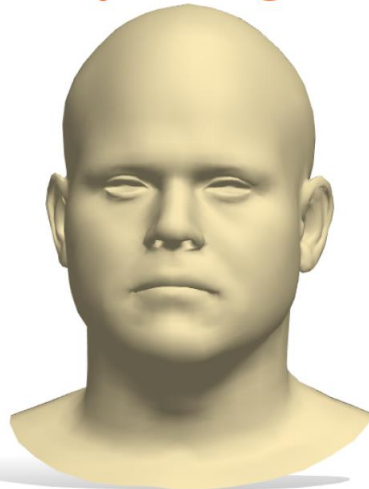
Application: Style transfer

$$\min_{\mathbf{v}} \|\text{Spec}(\mathcal{X}_{\text{style}}) - \rho(\mathbf{v})\|_2^2 + w \|\mathbf{v} - E(\mathcal{X}_{\text{pose}})\|_2^2$$

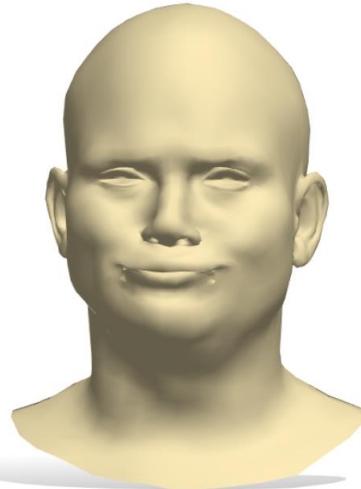
pose target



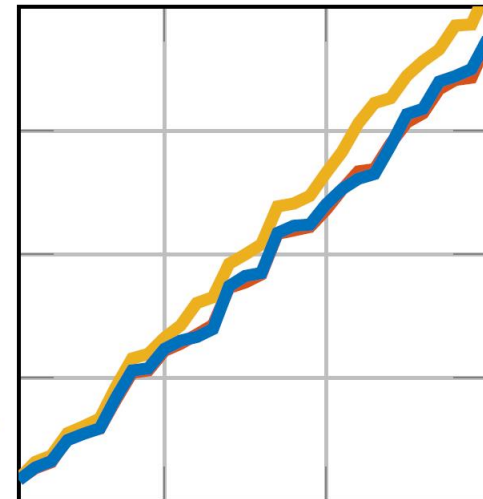
style target



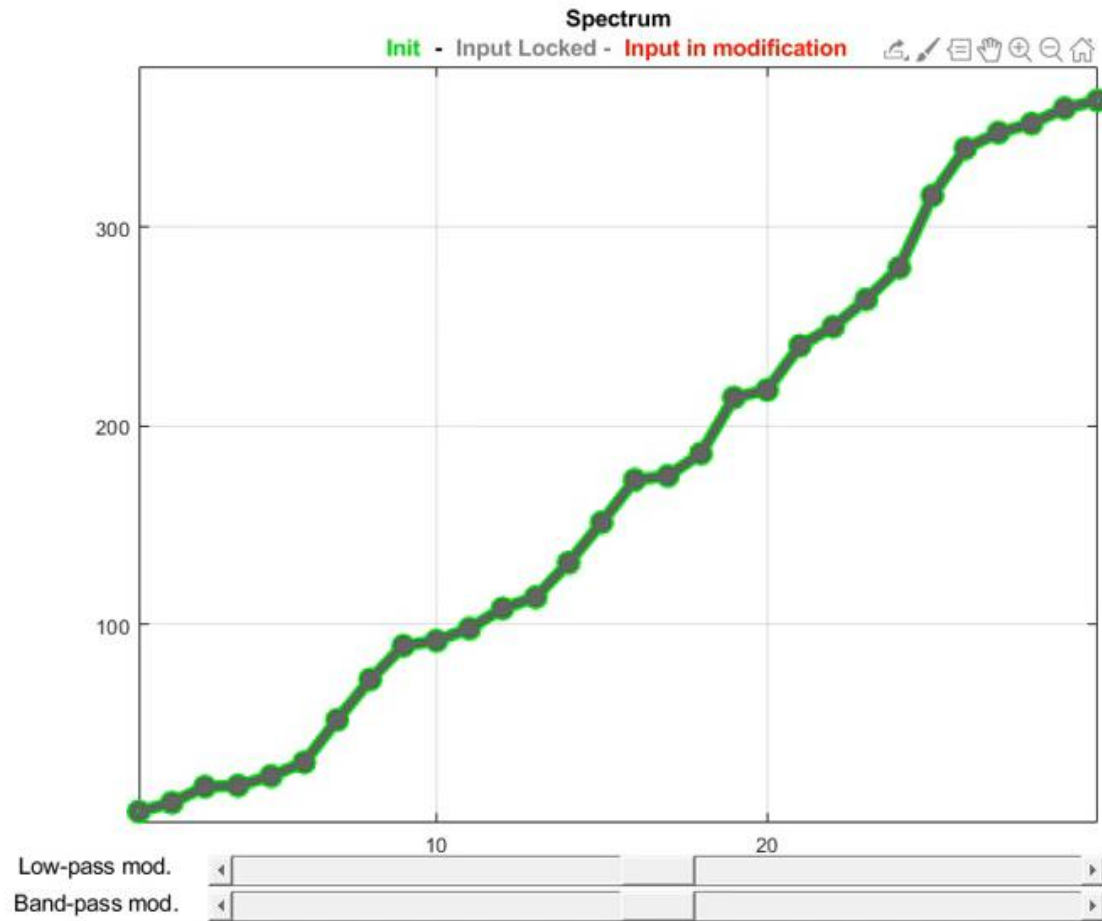
our result



eigenvalues



Application: Shape exploration

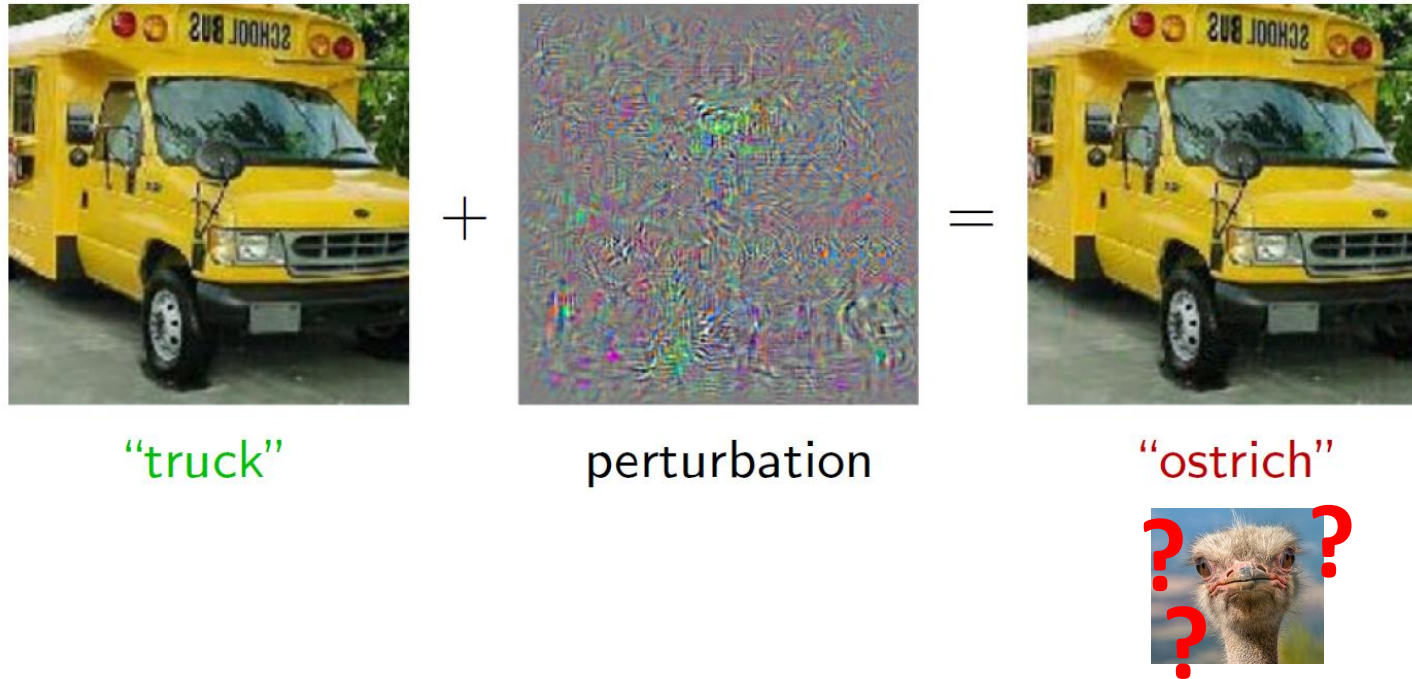


Output shape from $D(\pi(e))$



Adversarial attacks

Adversarial perturbations



The perturbation should be **undetectable** and can be explicitly optimized for.

Malicious attacks



“speed limit 50mph”

Example of a **malicious** attack on a visual classifier

Targeted attacks

Given an input sample \mathbf{x} , a classifier C , and a **target** class t , consider:

$$\begin{aligned} \min_{\mathbf{x}' \in [0,1]^n} \quad & \|\mathbf{x} - \mathbf{x}'\|_2^2 \\ \text{s.t.} \quad & C(\mathbf{x}') = t \end{aligned}$$

