# Unsupervised Learning and Inverse Problems with Deep Neural Networks
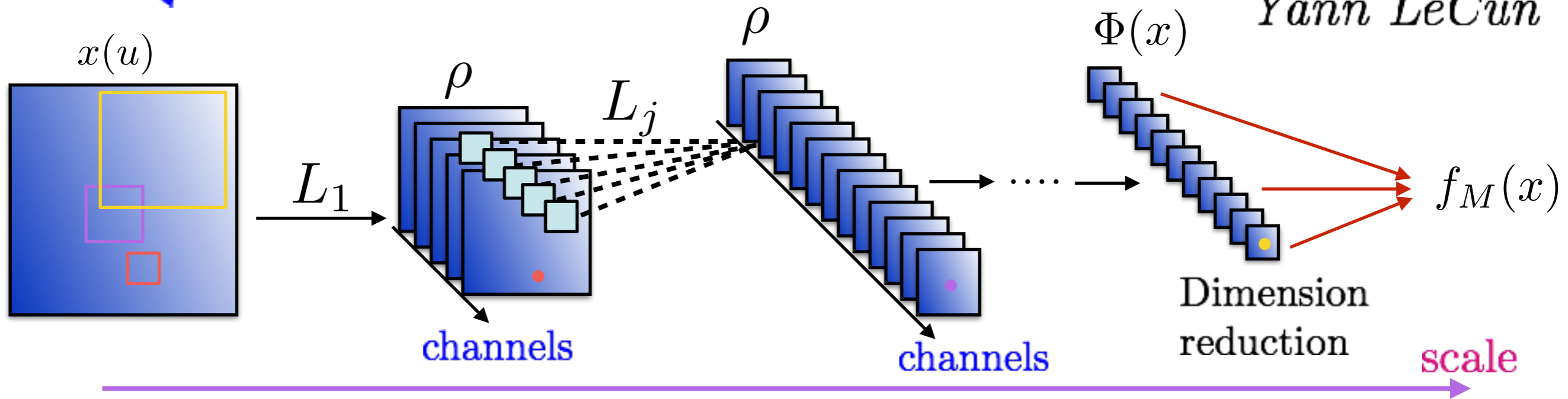
*Joan Bruna, Stéphane Mallat,*
*Ivan Dokmanic, Martin de Hoop*

**École Normale Supérieure**
www.di.ens.fr/data

# Deep Convolutional Networks



$Yann\ LeCun$

$L_j$ is a sum of spatial convolutions across channels, subsampling

$\rho(u)$ is a scalar non-linearity: $\max(u, 0)$ or $|u|$ or ...

**Part I** Architecture Simplification: wavelet scattering

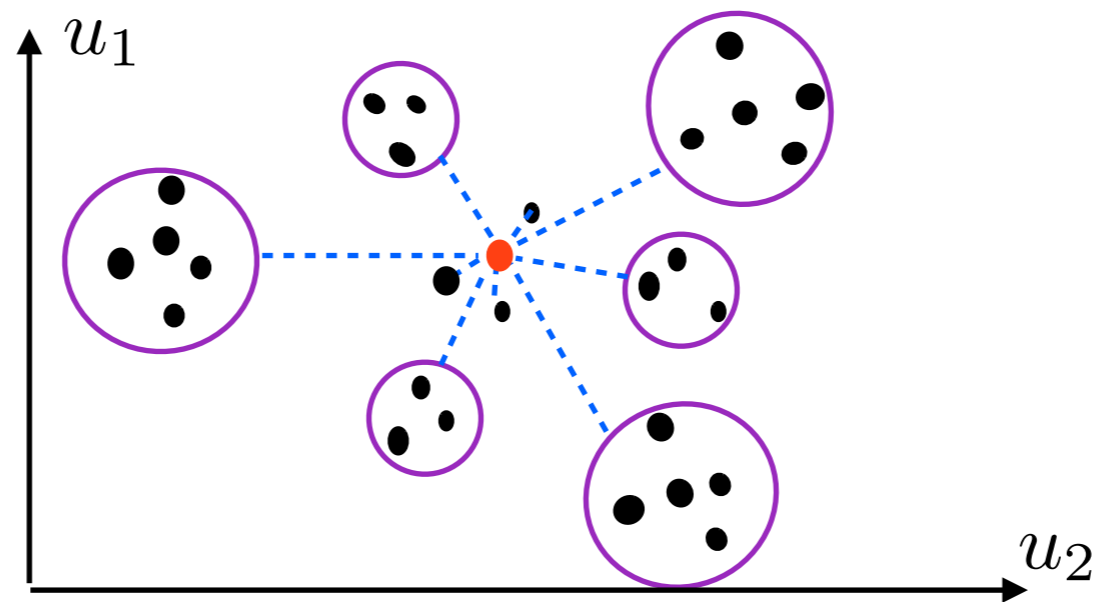**Part II** Unsupervised learning: generative models

**Part III** Inverse problems

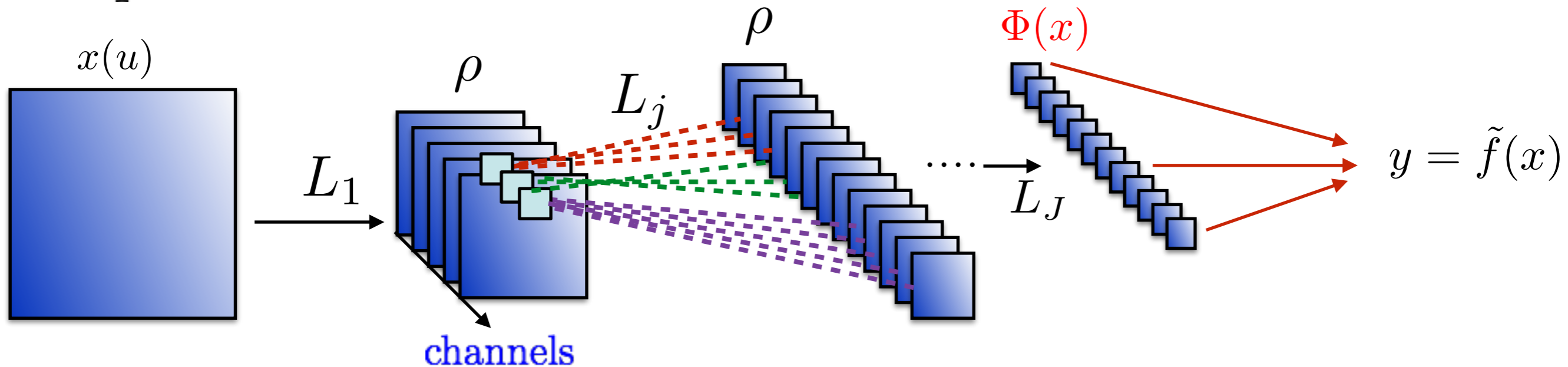- Why can we learn despite the curse of dimensionality ?

  Multiscale structures/interactions

Interactions de $d$ variables $x(u)$: pixels, particules, agents...



Regroupement of $d$ interactions in $O(\log d)$

Simplified architecture:
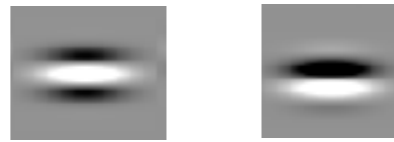


Cascade of convolutions: no channel connections
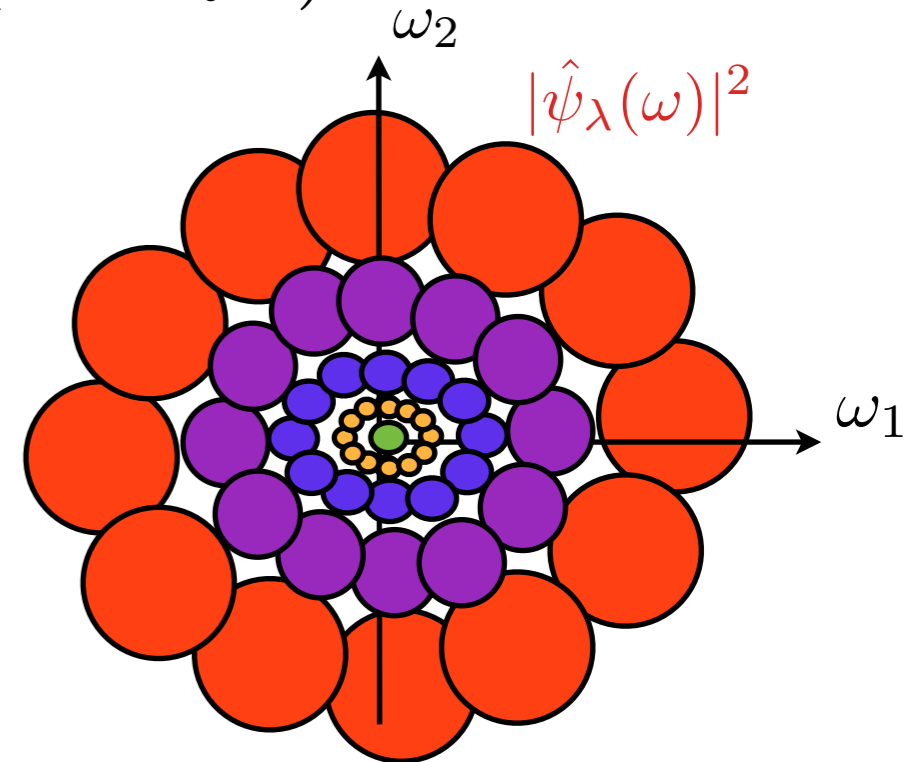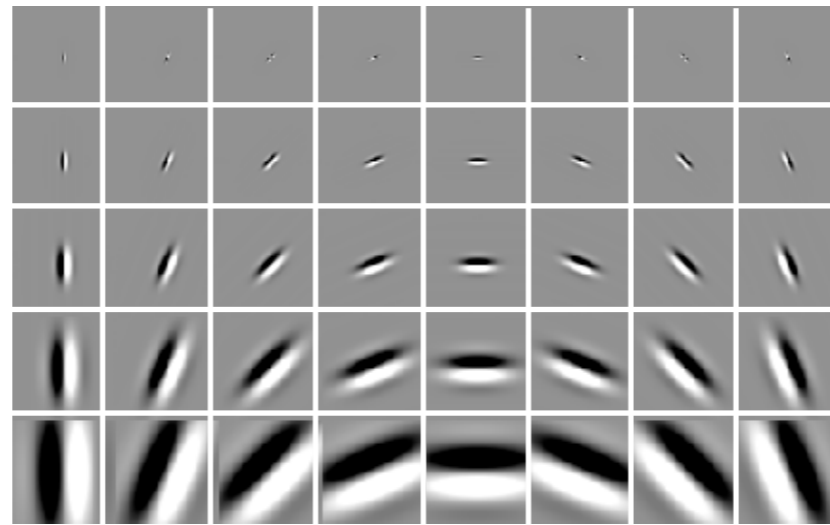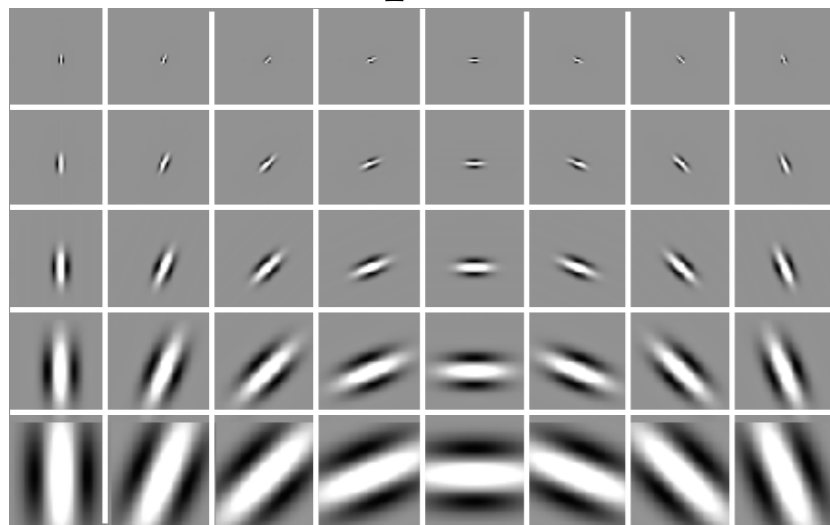
predefined wavelet filters

