



Learning Music, Images and Physics with Deep Neural Networks



*Joan Bruna, Matthew Hirn, Stéphane Mallat
Vincent Lostanlen, Edouard Oyallon, Nicolas Poilvert,
Laurent Sifre, Irène Waldspurger*

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

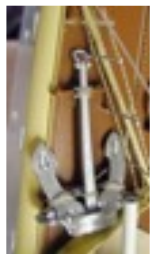
High Dimensional Learning

- High-dimensional $x = (x(1), \dots, x(d)) \in \mathbb{R}^d$:
- **Classification:** estimate a class label $f(x)$
given n sample values $\{x_i, y_i = f(x_i)\}_{i \leq n}$

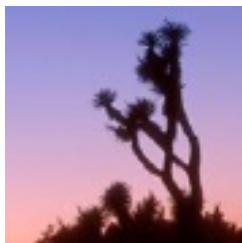
Image Classification $d = 10^6$

Huge variability
inside classes

Anchor



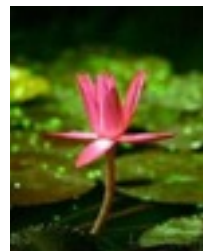
Joshua Tree



Beaver



Lotus



Water Lily

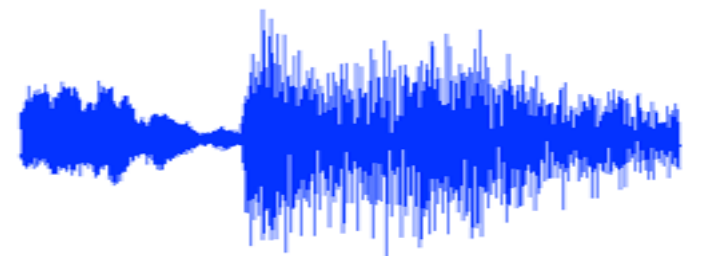
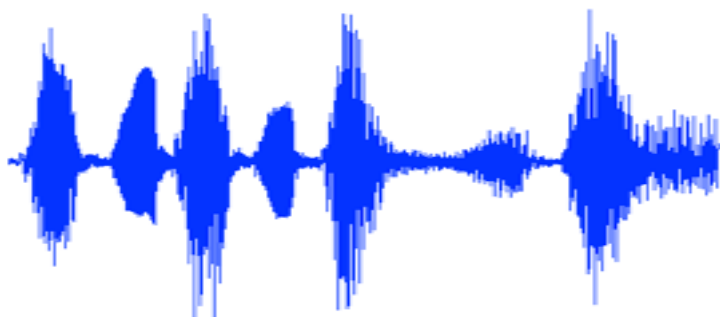
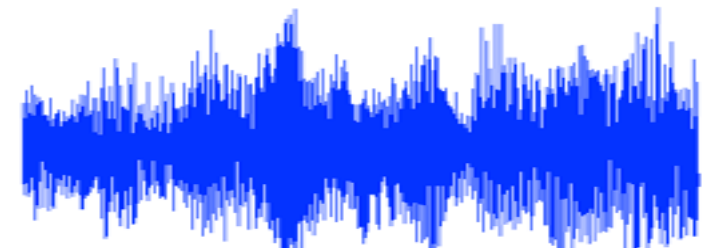
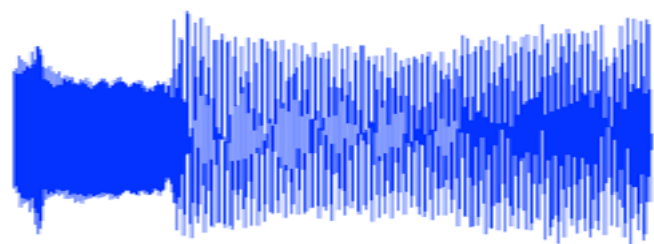


High Dimensional Learning

- High-dimensional $x = (x(1), \dots, x(d)) \in \mathbb{R}^d$:
- **Classification:** estimate a class label $f(x)$
given n sample values $\{x_i, y_i = f(x_i)\}_{i \leq n}$