We call \mathbf{x}' an **adversarial example**.

Relax the difficult constraint to a penalty term:

$$\min_{\mathbf{x}' \in [0,1]^n} \|\mathbf{x} - \mathbf{x}'\|_2^2 + cL(\mathbf{x}', t)$$

Targeted attacks

A more general approach is given by:

$$\min_{\boldsymbol{\delta} \in [0,1]^n} d(\mathbf{x}, \mathbf{x} + \boldsymbol{\delta}) + c f(\mathbf{x} + \boldsymbol{\delta})$$

where the **perturbation** $\boldsymbol{\delta}$ appears explicitly, and d is some **distance**

f is such that $C(\mathbf{x} + \boldsymbol{\delta}) = t$ if and only if $f(\mathbf{x} + \boldsymbol{\delta}) \leq 0$.

$$f_1(x') = -\text{loss}_{F,t}(x') + 1$$

$$f_2(x') = (\max_{i \neq t} (F(x')_i) - F(x')_t)^+$$

$$f_3(x') = \text{softplus}(\max_{i \neq t} (F(x')_i) - F(x')_t) - \log(2)$$

$$f_4(x') = (0.5 - F(x')_t)^+$$

$$f_5(x') = -\log(2F(x')_t - 2)$$

$$f_6(x') = (\max_{i \neq t} (Z(x')_i) - Z(x')_t)^+$$

$$f_7(x') = \text{softplus}(\max_{i \neq t} (Z(x')_i) - Z(x')_t) - \log(2)$$

See:

Carlini and Wagner, 2016

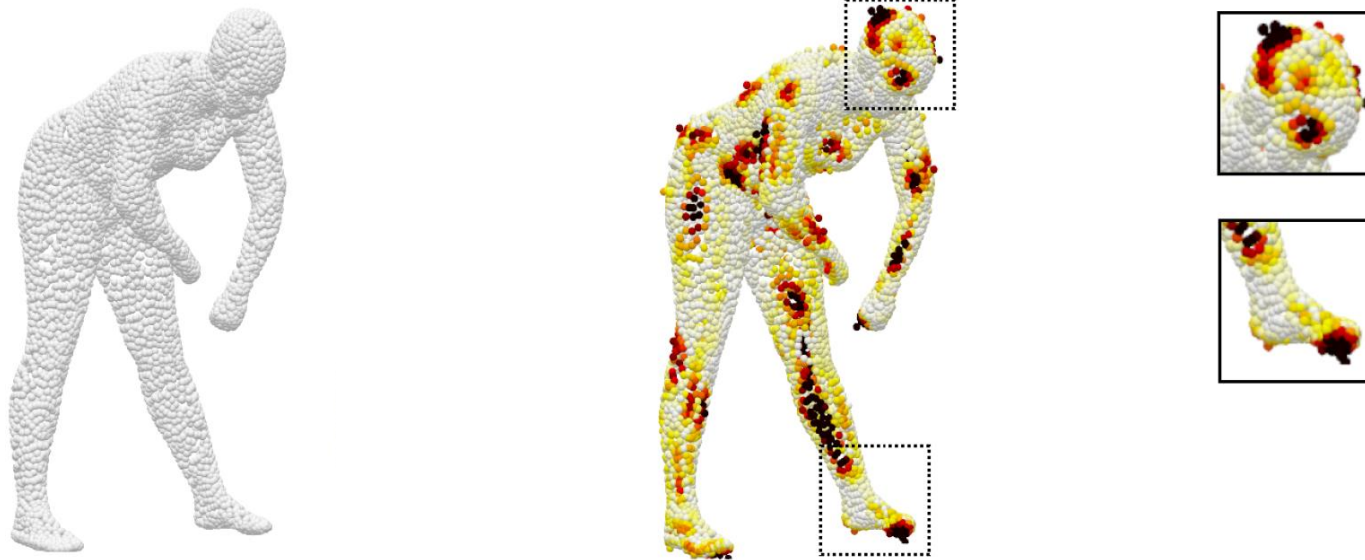
“Towards evaluating the robustness of neural networks”

Surface attacks

A perturbation \mathbf{V} is a **displacement field**:

$$\mathbf{X}' = \mathbf{X} + \mathbf{V}$$

Arbitrary displacement can lead to noticeable adversarial **jittering**:

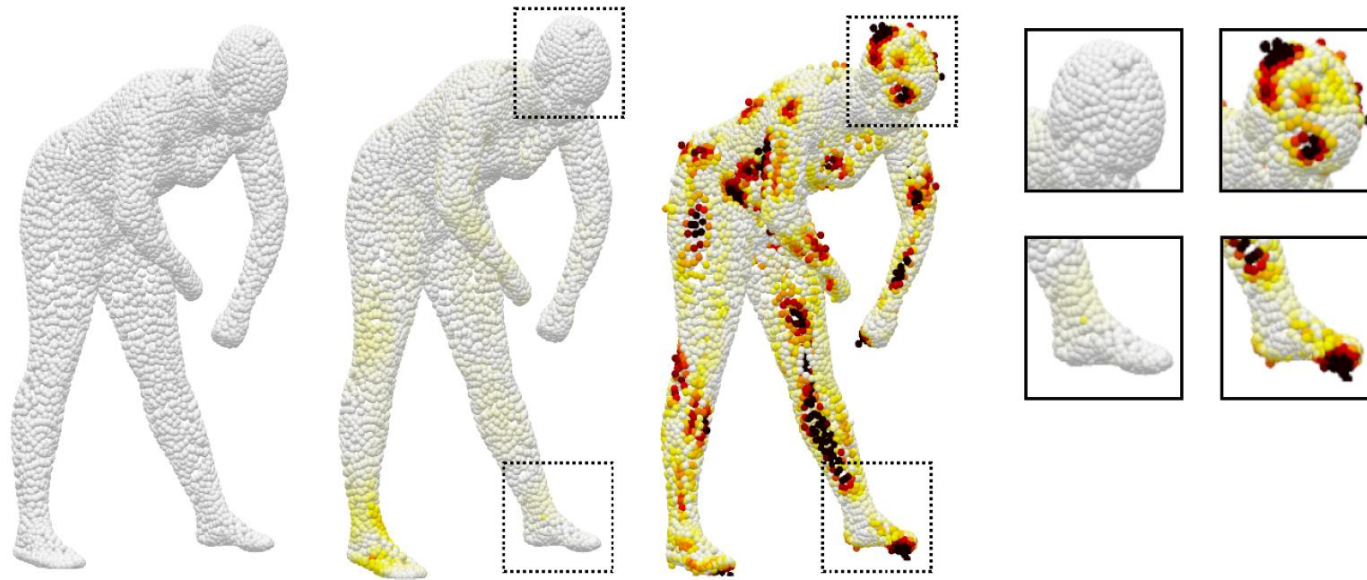


Surface attacks

A perturbation \mathbf{V} is a **displacement field**:

$$\mathbf{X}' = \mathbf{X} + \mathbf{V}$$

Arbitrary displacement can lead to noticeable adversarial **jittering**:



Idea: Regularize the displacement to make it less noticeable.

Band-limited perturbations

Represent the perturbation in the **truncated** Laplacian eigenbasis:

$$\mathbf{V} = \Phi \mathbf{v}$$

Theorem 1 [ABK15] For any given choice of $k \geq 1$ and any function $f \in \mathcal{F}(\mathcal{X})$, the inequality:

$$\|f - \sum_{i=1}^k \langle \psi_i, f \rangle \psi_i\|^2 \leq \alpha \frac{\|\nabla f\|^2}{\lambda_{k+1}} \quad (4)$$

holds for $\alpha = 1$ whenever one chooses ψ_i to be the Laplacian eigenfunctions, while tightening the bound with $0 \leq \alpha < 1$ is not possible for *any* sequence of orthogonal functions $\{\psi_i \in \mathcal{F}(\mathcal{X})\}$.