- Wavelet filter $\psi(u)$: 

rotated and dilated: $\psi_{2^j,\theta}(u) = 2^{-j}\,\psi(2^{-j}r_\theta u)$

real parts      imaginary parts      $|\hat{\psi}_\lambda(\omega)|^2$
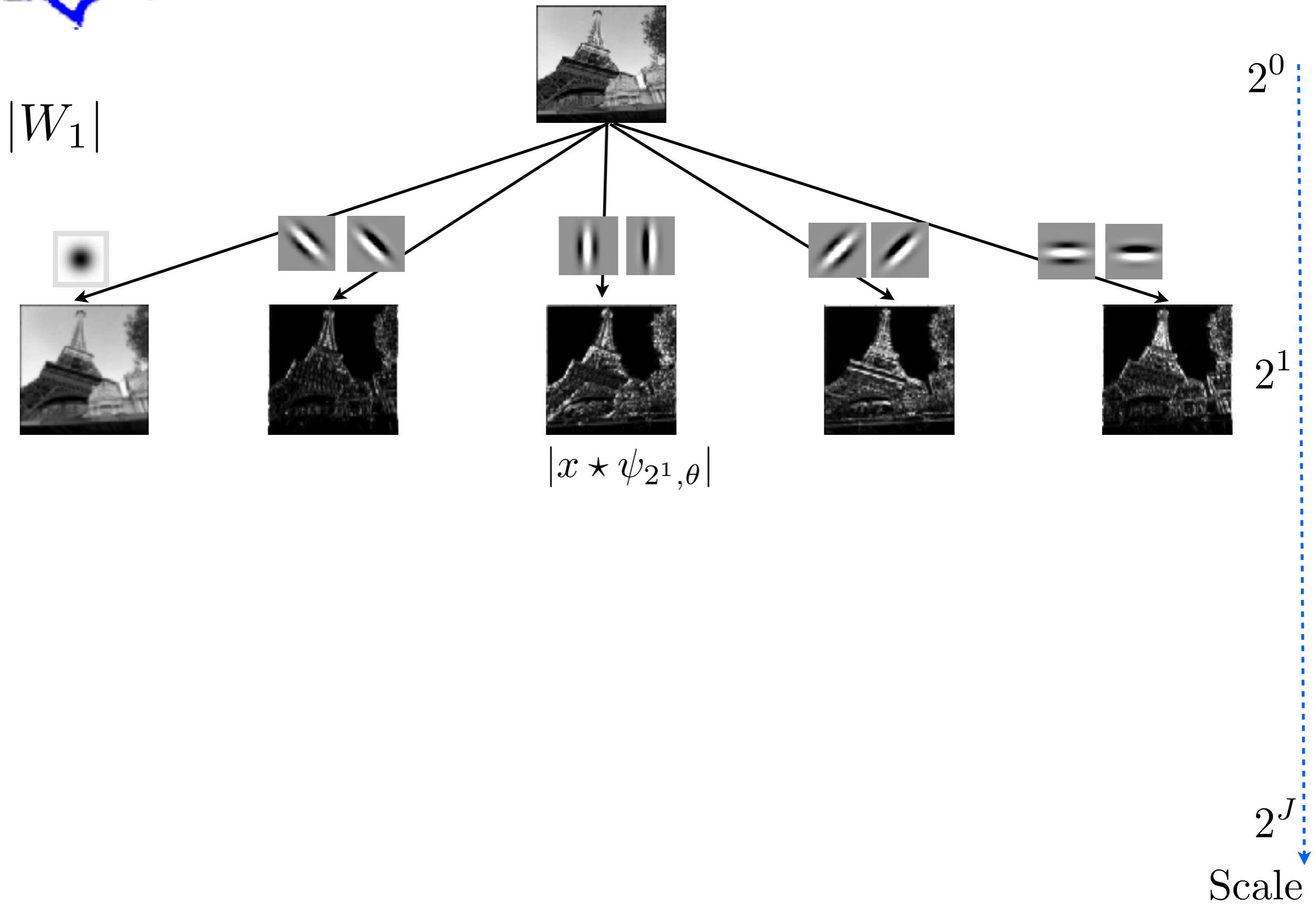


$$x \star \psi_{2^j,\theta}(u) = \int x(v)\,\psi_{2^j,\theta}(u-v)\,dv$$

- Wavelet transform: $Wx = \begin{pmatrix} x \star \phi_{2^J}(u) \\ x \star \psi_{2^j,\theta}(u) \end{pmatrix}_{j \le J, \theta}$   : average

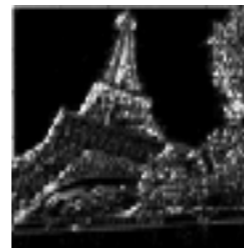  : higher frequencies

Preserves norm: $\|Wx\|^2 = \|x\|^2$ .

$|W_1|$

$|x \star \psi_{2^1,\theta}|$
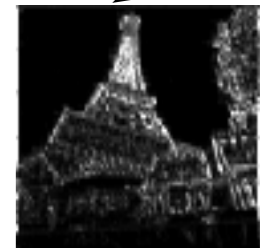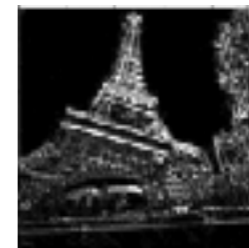
$2^0$

$2^1$

$2^J$

Scale

# Wavelet Filter Bank



$$\rho(\alpha) = |\alpha|$$

$$|W_1|$$

$$x(u)$$

$$2^0$$

$$|x \star \psi_{2^1,\theta}|$$

$$2^1$$

$$|x \star \psi_{2^2,\theta}|$$

$$2^2$$

$$|x \star \psi_{2^j,\theta}|$$

$$2^J$$

Scale

# Wavelet Scattering Network



$$S_J = \rho\, W_1 \quad \rho\, W_2 \quad \cdots \quad \rho\, W_J$$

$$\rho(\alpha) = |\alpha| \qquad S_J x = \left\{ \big|\,\|\,|x \star \psi_{\lambda_1}| \star \psi_{\lambda_2} \star ...\,| \star \psi_{\lambda_m}\big| \star \phi_J \right\}_{\lambda_k}$$

Interactions across scales

$$S_J x = \begin{pmatrix} x \star \phi_{2^J} \\ |x \star \psi_{\lambda_1}| \star \phi_{2^J} \\ ||x \star \psi_{\lambda_1}| \star \psi_{\lambda_2}| \star \phi_{2^J} \\ |||x \star \psi_{\lambda_2}| \star \psi_{\lambda_2}| \star \psi_{\lambda_3}| \star \phi_{2^J} \\ ... \end{pmatrix}_{\lambda_1, \lambda_2, \lambda_3, ...} = \ldots |W_3| \, |W_2| \, |W_1| \, x$$

**Lemma**: $\| |W_k, D_\tau| \| = \| |W_k| D_\tau x' D_\tau W_k| \| \leq C' \|\nabla\tau\|_\infty$

**Theorem**: *For appropriate wavelets, a scattering is*

*contractive* $\|S_J x - S_J y\| \leq \|x - y\|$ $(\mathbf{L^2}$ *stability)*

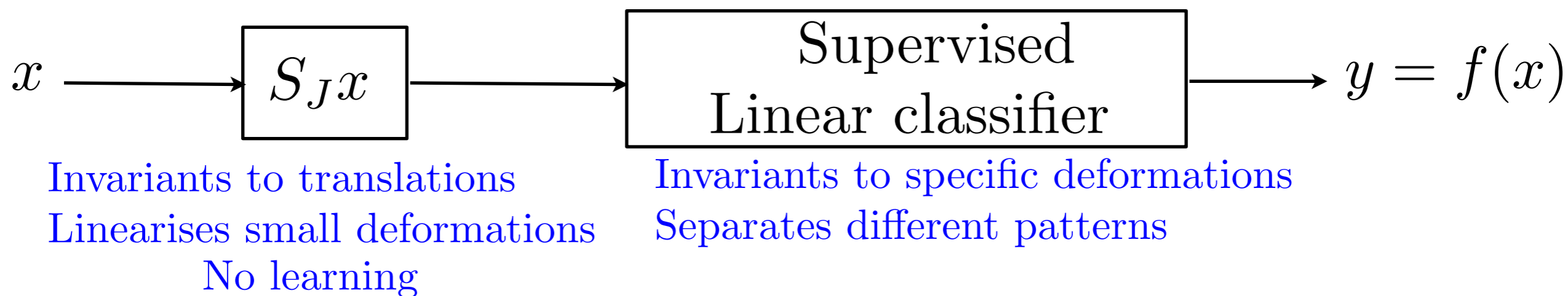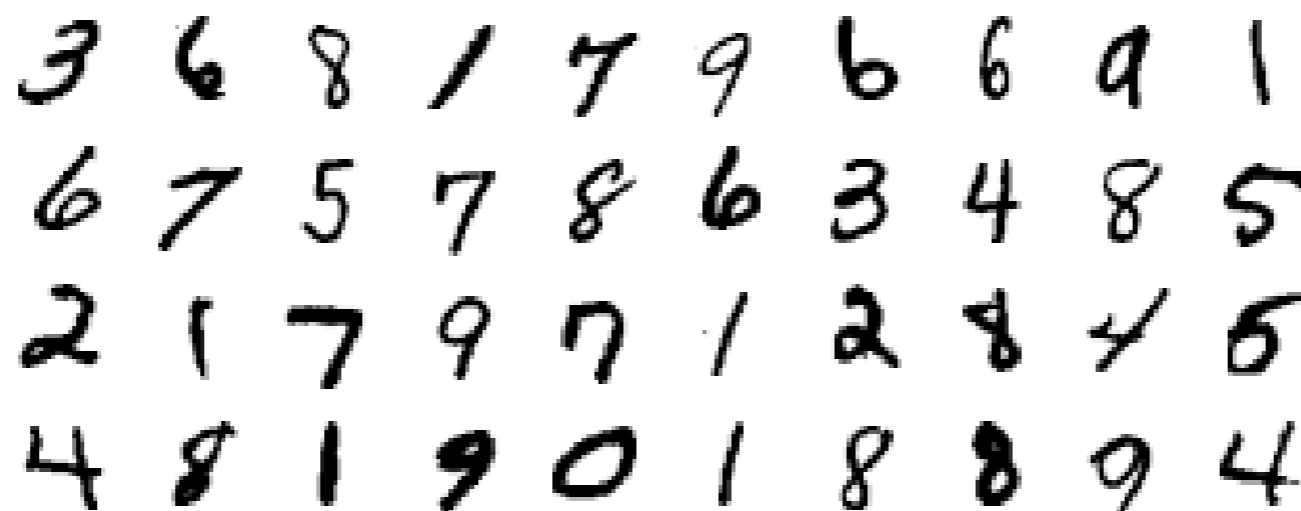*preserves norms* $\|S_J x\| = \|x\|$

*translations invariance and deformation stability:*

*if* $D_\tau x(u) = x(u - \tau(u))$ *then*

$$\lim_{J \to \infty} \|S_J D_\tau x - S_J x\| \leq C \|\nabla\tau\|_\infty \|x\|$$

# Digit Classification: MNIST

*Joan Bruna*



$$x \longrightarrow \boxed{S_J x} \longrightarrow \boxed{\begin{array}{c} \text{Supervised} \\ \text{Linear classifier} \end{array}} \longrightarrow y = f(x)$$

Invariants to translations
Linearises small deformations
No learning

Invariants to specific deformations
Separates different patterns

## Classification Errors

| Training size | Conv. Net. | Scattering |
|---|---|---|
| 50000 | 0.4% | 0.4% |

LeCun et. al.

*joint work with Joan Bruna*

## Unsupervised learning:

Approximate the probability distribution $p(x)$ of $X \in \mathbb{R}^d$ given $P$ realisations $\{x_i\}_{i \leq P}$ with potentially $P = 1$

*Which class of processes can we approximate ?*

- Ergodic versus non-ergodic (long-range dependance)

- Capture non-Gaussianity: geometry of realisations

Scattering/Deep Net. of a stationary process $X(t)$

$$S_J X = \begin{pmatrix} X \star \phi_{2^J}(t) \\ |X \star \psi_{\lambda_1}| \star \phi_{2^J}(t) \\ ||X \star \psi_{\lambda_1}| \star \psi_{\lambda_2}| \star \phi_{2^J}(t) \\ |||X \star \psi_{\lambda_2}| \star \psi_{\lambda_2}| \star \psi_{\lambda_3}| \star \phi_{2^J}(t) \\ ... \end{pmatrix}_{\lambda_1, \lambda_2, \lambda_3, ...} : \text{ stationary vector}$$

Scattering transform of a stationary vector $X \in \mathbb{R}^d$

maximum scale: $2^J = d$

$$S_J X = \begin{pmatrix} d^{-1} \sum_{u=1}^{d} X(u) \\ d^{-1} \| X \star \psi_{\lambda_1} \|_1 \\ d^{-1} \| |X \star \psi_{\lambda_1}| \star \psi_{\lambda_2} \|_1 \\ d^{-1} \| ||X \star \psi_{\lambda_2}| \star \psi_{\lambda_2}| \star \psi_{\lambda_3} \|_1 \\ ... \end{pmatrix}_{\lambda_1, \lambda_2, \lambda_3, ...}$$

$d \to \infty$

Central limit theorem
with "weak" ergodicity conditions

$$\mathbb{E}(SX) = \begin{pmatrix} \mathbb{E}(X) \\ \mathbb{E}(|X \star \psi_{\lambda_1}|) \\ \mathbb{E}(||X \star \psi_{\lambda_1}| \star \psi_{\lambda_2}|) \\ \mathbb{E}(|||X \star \psi_{\lambda_2}| \star \psi_{\lambda_2}| \star \psi_{\lambda_3}|) \\ ... \end{pmatrix}_{\lambda_1, \lambda_2, \lambda_3, ...} : \text{scattering moments.}$$

Scattering transform of a stationary vector $X \in \mathbb{R}^d$

maximum scale: $2^J = d$

$$S_J X = \begin{pmatrix} d^{-1} \sum_{u=1}^{d} X(u) \\ d^{-1} \| X \star \psi_{\lambda_1} \|_1 \\ d^{-1} \| |X \star \psi_{\lambda_1}| \star \psi_{\lambda_2} \|_1 \\ d^{-1} \| ||X \star \psi_{\lambda_2}| \star \psi_{\lambda_2}| \star \psi_{\lambda_3} \|_1 \\ \dots \end{pmatrix}_{\lambda_1, \lambda_2, \lambda_3, \dots}$$

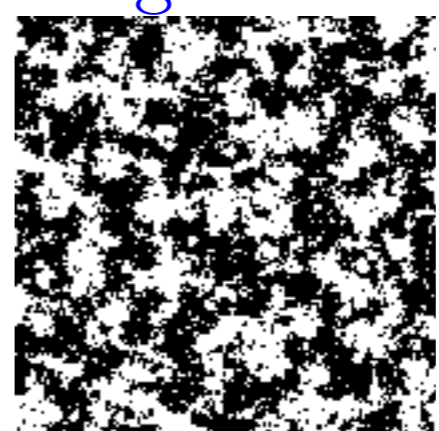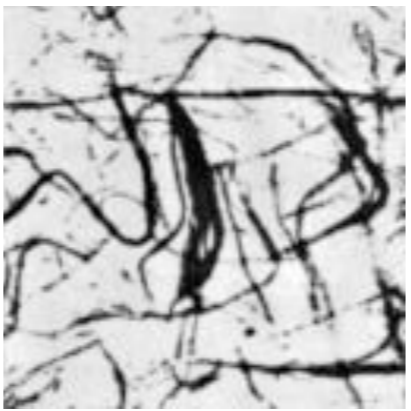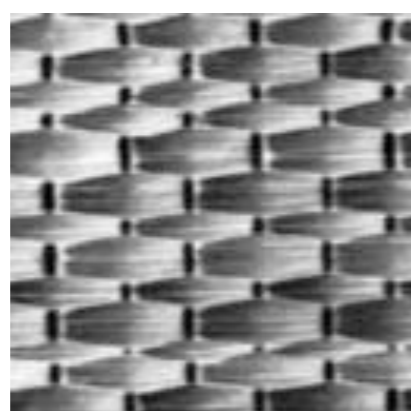- Reconstruction: compute $\tilde{X}$ which satisfies

$$S_J \tilde{X} \approx S_J X$$

with random initialisation and gradient descent.