Audio: instrument recognition

Huge variability
inside classes



High Dimensional Learning

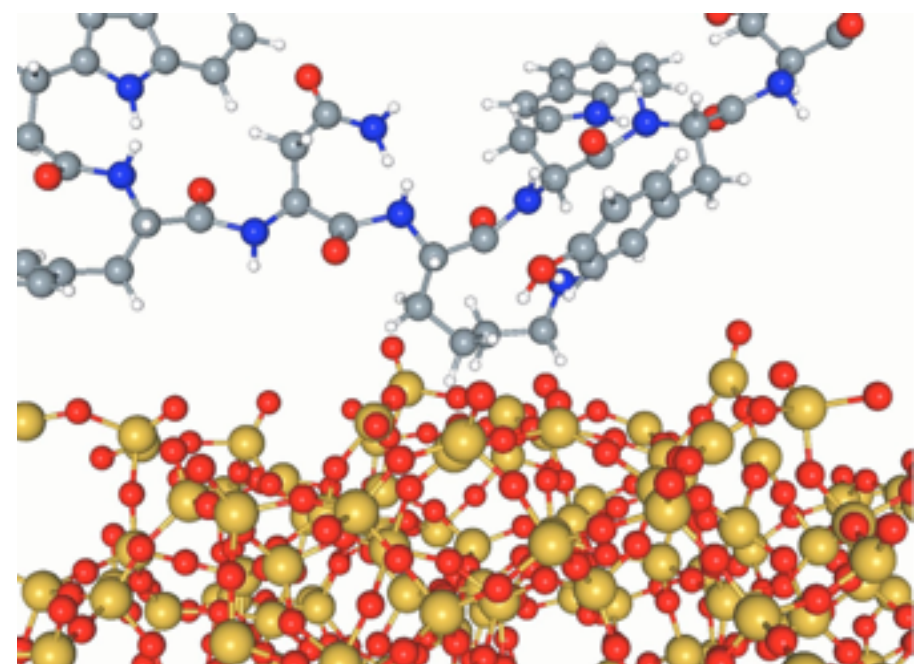
- High-dimensional $x = (x(1), \dots, x(d)) \in \mathbb{R}^d$:
- **Regression:** approximate a *functional* $f(x)$
given n sample values $\{x_i, y_i = f(x_i) \in \mathbb{R}\}_{i \leq n}$

Physics: energy $f(x)$ of a state vector x

Astronomy

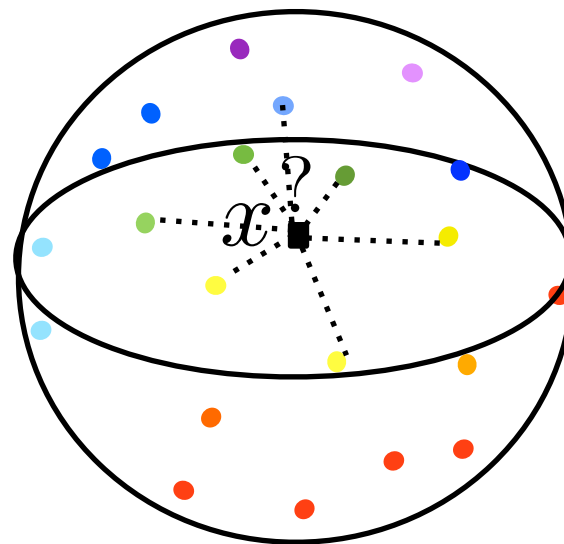


Quantum Chemistry



Curse of Dimensionality

- $f(x)$ can be approximated from examples $\{x_i, f(x_i)\}_i$ by local interpolation if f is regular and there are close examples:



- Need ϵ^{-d} points to cover $[0, 1]^d$ at a Euclidean distance ϵ
 $\Rightarrow \|x - x_i\|$ is always large