Band-limited perturbations

Represent the perturbation in the **truncated** Laplacian eigenbasis:

$$\mathbf{V} = \Phi \mathbf{v}$$

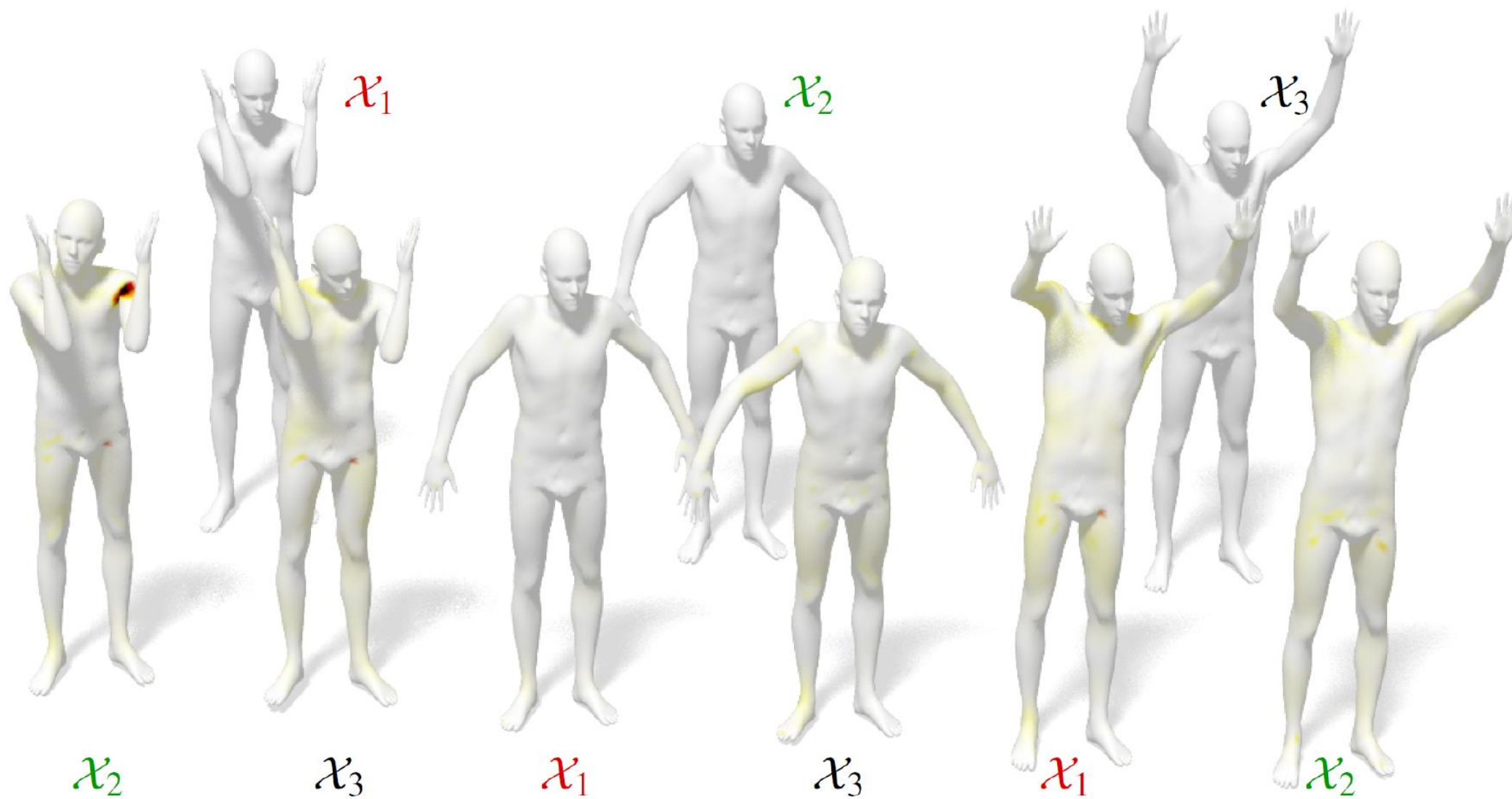
Theorem 1 [ABK15] For any given choice of $k \geq 1$ and any function $f \in \mathcal{F}(\mathcal{X})$, the inequality:

$$\|f - \sum_{i=1}^k \langle \psi_i, f \rangle \psi_i\|^2 \leq \alpha \frac{\|\nabla f\|^2}{\lambda_{k+1}} \quad (4)$$

holds for $\alpha = 1$ whenever one chooses ψ_i to be the Laplacian eigenfunctions, while tightening the bound with $0 \leq \alpha < 1$ is not possible for *any* sequence of orthogonal functions $\{\psi_i \in \mathcal{F}(\mathcal{X})\}$.