# Texture Reconstructions
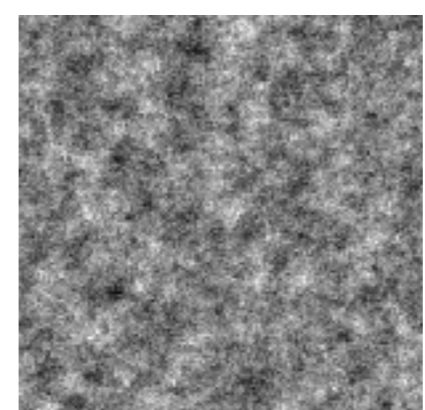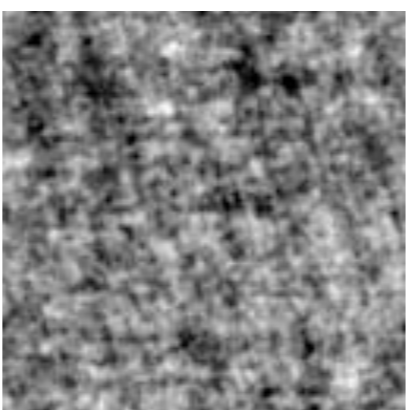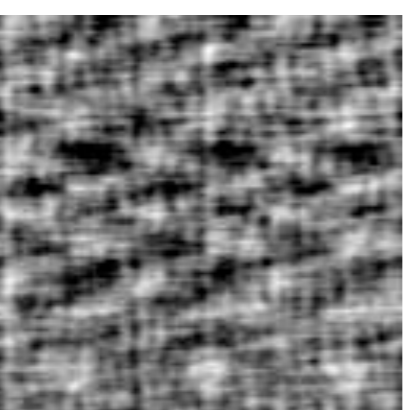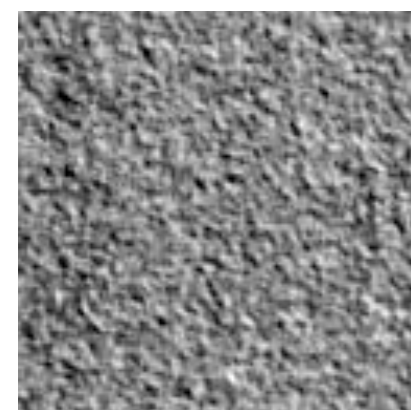
*Joan Bruna*

Texture of $d$ pixels

Ising-critical    Turbulence 2D
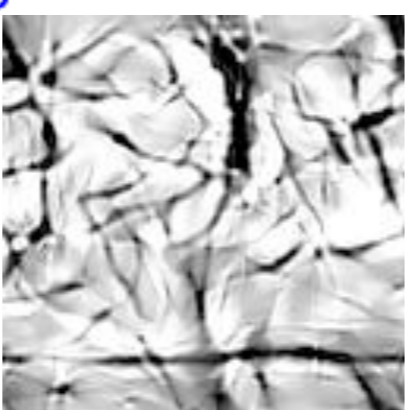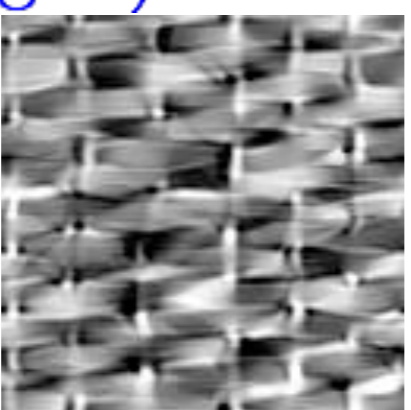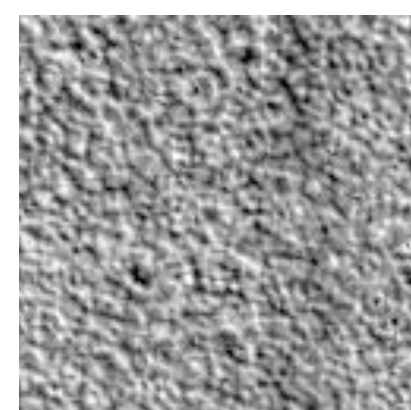


Gaussian process model with $d$ second order moments



Reconstructions from $\|X \star \psi_{\lambda_1}\|_1$ and $\||X \star \psi_{\lambda_1}| \star \psi_{\lambda_2}\|_1$
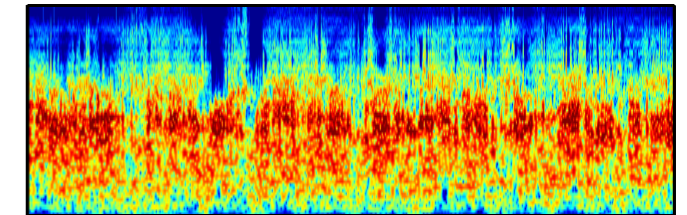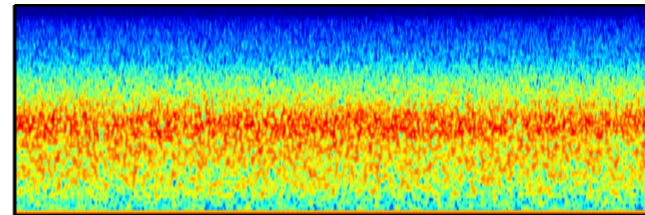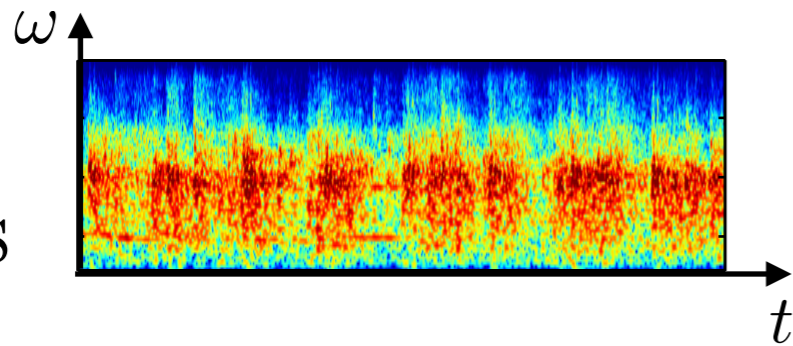
$O(\log^2 d)$ scattering coefficients
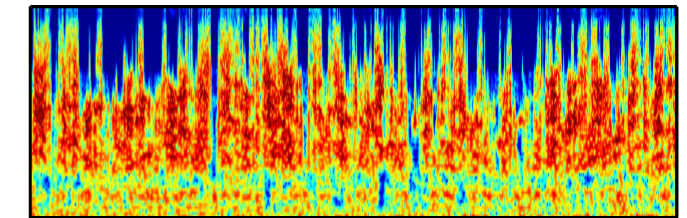
*Joan Bruna*

Original

Gaussian
in time

Scattering
Order 2

Applauds

Paper

Cocktail Party

- A representation $\Phi(x) = \{\phi_k(x)\}_{k \leq K}$ with $x \in \mathbb{R}^d$

- Canonical distribution $p(x)$ of $X$ satisfies

$$\mu_k = \mathbb{E}(\phi_k X) = \int \phi_k(x)\, p(x)\, dx$$

with maximum entropy: $H(p) = -\int p(x) \log p(x)\, dx$

$$\Rightarrow \quad p(x) = Z^{-1} \exp\left(-\sum_k \theta_k \phi_k(x)\right)$$

*Gaussian, Markov random field models*

- **Problem:** in other cases we can't compute the $\theta_k$.

- If concentration: $\mathrm{Prob}\Big(|\Phi X - \mu| < \epsilon\Big) \xrightarrow[d \to \infty]{} 1$

$$\text{with } \mu = \mathbb{E}(\Phi X)$$

$$\mathbb{R}^d \xrightarrow{\quad\quad\Phi\quad\quad} \Omega \subset \mathbb{R}^L$$

$\Phi^{-1}$

$\Phi^{-1}(\mu)$

$\mu$

A microcanonical model $\tilde{X}$ has a distribution $\tilde{p}$ of maximum entropy conditioned to $\Phi\tilde{X} = \mu$ which is uniform in $\Phi^{-1}(\mu)$ (if compact)

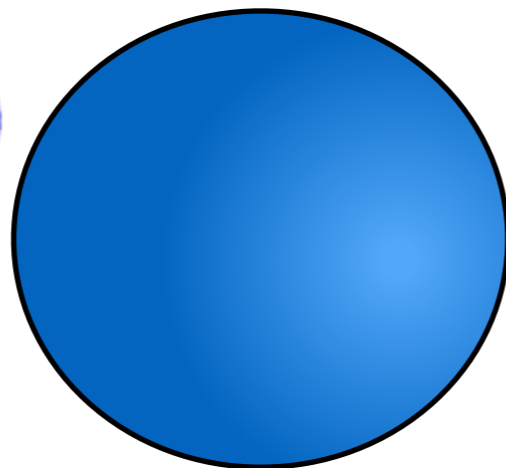- Sphere in $\mathbb{R}^d$ $\quad \Phi x = d^{-1/2}\|x\|_2 = \left(d^{-1}\sum_{k=1}^{d}|x(k)|^2\right)^{1/2} = \mu$

$\Phi^{-1}(\mu)$

not a low-dimensional manifold !



*Borel 1914*
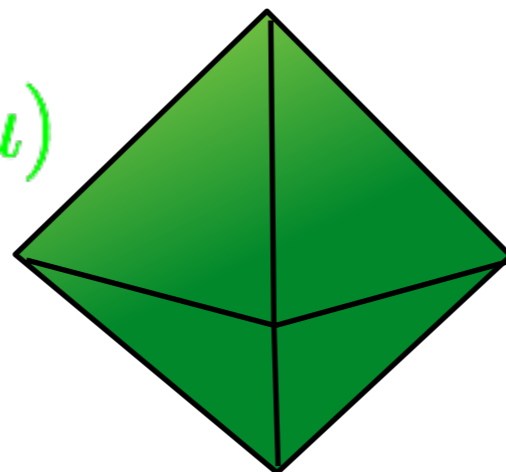*Diaconis, Freedman 1987*

$$\tilde{X}(1),...,\tilde{X}(K) \xrightarrow[d\to\infty]{} \text{.i.i.d Gaussian} \quad \sim e^{-u^2/2\sigma^2}$$

- Simplex in $\mathbb{R}^d$ $\quad \Phi x = d^{-1}\|x\|_1 = d^{-1}\sum_{k=1}^{d}|x(k)| = \mu$

$\Phi^{-1}(\mu)$



*Diaconis, Freedman 1987*

$$\tilde{X}(1),...,\tilde{X}(K) \xrightarrow[d\to\infty]{} \text{i.i.d Exponential} \quad \sim e^{-\lambda|u|}$$

- Scattering coefficients of order 0, 1 and 2; up to scale $2^J$

$$\Phi x = \left\{ d^{-1} \sum_u x(u) \ , \ d^{-1} \|x \star \psi_{\lambda_1}\|_1 \ , \ d^{-1} \||x \star \psi_{\lambda_1}| \star \psi_{\lambda_2}\|_1 \right\}$$

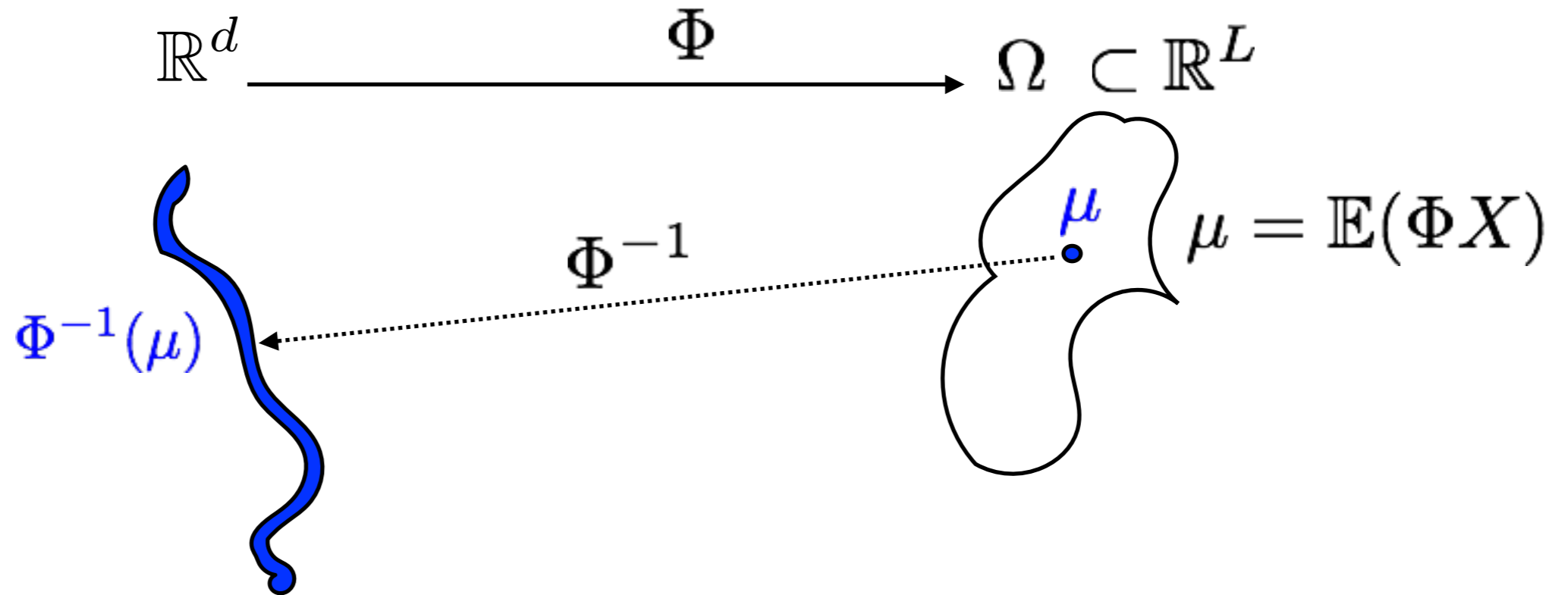$\Phi^{-1}(\mu)$ is an intersection of about $J^2$ polytopes in $\mathbb{R}^d$