Huge variability
inside classes

Learning by Euclidean Embedding

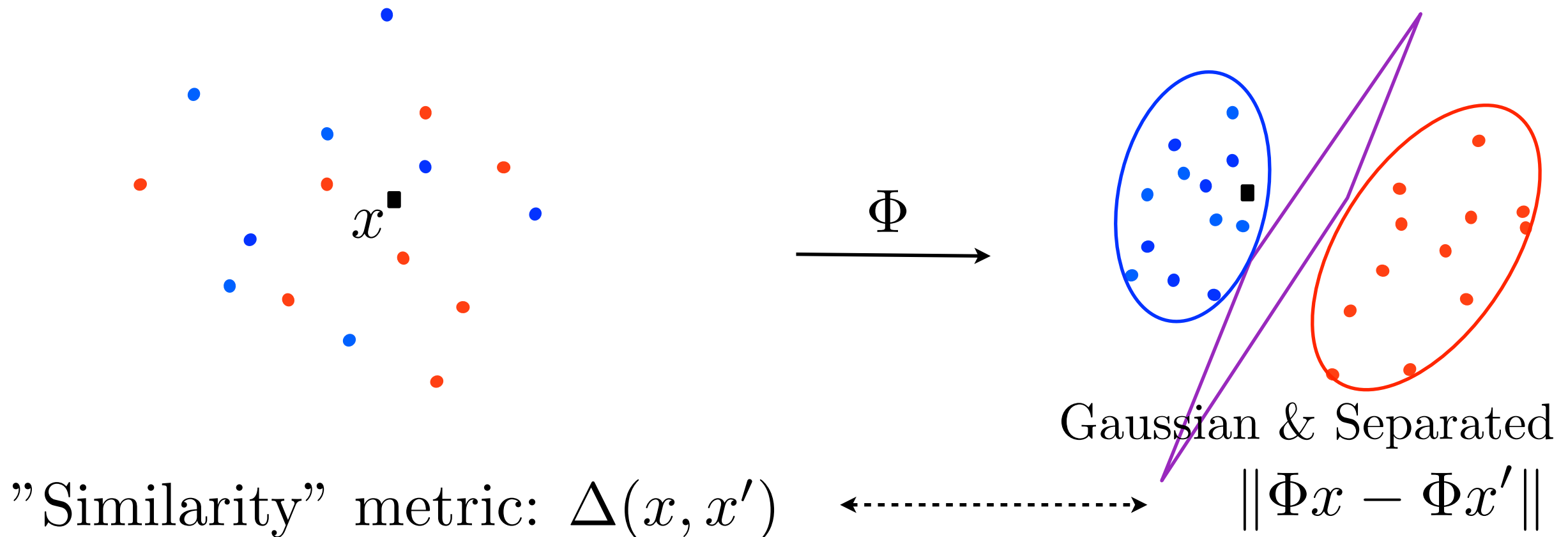
Data: $x \in \mathbb{R}^d$

$\|x - x'\|$: non-informative

Representation

$\Phi x \in \mathcal{H}$

Linear Classifier



"Similarity" metric: $\Delta(x, x')$

Gaussian & Separated

$\|\Phi x - \Phi x'\|$

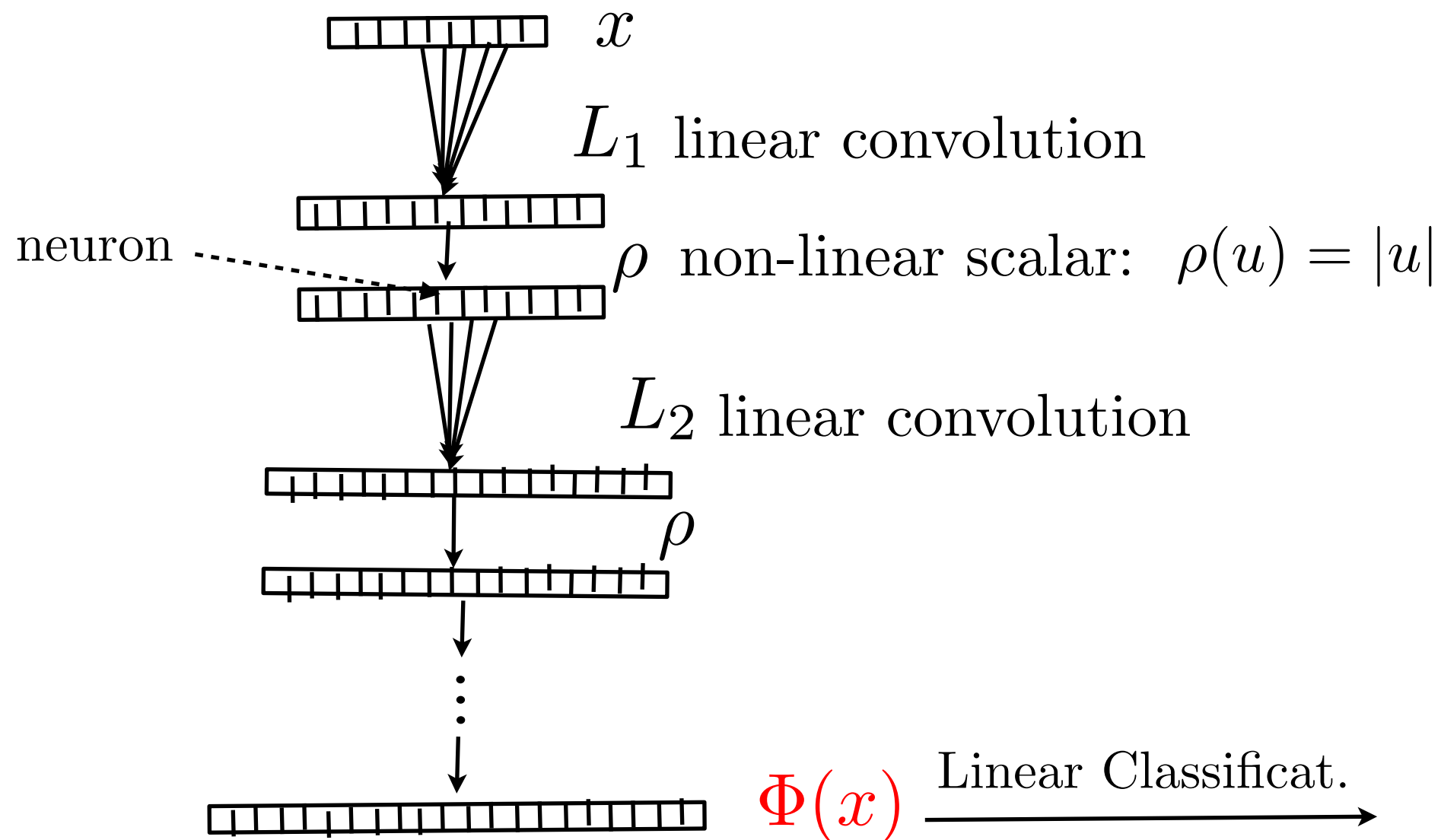
Equivalent Euclidean metric:

$$C_1 \|\Phi x - \Phi x'\| \leq \Delta(x, x') \leq C_2 \|\Phi x - \Phi x'\|$$

How to define Φ ?

Deep Convolution Networks

- The revival of an old (1950) idea: *Y. LeCun, G. Hinton*



Optimize the L_k with **support constraints**: over 10^9 parameters
Exceptional results for *images, speech, bio-data* classification.
Products by FaceBook, IBM, Google, Microsoft, Yahoo...

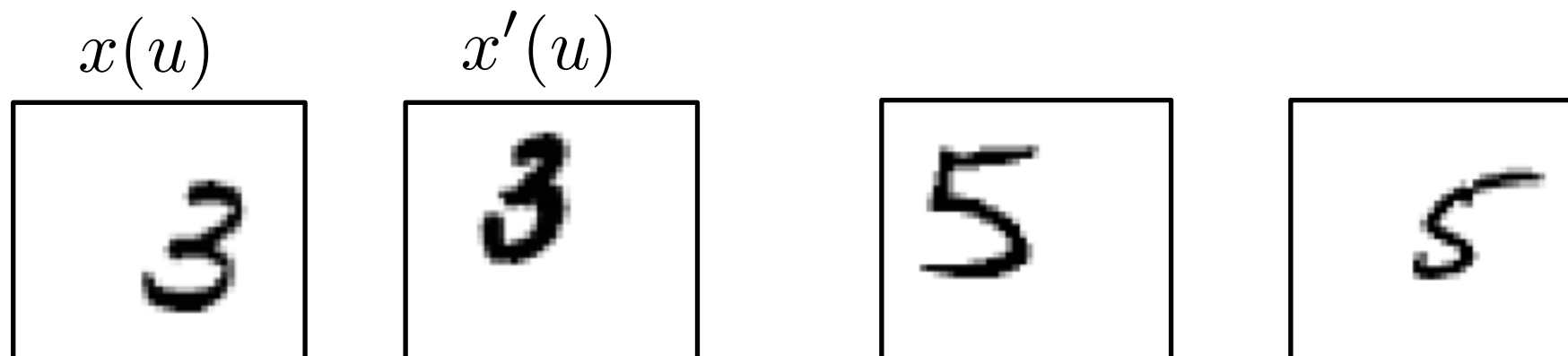
Why does it work so well ?