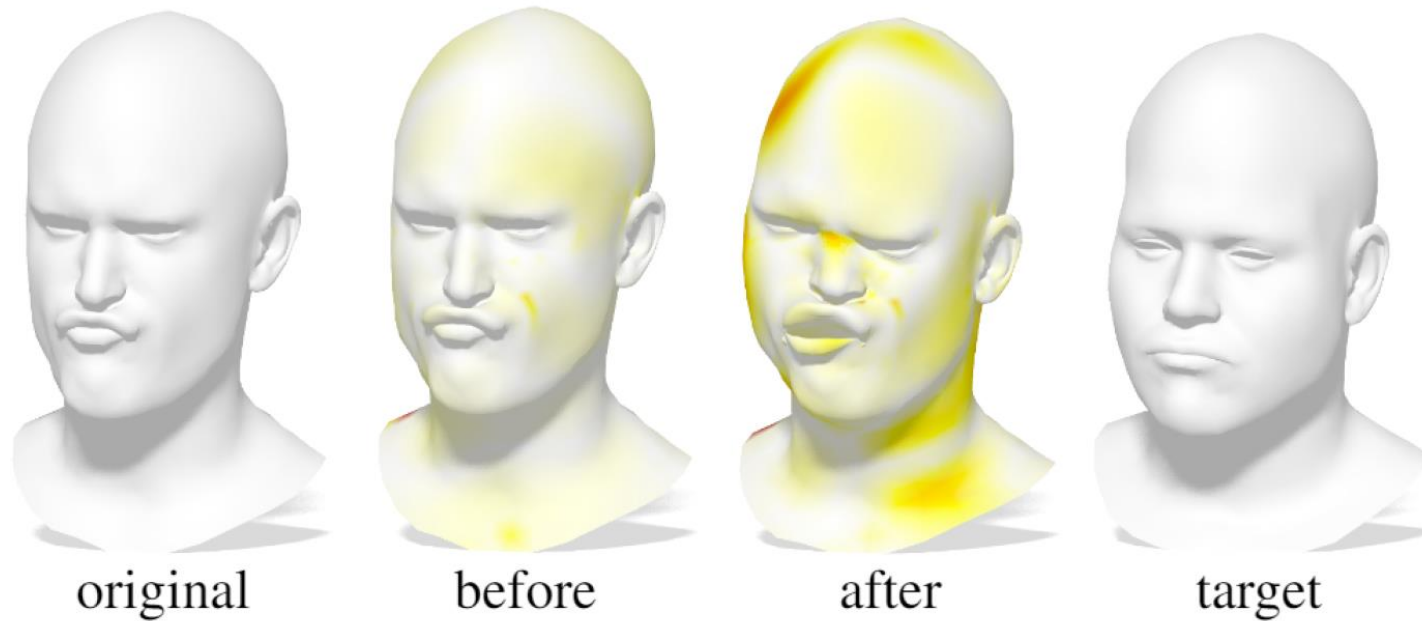


Examples



Adversarial training

Using adversarial examples as **training** data improves the **robustness** of the attacked learning model:

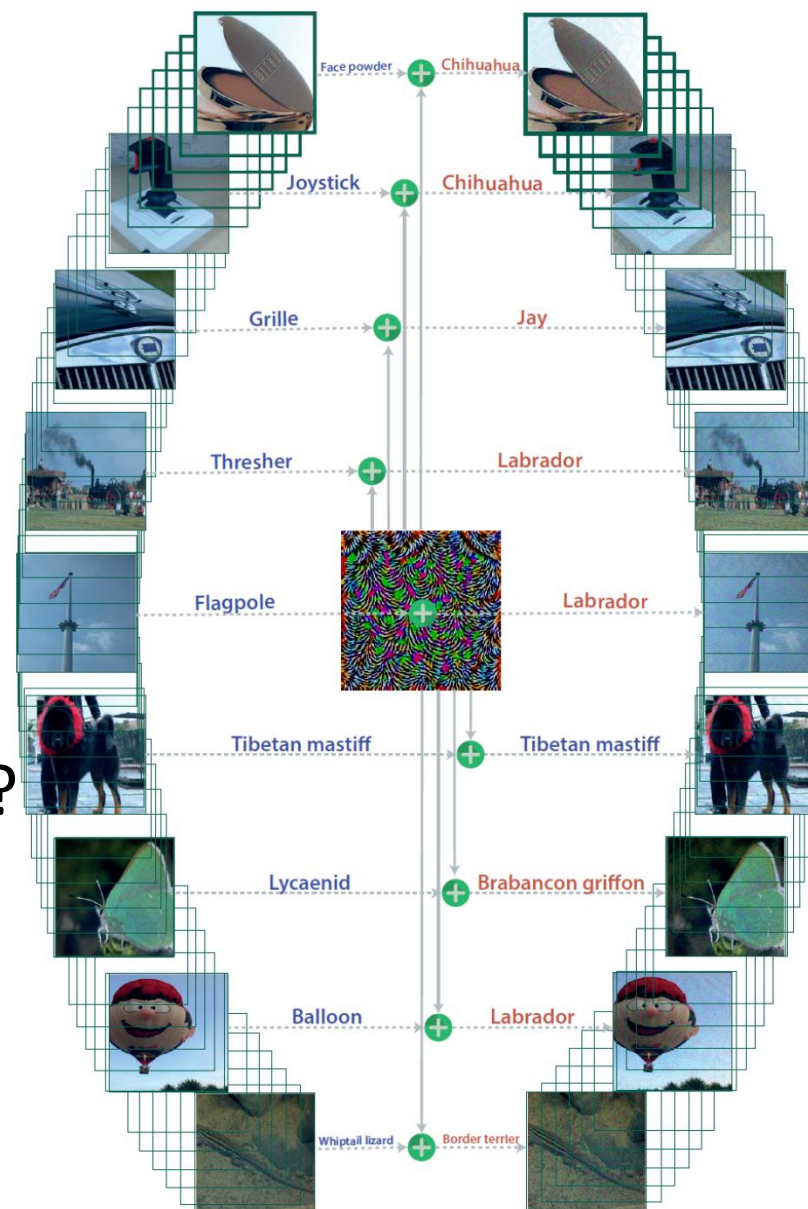


Universal perturbations

Image-agnostic perturbations are known to exist.

What about surfaces and point clouds?

Can we even define **a single spatial perturbation** for an entire collection of shapes?