Complex high-dimensional geometry

- Reproduces $\mathbf{l^2}$ norms

$$d^{-1}\|x \star \psi_{\lambda_1}\|_2^2 = d^{-2}\|x \star \psi_{\lambda_1}\|_1^2 + \sum_{\lambda_2} d^{-2}\||x \star \psi_{\lambda_1}| \star \psi_{\lambda_2}\|_2^2 + \text{higher order}$$

Specify $\{\|x \star \psi_{\lambda_1}\|_2\}_{\lambda_1}$: intersection of $\mathbf{l^2}$ balls
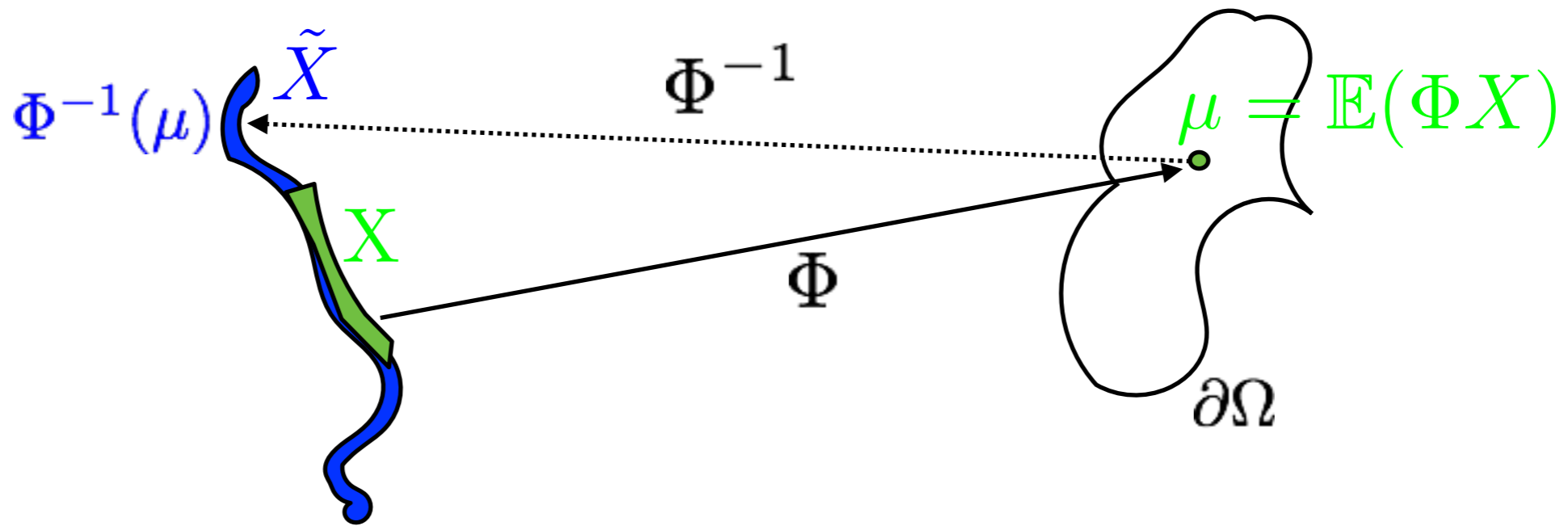
# Microcanonical Scattering



$$\mathbb{R}^d \xrightarrow{\;\;\Phi\;\;} \Omega \subset \mathbb{R}^L$$

$\Phi^{-1}$

$\Phi^{-1}(\mu)$

$\mu = \mathbb{E}(\Phi X)$

**Proposition** If $X(u)$ is stationary and

$X(u)$ and $X(v)$ are independent for $|u - v| \geq \Delta$

then $\quad \lim_{d \to \infty} \mathbb{E}(\|\Phi X - \mu\|^2) = 0$
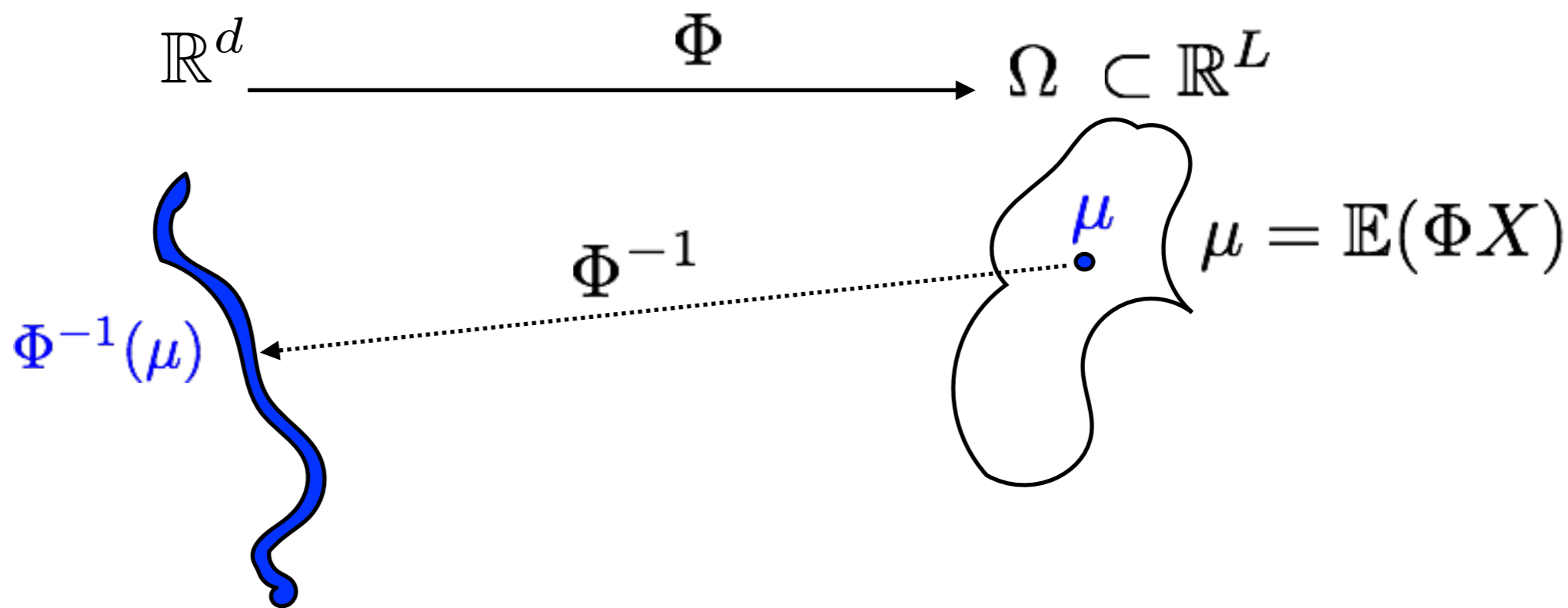
**Theorem** If $X(u)$ is stationary and

$X(u)$ and $X(v)$ are independent for $|u - v| \geq \Delta$

If Typical of $\tilde{X}$ is typical of $X$

and $\quad \lim_{d \to \infty} \mathbb{E}(|d^{-1} \log p(\tilde{X}) - H(p)|^2) = 0$ then

$\tilde{X}(1), ..., \tilde{X}(K)$ converges in probability to $X(1), ..., X(K)$

$$\mathbb{R}^d \xrightarrow{\quad\Phi\quad} \Omega \subset \mathbb{R}^L$$

$$\mu = \mathbb{E}(\Phi X)$$

$$\Phi^{-1}$$

$$\Phi^{-1}(\mu)$$

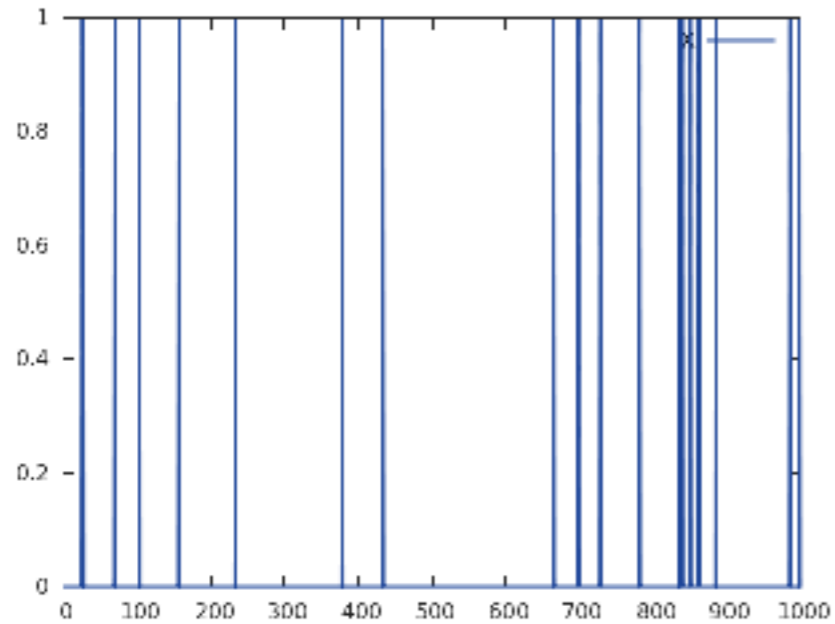If $X$ is Gaussian stationary
with a bounded and regular spectrum
then for a scattering with appropriate wavelets
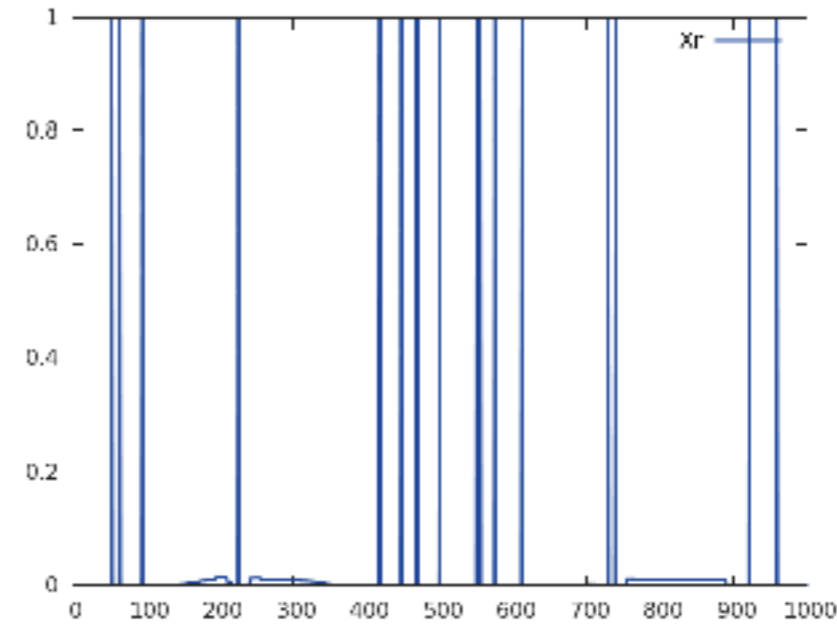$\tilde{X}(1), ..., \tilde{X}(K)$ converges in probability to $X(1), ..., X(K$
up to an arbitrary small error $\epsilon$

Bernoulli $X$

Scattering Microcanonical $\tilde{X}$

Concentration of $\Phi X$  Typical of $\tilde{X}$ is typical of $X$

Why ?

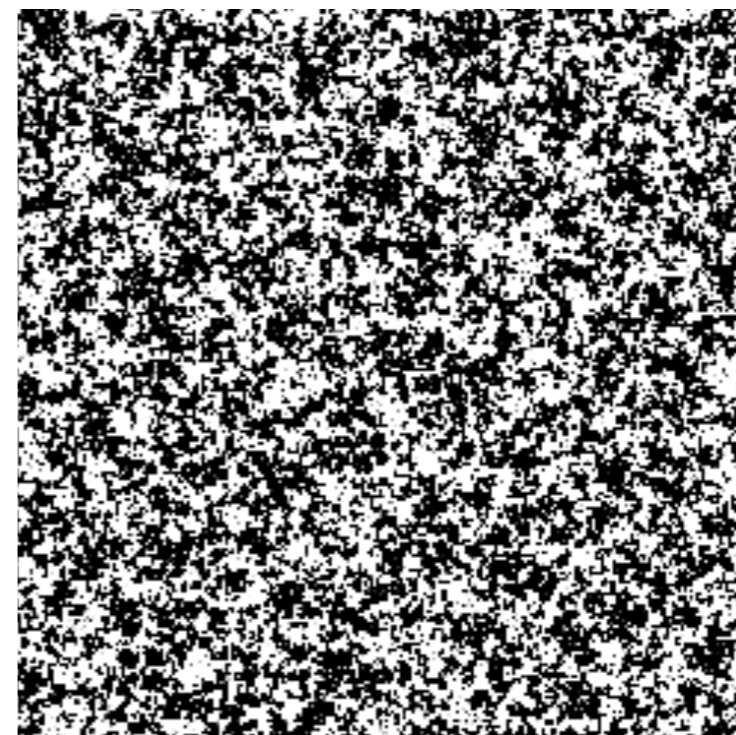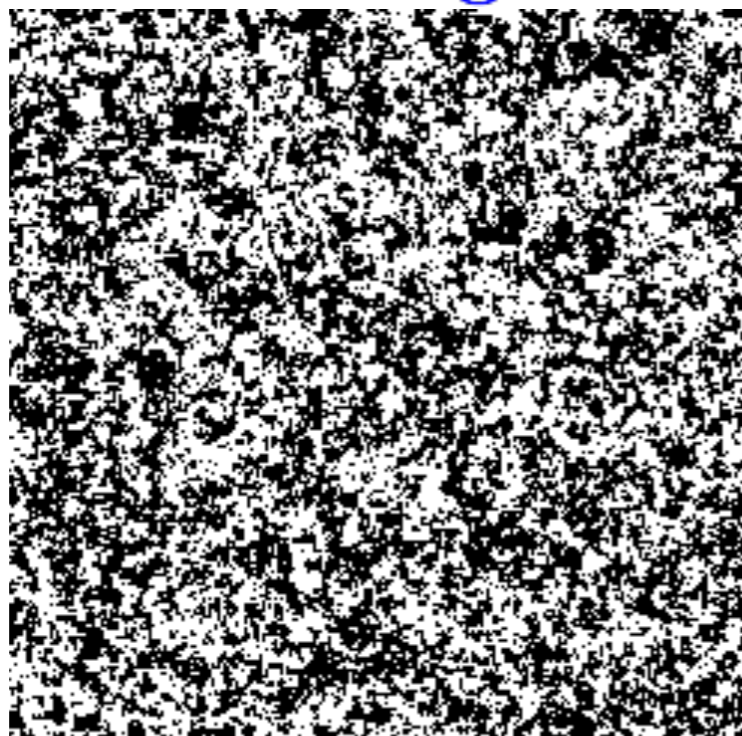$$x(u) \in \{0,1\} \qquad p(x) = Z^{-1} \exp\left(\frac{1}{T} \sum_{(u,u') \in C_I} x(u)\, x(u')\right)$$