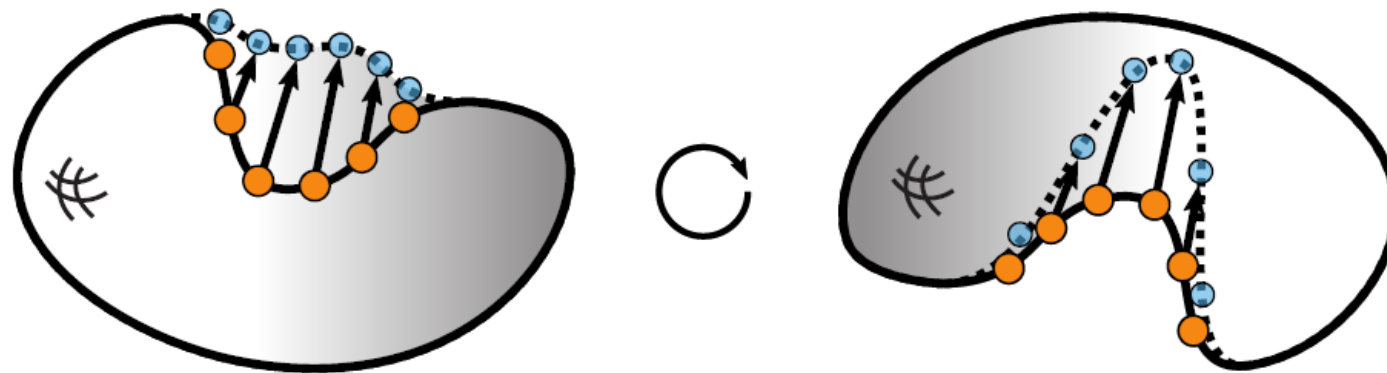


Universal spatial perturbations

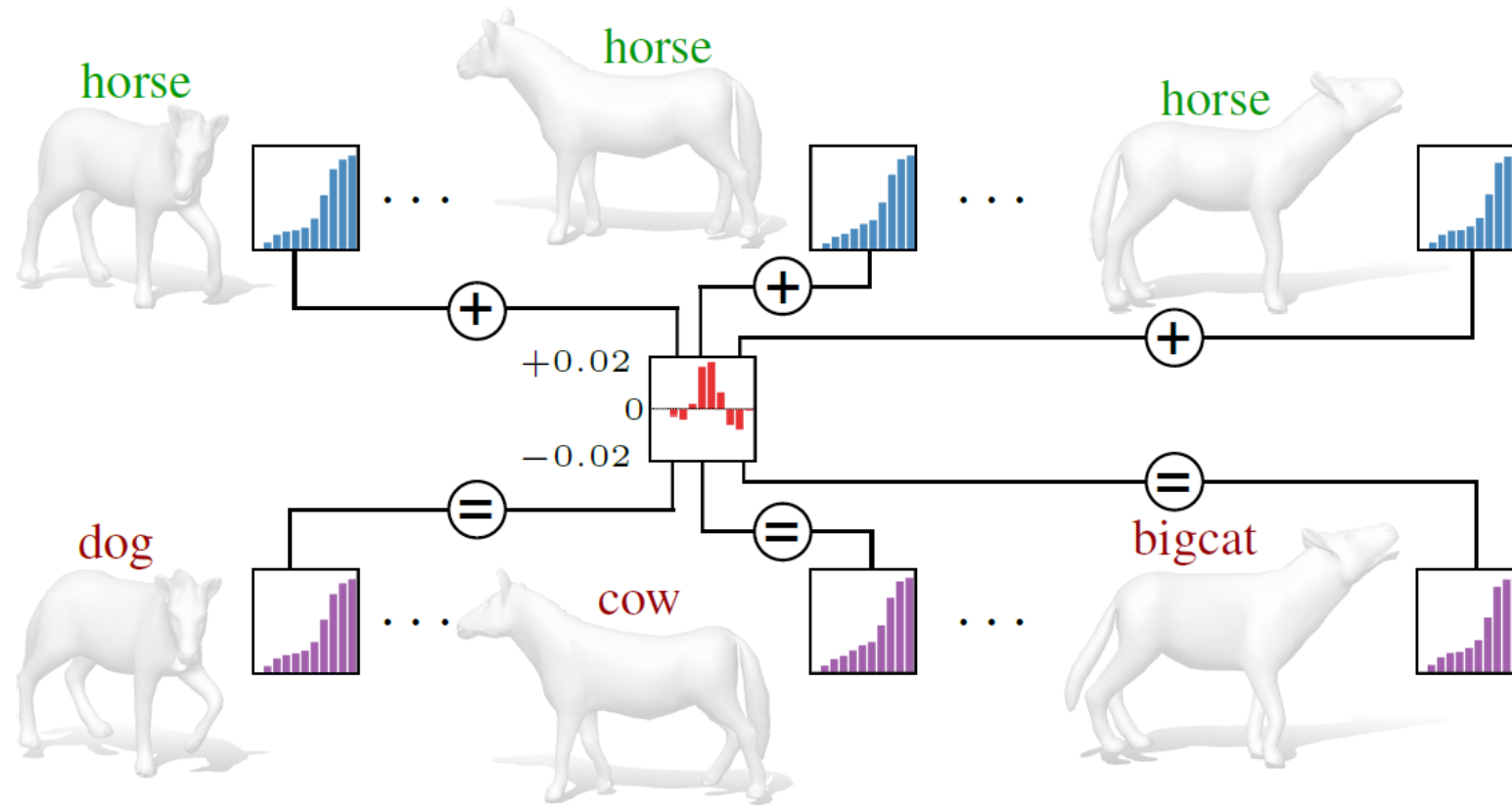
No!

Each shape is its own domain.

Any spatial perturbation only applies to the domain where it is defined.



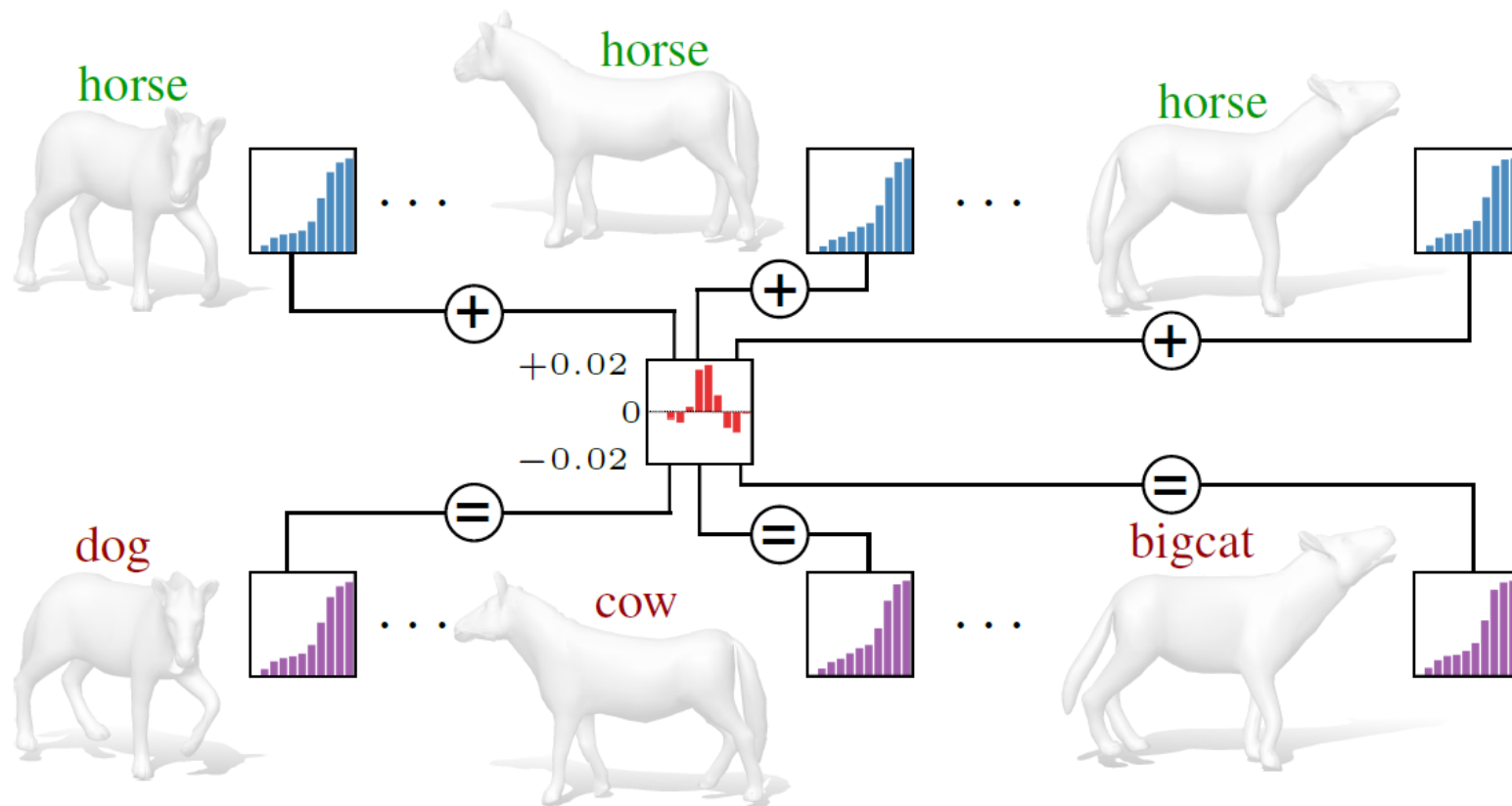
Universal spectral perturbations



Universal spectral perturbations

$$\min_{\substack{\rho \in \mathbb{R}^k \\ \mathcal{P}_i}} \sum_{X_i \in \mathcal{S}} \|\sigma(X_i)(1 + \rho) - \sigma(\mathcal{P}_i(X_i))\|_2^2$$

s.t. $\mathcal{C}(\mathcal{P}_i(X_i)) \neq \mathcal{C}(X_i) \quad \forall X_i \in \mathcal{S}$