Ising $X$ for $T \geq T_{critic}$ 　　　 Microcanonical Scat $\tilde{X}$

Ergodic

Concentration of $\Phi X$ 　　 Typical of $\tilde{X}$ is typical of $X$

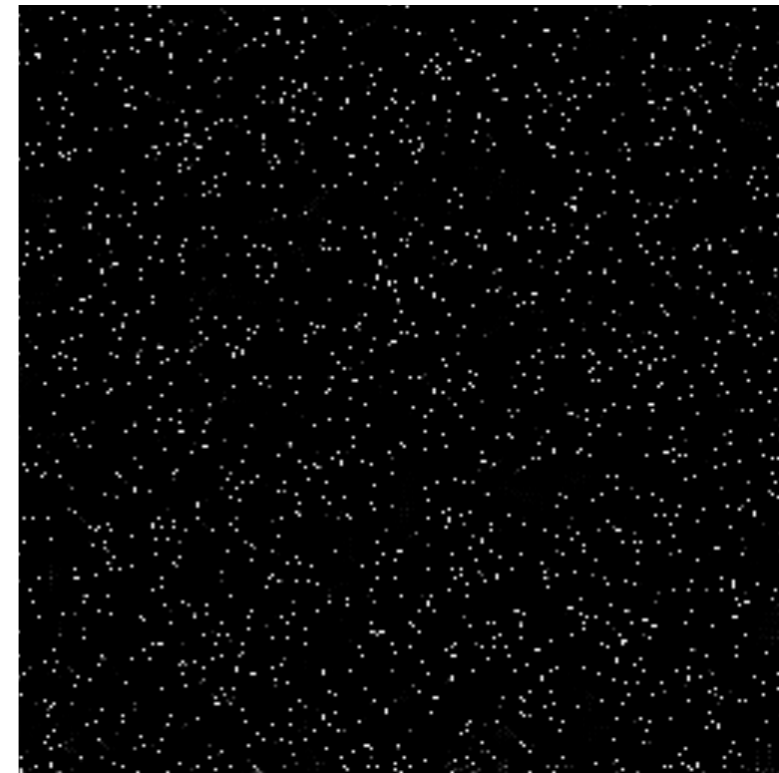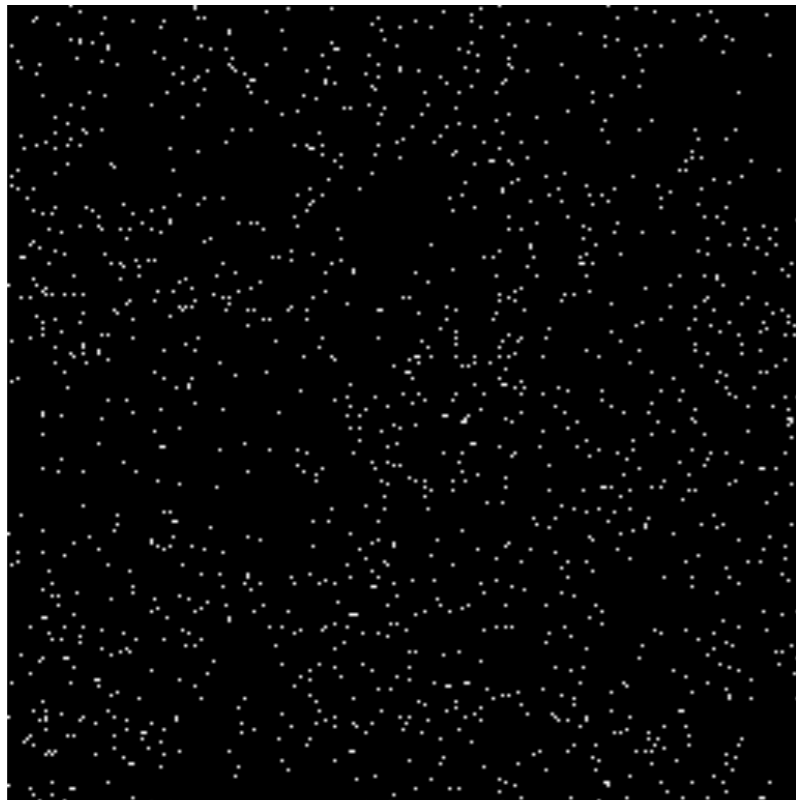| $d$ | $\dfrac{\mathbb{E}(\|\Phi(X) - \mathbb{E}\Phi(X)\|^2)}{\|\mathbb{E}\Phi(X)\|^2}$ | $\dfrac{\mathbb{E}(|d^{-1}\log p(\tilde{X}) - H(p)|^2)}{H(p)^2}$ |
|---|---|---|
| $2^{12}$ | $3 \cdot 10^{-4}$ | $1 \cdot 10^{-5}$ |
| $2^{14}$ | $1 \cdot 10^{-4}$ | $5 \cdot 10^{-6}$ |

Why ?

Bernoulli with random density $\lambda(u)$

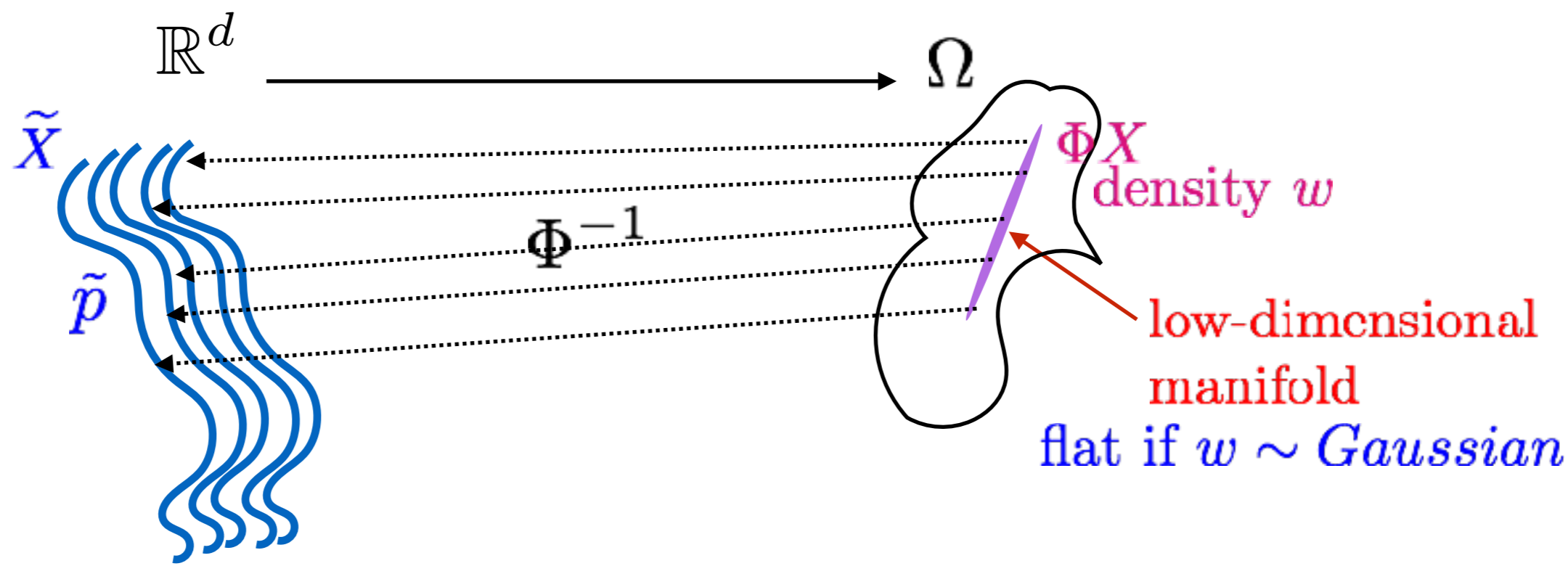Cox $X$ Ergodic

Microcanonical Scat $\tilde{X}$



Concentration of $\Phi X$ — Typical of $\tilde{X}$ is typical of $X$

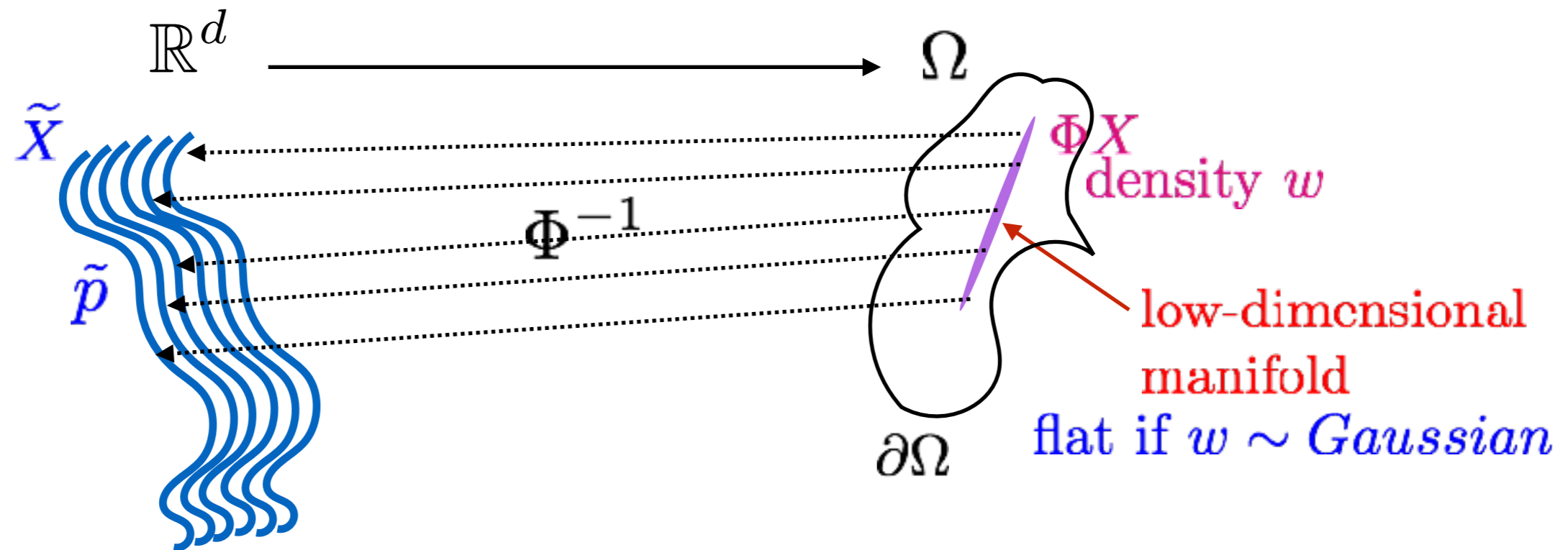| $d$ | $\dfrac{\mathbb{E}(\|\Phi(X)-\mathbb{E}\Phi(X)\|^2)}{\|\mathbb{E}\Phi(X)\|^2}$ | $\dfrac{\mathbb{E}(|d^{-1}\log p(\tilde{X})-H(p)|^2)}{H(p)^2}$ |
|---|---|---|
| $2^{12}$ | $3\cdot 10^{-4}$ | |
| $2^{14}$ | $1\cdot 10^{-4}$ | |

- Non-ergodicity: $\Phi(X)$ does not concentrate in all directions



Maximum entropy conditioned to $\Phi\tilde{X}$ having a density $w$

micro canonical mixture $\tilde{X}$ weighted by the density $w$ of $\Phi X$

- Non-ergodicity: $\Phi(X)$ does not concentrate in all directions

$\mathbb{R}^d \longrightarrow \Omega$

$\widetilde{X}$

$\Phi^{-1}$

$\tilde{p}$

$\Phi X$
density $w$

low-dimensional
manifold

$\partial \Omega$

flat if $w \sim Gaussian$

**Theorem** A microcanonical mixture has a density $\tilde{p}$ with

$$\tilde{p}(x) = \frac{w(\Phi x)}{h(\Phi x)}$$

with  $h(y) = \int_{\Phi^{-1}(y)} |J_L \Phi x|^{-1} \, d\mathcal{H}^{d-L}(x)$

which is singular only if $\Phi x \in \partial \Omega$

- Multifractal processes with stationary increment have non-ergodic low-frequencies: long-range correlations.

- Wavelet coefficients $X \star \psi_\lambda(u)$ decorrelate at larger scales

- Scattering coefficients of order 0, 1 and 2:

$$\Phi X = \left\{ d^{-1} \sum_u X(u), d^{-1} \|X \star \psi_{\lambda_1}\|_1, d^{-1} \||X \star \psi_{\lambda_1}| \star \psi_{\lambda_2}\|_1 \right\}$$

$\qquad\qquad$ non-ergodic $\qquad\qquad\qquad$ ergodic $\qquad\qquad\qquad$ ergodic

$\Rightarrow$ one-dimensional mixture weight $w$ (non-ergodic part) can be estimated from few examples: manifold.

$$p(x) = Z^{-1} \exp\left(\frac{1}{T} \sum_{(u,u')\in C_I} x(u)\, x(u')\right)$$

Ising $X$ for $T = T_{critic}$
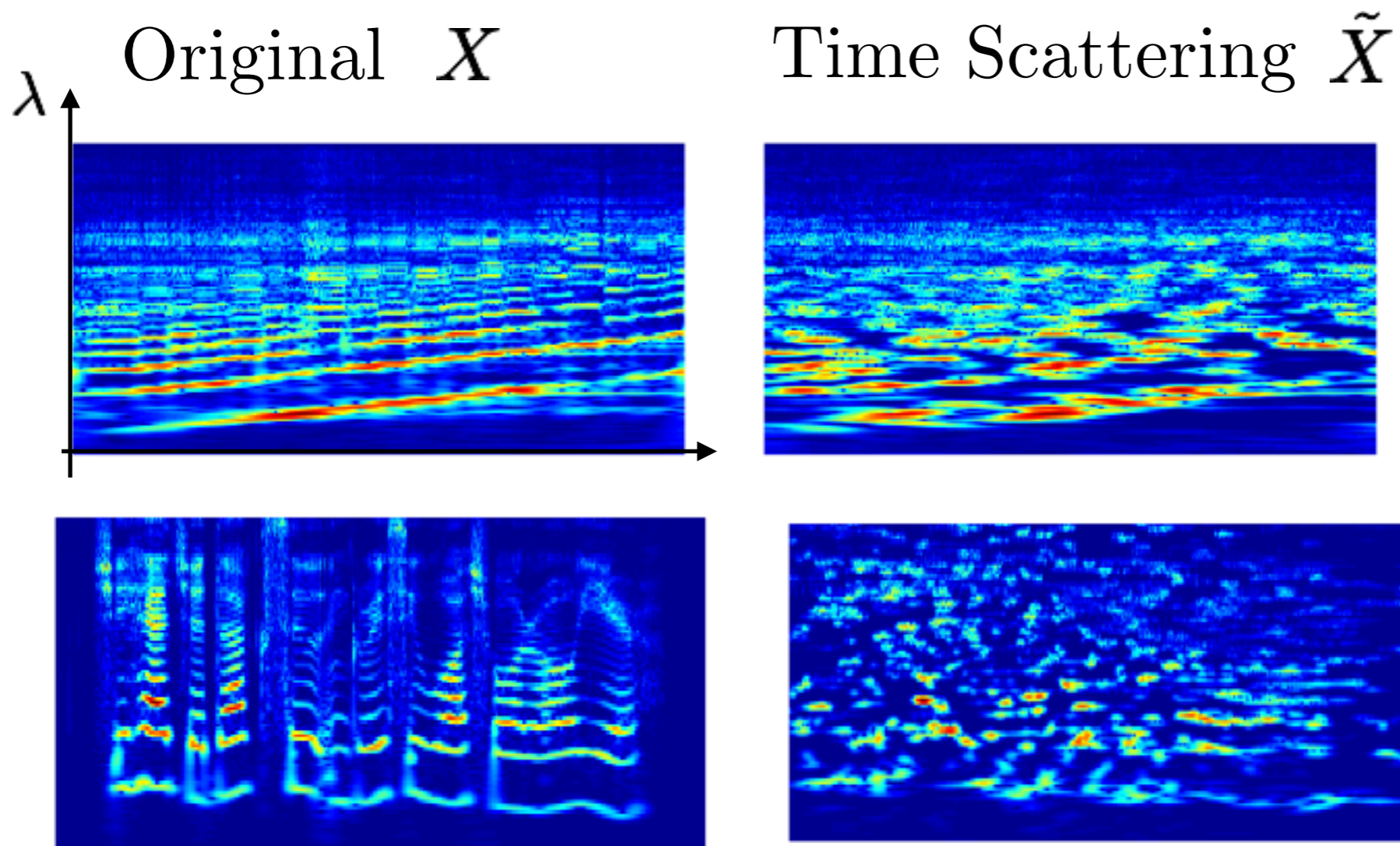
Non ergodic

Microcanonical Scat $\tilde{X}$





Concentration of $\Phi X$ without low-freq.

Typical of $\tilde{X}$ is typical of $X$

| $d$ | $\dfrac{\mathbb{E}(\|\Phi(X)-\mathbb{E}\Phi(X)\|^2)}{\|\mathbb{E}\Phi(X)\|^2}$ | $\dfrac{\mathbb{E}(|d^{-1}\log p(\tilde{X})-H(p)|^2)}{H(p)^2}$ |
|---|---|---|
| $2^{12}$ | $8 \cdot 10^{-3}$ | $2 \cdot 10^{-3}$ |
| $2^{14}$ | $2.5 \cdot 10^{-3}$ | $2 \cdot 10^{-4}$ |

# Failures of Audio Synthesis

*J. Anden and V. Lostanl*

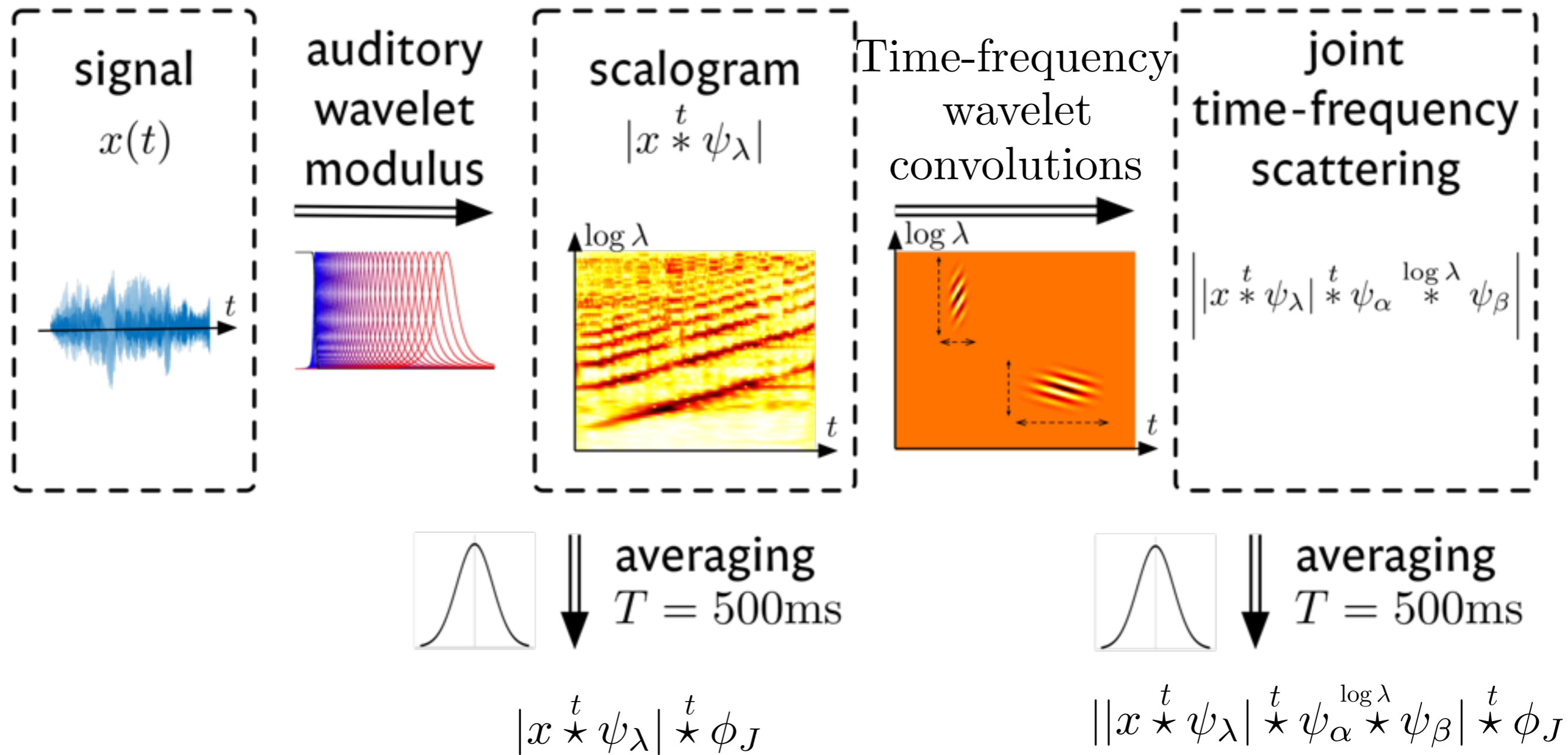Original  $X$　　　Time Scattering  $\tilde{X}$



Typical of $\tilde{X}$ is not typical of $X$

- Missing frequency connections $\Rightarrow$ misalignments

$\Rightarrow$ incorporate two-dimensional translations in time-frequency

*J. Anden and V. Lostanlen*

# Joint Time-Frequency Scattering
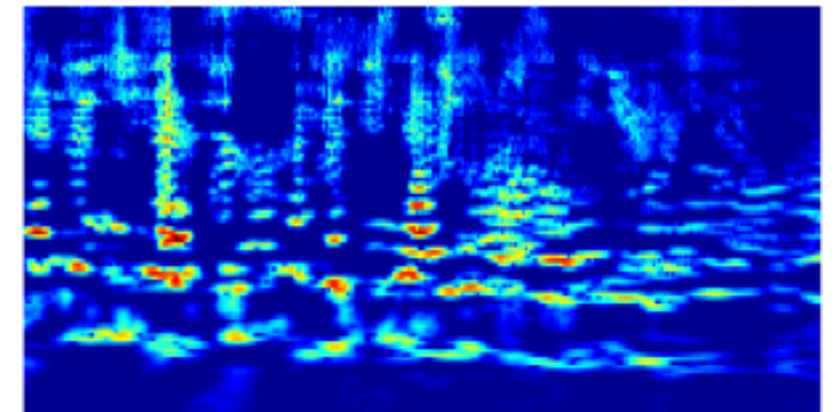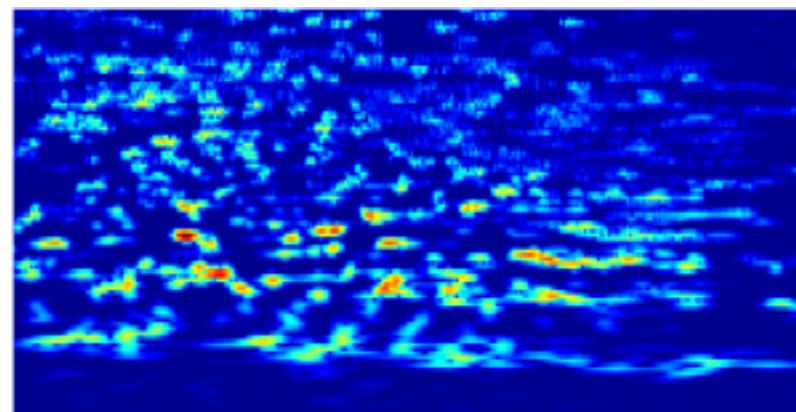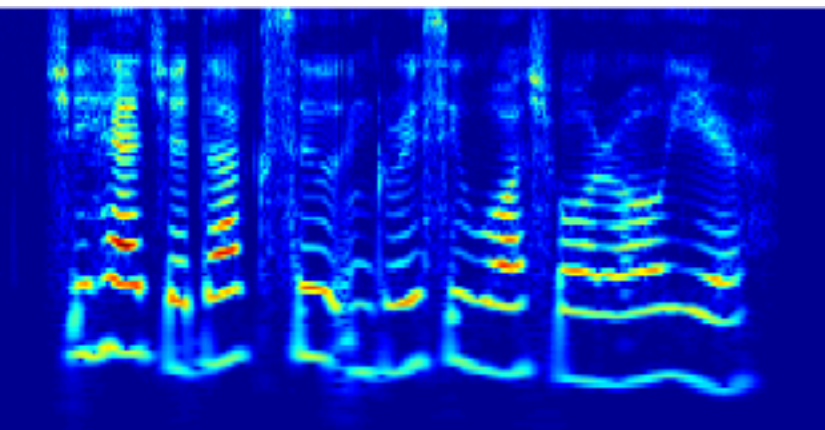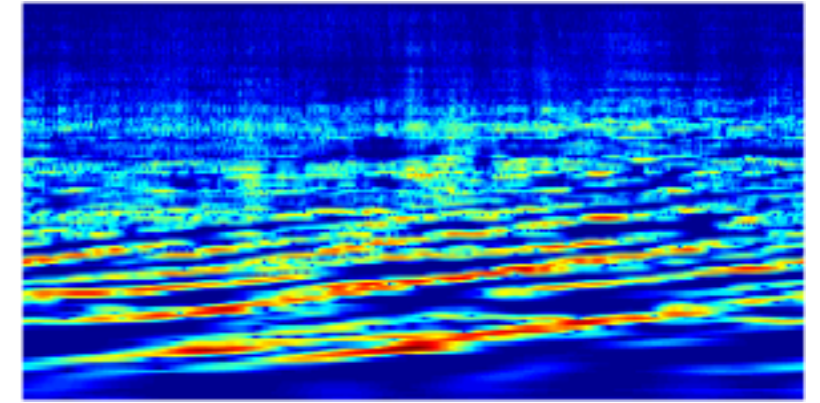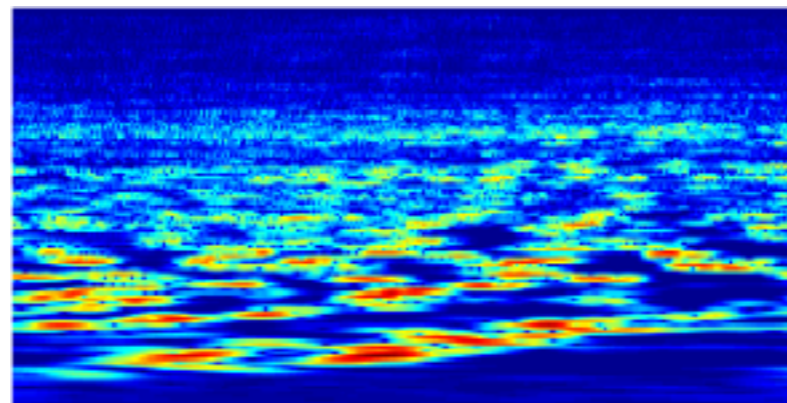
*J. Anden and V. Lostanl*

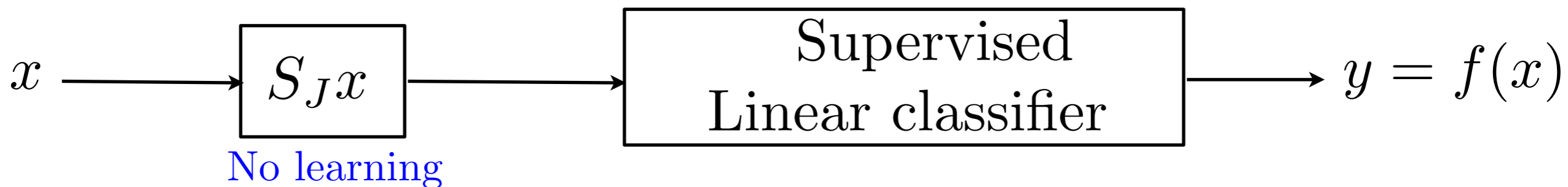Original　　　　Time Scattering　　　Time/Freq Scattering

$$x_j = \rho\, L_j\, x_{j-1}$$

- $L_j$ is a linear combination of convolutions and subsampling:

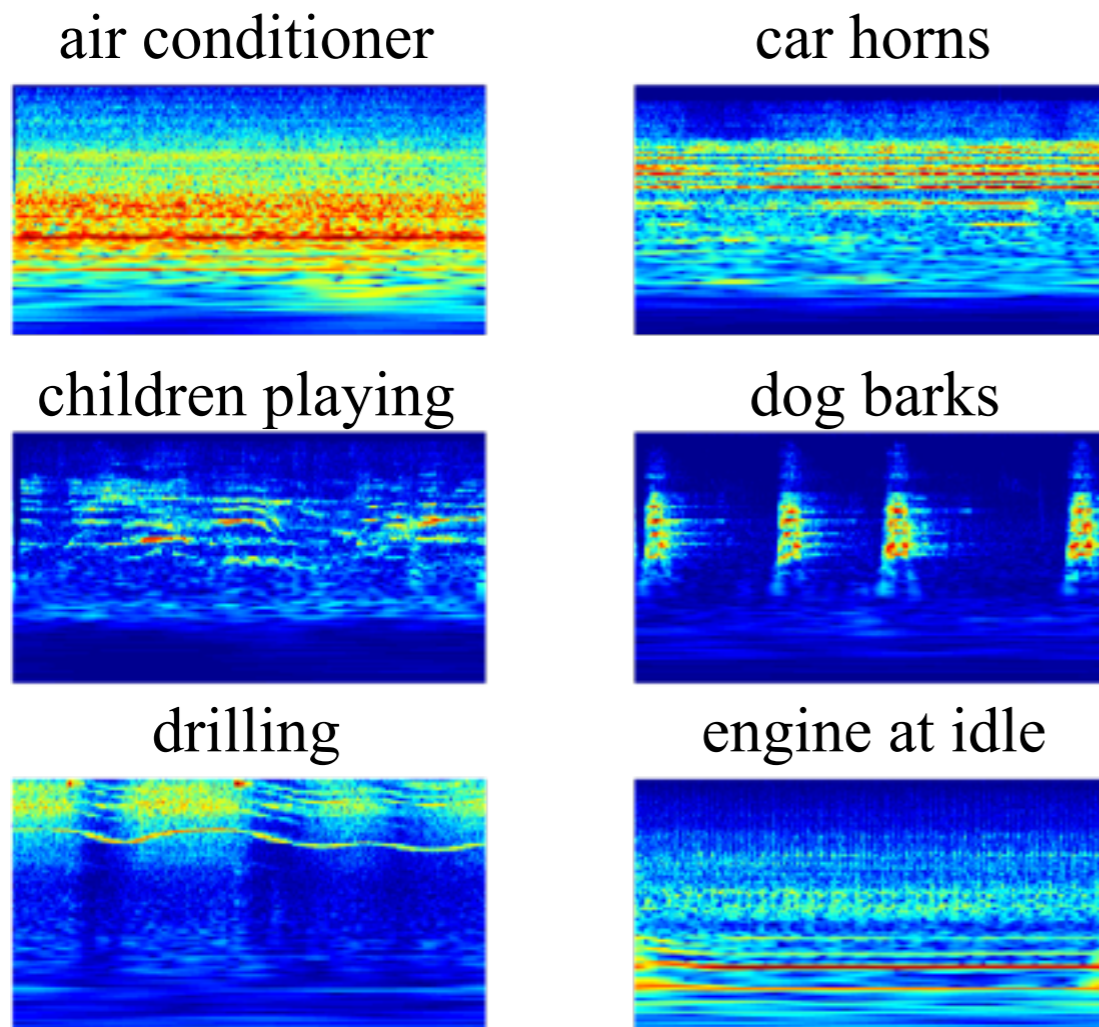$$x_j(u, k_j) = \rho\Big( \sum_k x_{j-1}(\cdot, k) \star h_{k_j,k}(u) \Big)$$

sum across channels

What is the role of channel connections ?

Invariant over groups of operators other than translations

# Environmental Sound Classification

*J. Anden and V. Lostanlen*

$$x \longrightarrow \boxed{S_J x} \longrightarrow \boxed{\begin{array}{c} \text{Supervised} \\ \text{Linear classifier} \end{array}} \longrightarrow y = f(x)$$

No learning

UrbanSound8k: 10 classes
8k training examples

class-wise average error

air conditioner

car horns

children playing

dog barks

drilling

engine at idle

| | |
|---|---|
| MFCC audio descriptors | 0,39 |
| time scattering | 0,27 |
| ConvNet (Piczak, MLSP 2015) | 0,26 |
| time-frequency scattering | 0,2 |

*Joan Bruna*

- Given $S_J x$ we want to compute $\tilde{x}$ such that:

$$S_J \tilde{x} = \begin{pmatrix} \tilde{x} \star \phi_{2^J} \\ |\tilde{x} \star \psi_{\lambda_1}| \star \phi_{2^J} \\ ... \\ |||\tilde{x} \star \psi_{\lambda_1}| \star ..| \star \psi_{\lambda_m}| \star \phi_{2^J} \end{pmatrix}_{\lambda_1,...,\lambda_m} = \begin{pmatrix} x \star \phi_{2^J} \\ |x \star \psi_{\lambda_1}| \star \phi_{2^J} \\ ... \\ |||x \star \psi_{\lambda_1}| \star ..| \star \psi_{\lambda_m}| \star \phi_{2^J} \end{pmatrix}_{\lambda_1,...,\lambda_m} = S_J x$$

We shall use $m = 2$.

- If $x(u)$ is a Dirac, or a straight edge or a sinusoid then $\tilde{x}$ is equal to $x$ up to a translation.
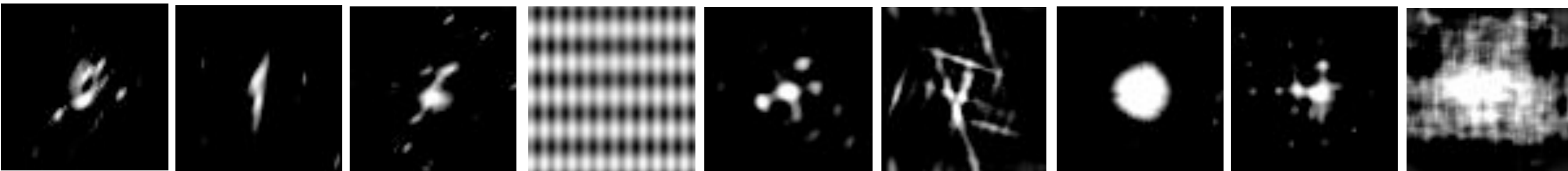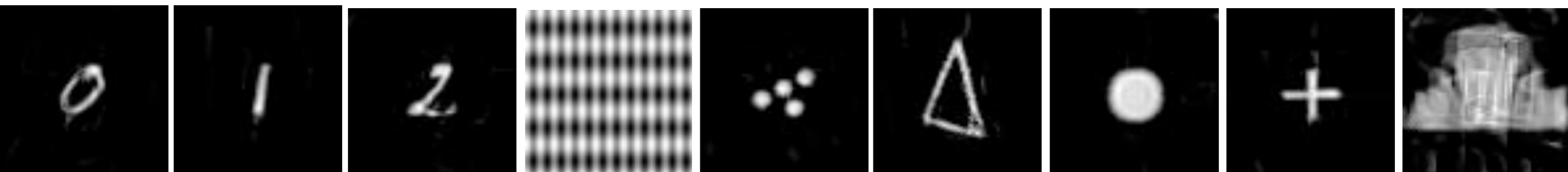
*Joan Bruna*

With a gradient descent algorithm:

Original images of $N^2$ pixels:



$m = 1, 2^J = N$: reconstruction from $O(\log_2 N)$ scattering coeff.



$m = 2, 2^J = N$: reconstruction from $O(\log_2^2 N)$ scattering coeff.

# Multiscale Scattering Reconstructions

Original
Images
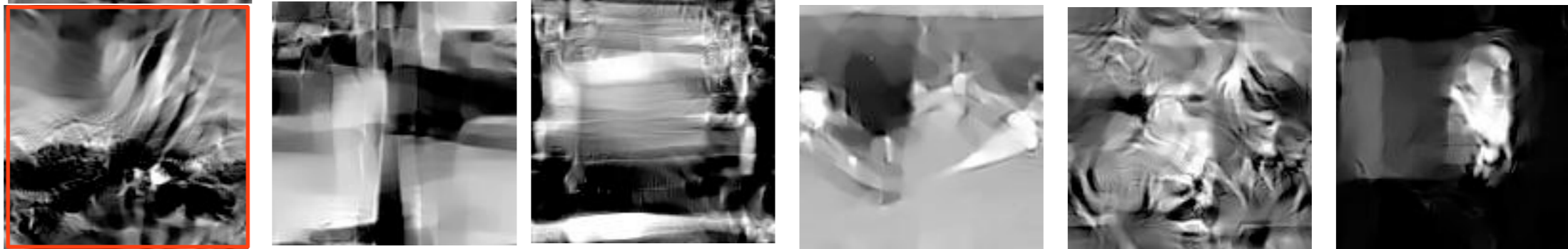$N^2$ pixels

Scattering
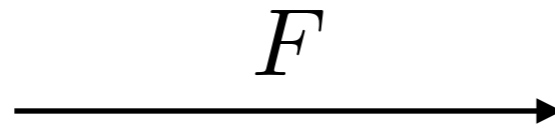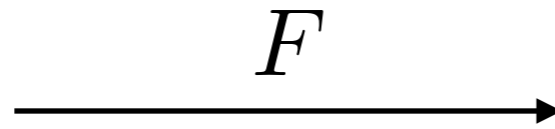Reconstruction
$2^J = 16$
$1.4\,N^2$ coeff.

$2^J = 32$
$0.5\,N^2$ coeff.

$2^J = 64$

$2^J = 128 = N$

$x$     $F$     $y$

- Best Linear Method: Least Squares estimate (linear interpolation):

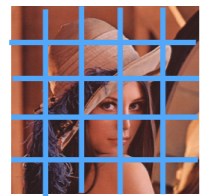$$\hat{y} = (\widehat{\Sigma}_x^{\dagger} \widehat{\Sigma}_{xy}) x$$

# Super-Resolution



$$x \xrightarrow{\quad F \quad} y$$

- Best Linear Method: Least Squares estimate (linear interpolation):
- State-of-the-art Methods:

$$\hat{y} = (\widehat{\Sigma}_x^\dagger \widehat{\Sigma}_{xy})x$$

  – Dictionary-learning Super-Resolution

  – CNN-based: Just train a CNN to regress from low-res to high-res.

  – They optimize cleverly a fundamentally unstable metric criterion:

$$\Theta^* = \arg\min_\Theta \sum_i \|F(x_i, \Theta) - y_i\|^2 \quad , \quad \hat{y} = F(x, \Theta^*)$$

$$S_{L,J}\, x \xrightarrow{\quad F \quad} S_{L-\alpha,J}\, x$$
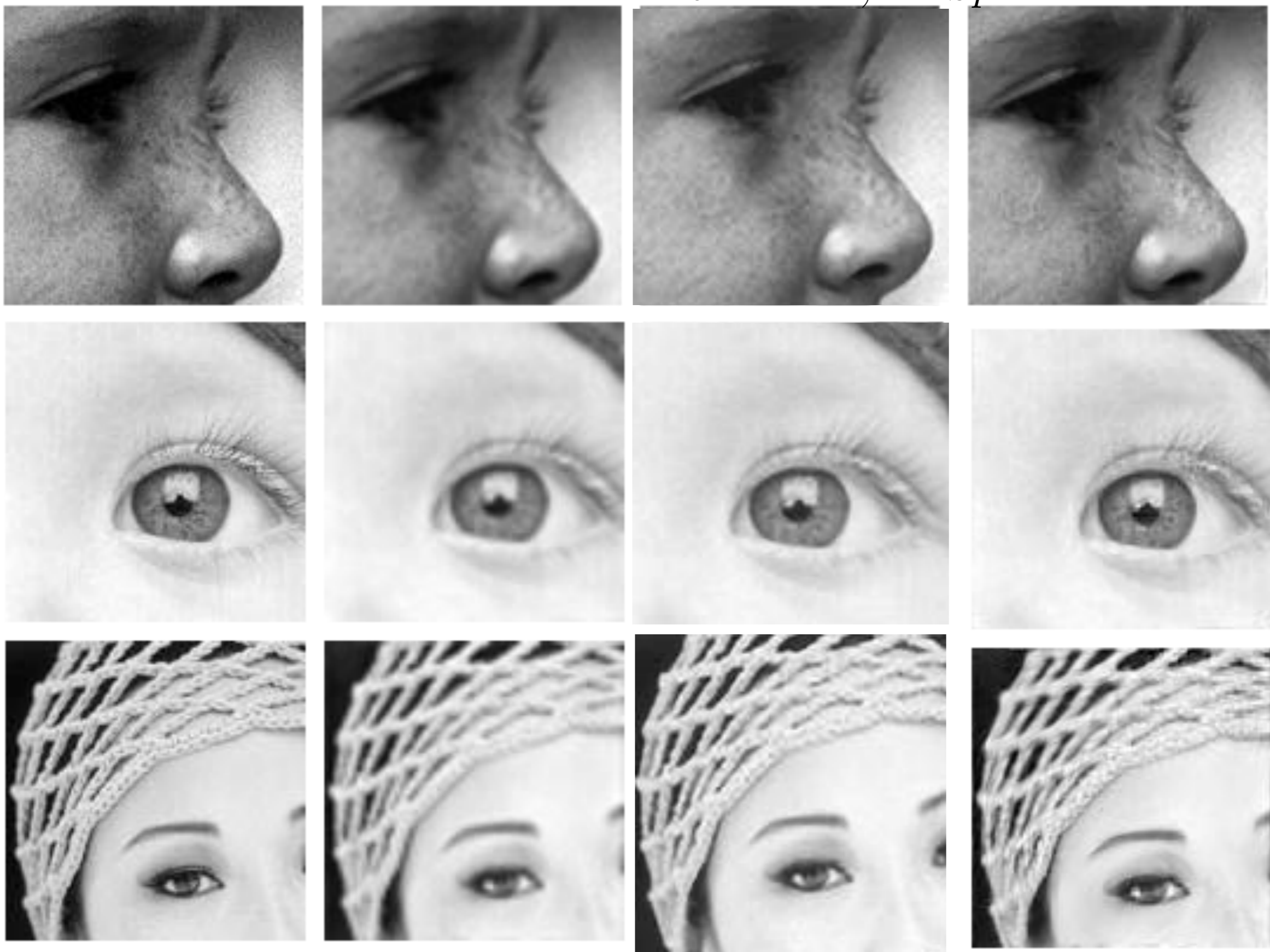
$$S_{L,J}x = \begin{pmatrix} x \star \phi_{2^J}(u) \\ |x \star \psi_{j_1,k_1}| \star \phi_{2^J}(u) \\ ||x \star \psi_{j_1,k_1}| \star \psi_{j_2,k_2}| \star \phi_{2^J}(u) \end{pmatrix}_{L \le j_1, j_2 \le J}$$