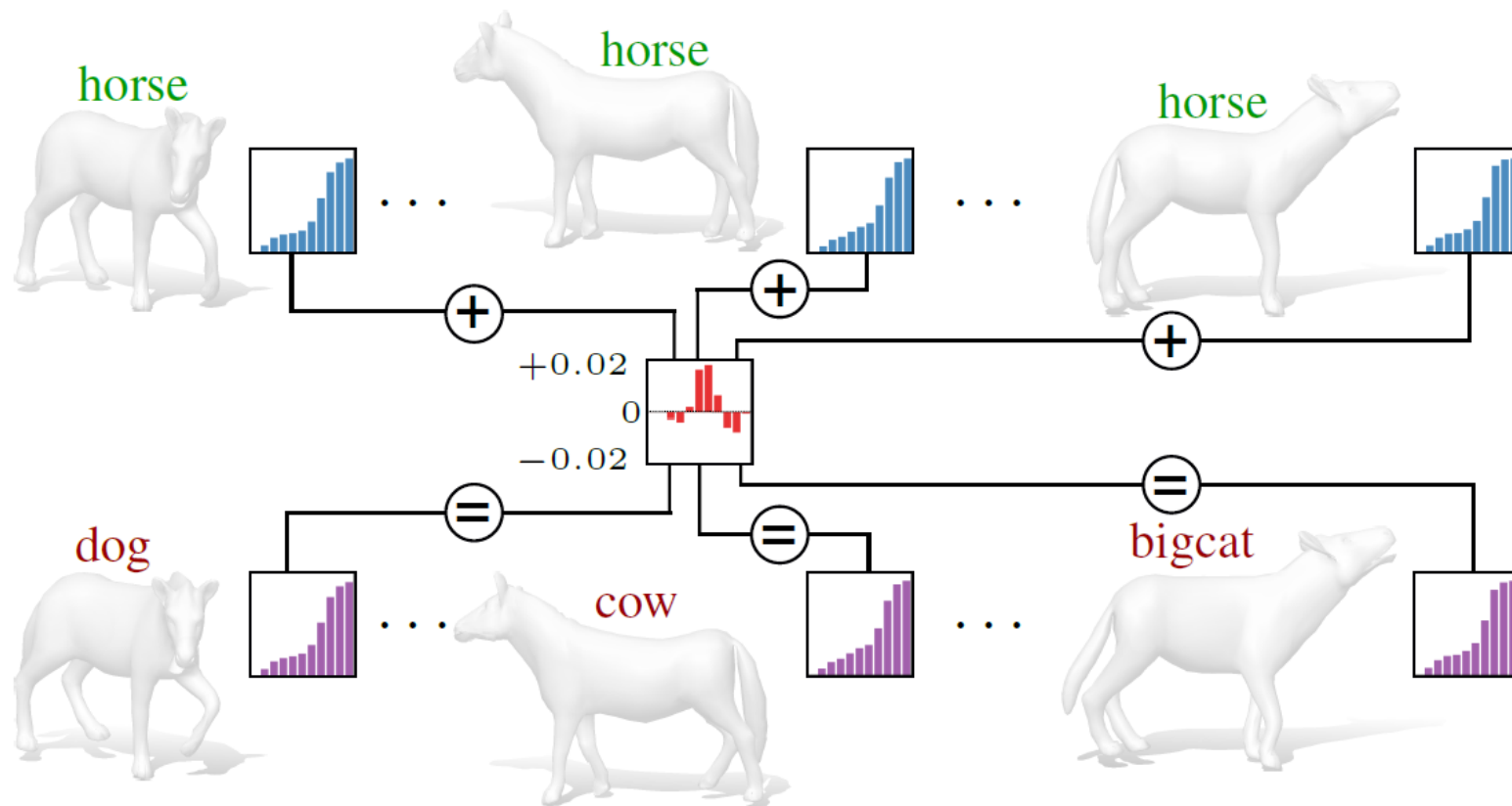


Universal spectral perturbations

$$\min_{\substack{\rho \in \mathbb{R}^k \\ \mathcal{P}_i}} \sum_{X_i \in \mathcal{S}} \|\sigma(X_i)(1 + \rho) - \sigma(\mathcal{P}_i(X_i))\|_2^2$$

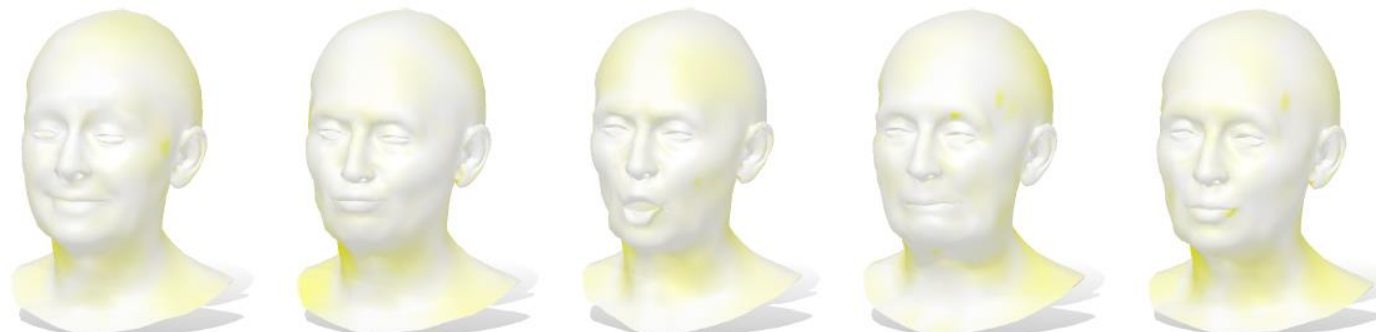
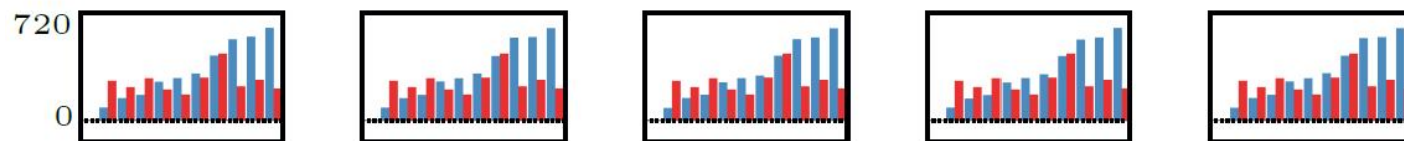
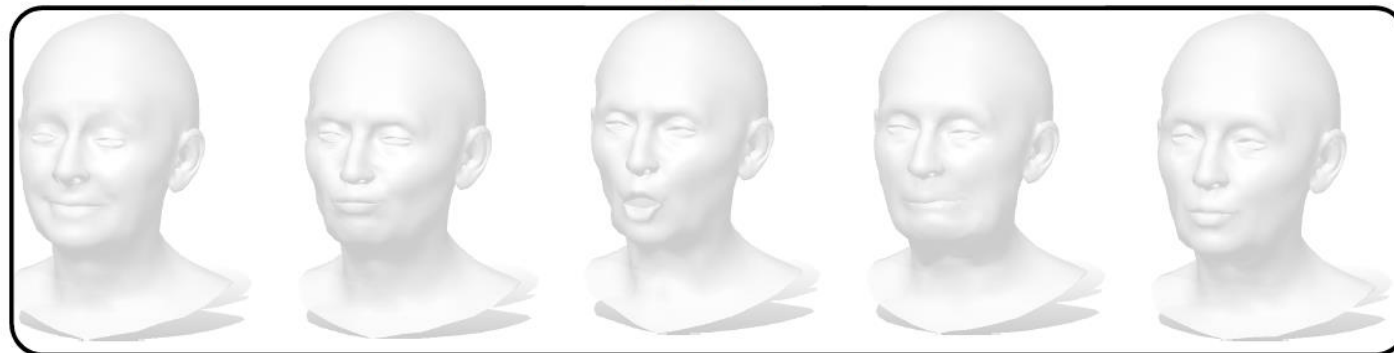
s.t. $\mathcal{C}(\mathcal{P}_i(X_i)) \neq \mathcal{C}(X_i) \quad \forall X_i \in \mathcal{S}$

$$\begin{array}{ccc} X_i & \xrightarrow{\sigma} & (\lambda^i) \\ \mathcal{P}_i \downarrow & & \downarrow \rho \\ \tilde{X}_i & \xrightarrow{\sigma} & (\tilde{\lambda}^i) \end{array}$$



Examples

ID 4



ID 2

ID 2

ID 2

ID 2

ID 2

Thanks for ~~hearing~~ listening!