- Linear estimation in the scattering domain

- No phase estimation: potentially worst PSNR

- Good image quality because of deformation stability

# Super-Resolution Results

*J. Bruna, P. Sprechmann*



Original          Linear Estimate          state-of-the-art          Scattering

# Super-Resolution Results

*J. Bruna, P. Sprechmann*



| Original | Best Linear Estimate | state-of-the-art | Scattering Estimate |

# Super-Resolution Results

*J. Bruna, P. Sprechmann*



Original

Best
Linear Estimate

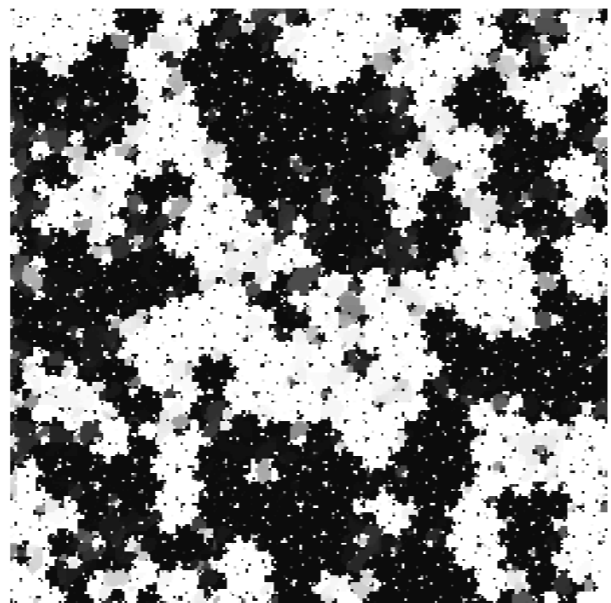state-of-the-art

Scattering
Estimate

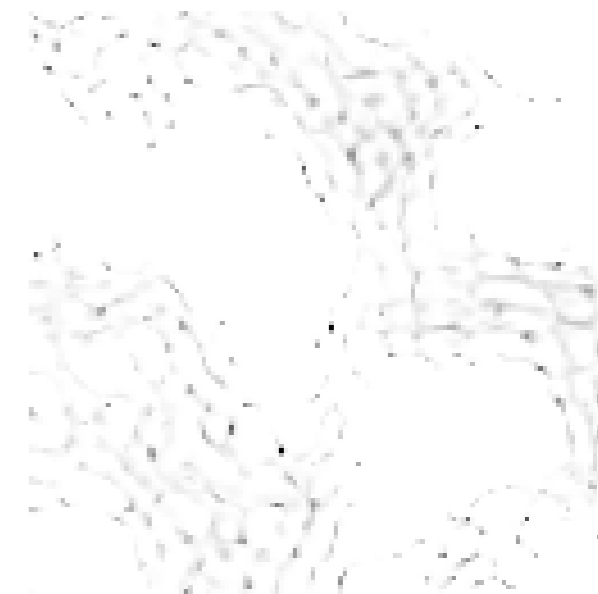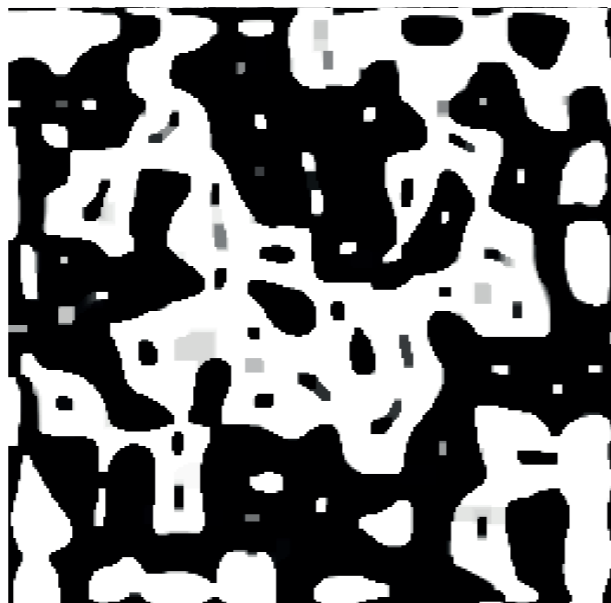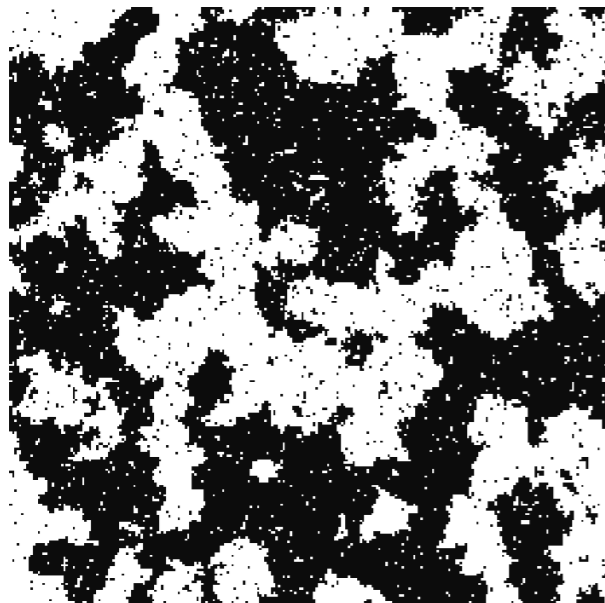# Super-Resolution Results

*I. Dokmanic, J. Bruna, M. De Hoop*



Original     TV Regularization     Original     $l^1$ Regularization
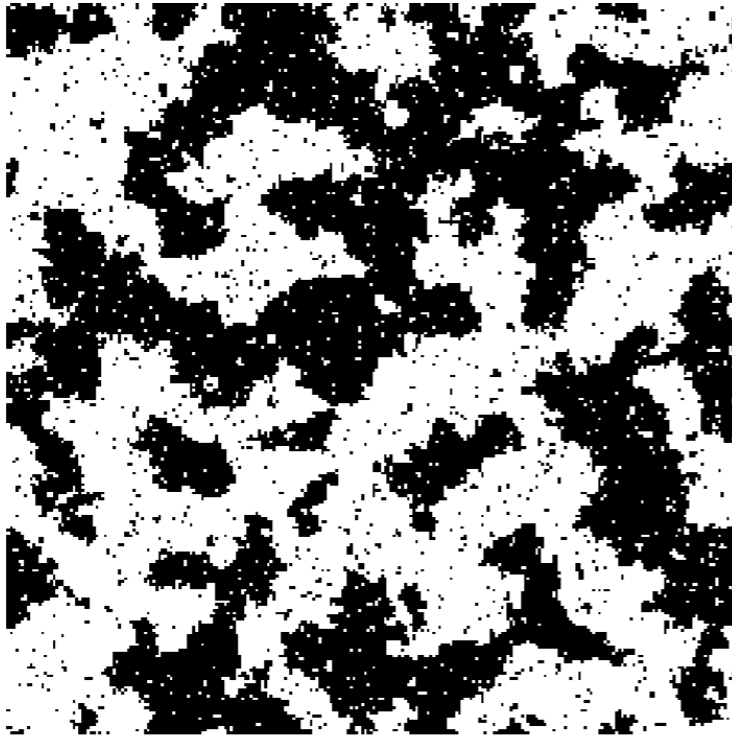
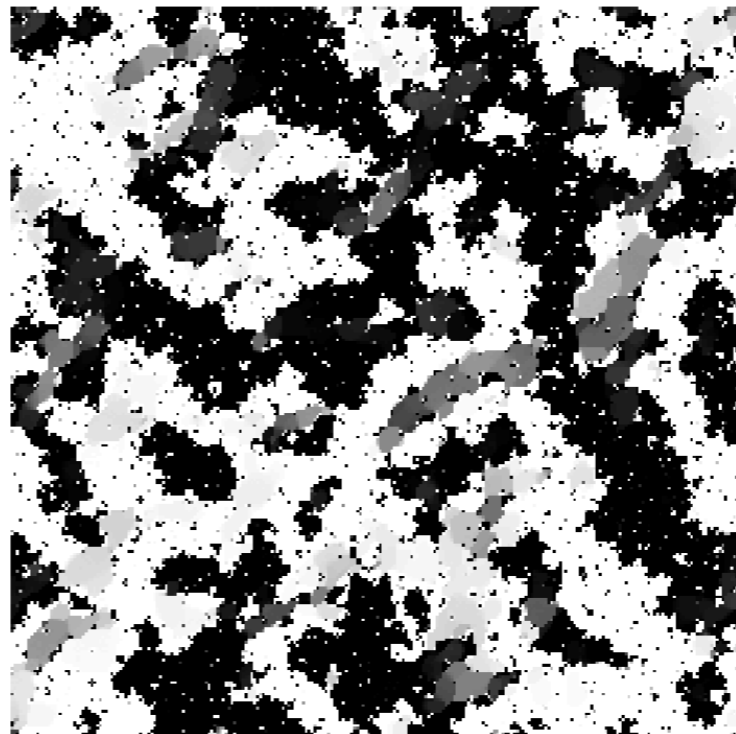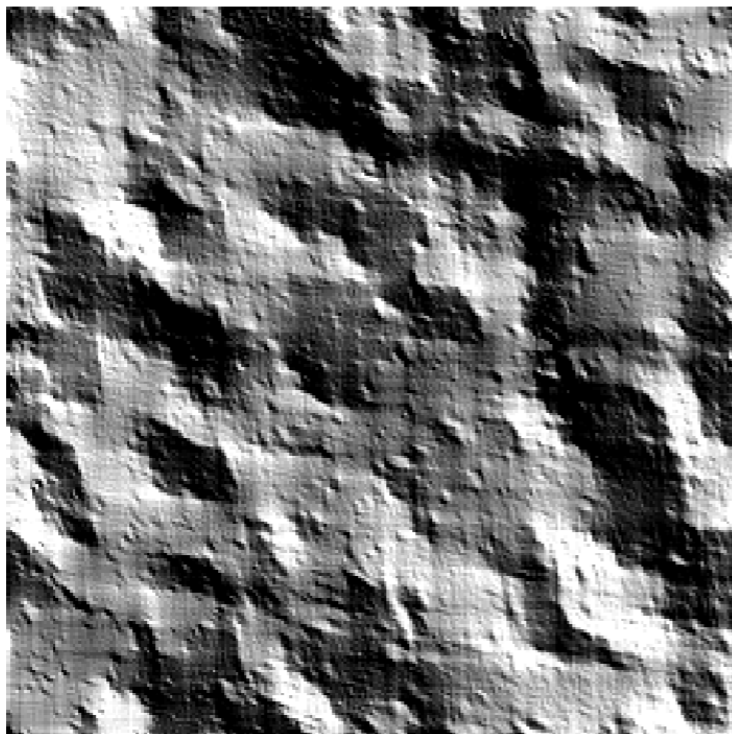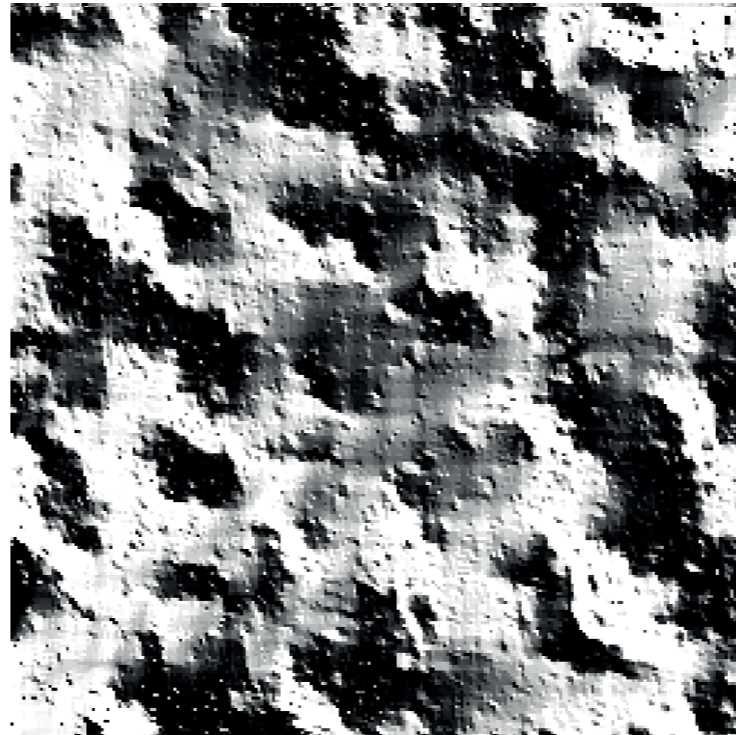Low-Resolution     Scattering     Low-Resolution     Scattering

*I. Dokmanic, J. Bruna, M. De Hoop*



Original

TV Regularization

Low-Resolution

Scattering

# Conclusions

- Deep convolutional networks have spectacular high-dimensional and generic approximation capabilities.

- New stochastic models of images for inverse problems.

- Outstanding mathematical problem to understand deep nets:
  - How to learn representations for inverse problems ?

*(Not) Understanding Deep Convolutional Networks*, arXiv 2016.