Overview

- Deep multiscale networks: invariant and stable metrics on groups
- Image classification
- Models of audio and image textures: information theory
- Learning physics: quantum chemistry energy regression

- Low-dimensional "geometric shapes"



Deformation metric: (classic mechanics) *Grenander*

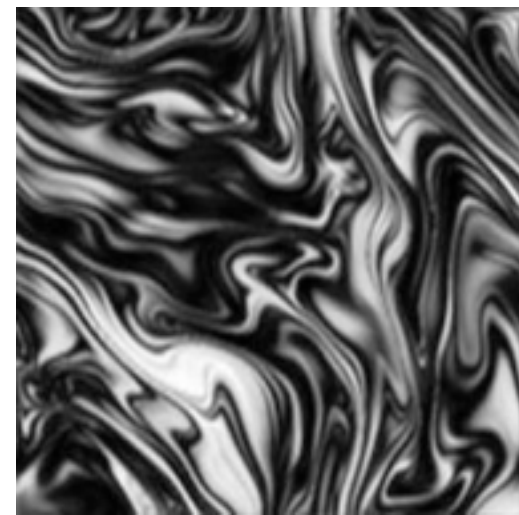
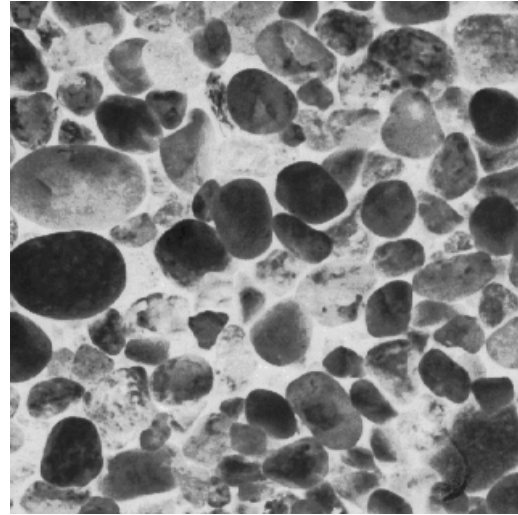
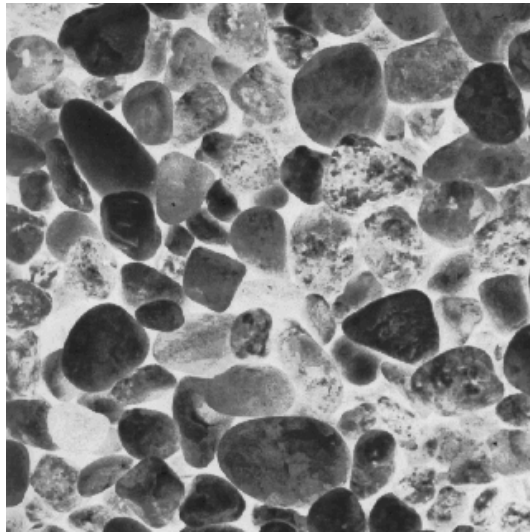
Diffeomorphism action: $D_\tau x(u) = x(u - \tau(u))$

$$\Delta(x, x') \sim \min_{\tau} \|D_\tau x - x'\| + \|\nabla \tau\|_\infty \|x\|$$

Invariant to translations

↓
diffeomorphism
amplitude

- High dimensional textures: $X(u)$ ergodic stationary processes



2D Turbulence

Highly non-Gaussian processes

- A Euclidean metric is a Maximum Likelihood on Gaussian models.
- Can we find Φ so that $\Phi(X)$ is nearly Gaussian, without losing information ?

- Stability to additive perturbations:

$$\|\Phi x - \Phi x'\| \leq C \|x - x'\|$$

- Invariance to translations:

$$x_c(u) = x(u - c) \Rightarrow \Phi(x_c) = \Phi(x)$$

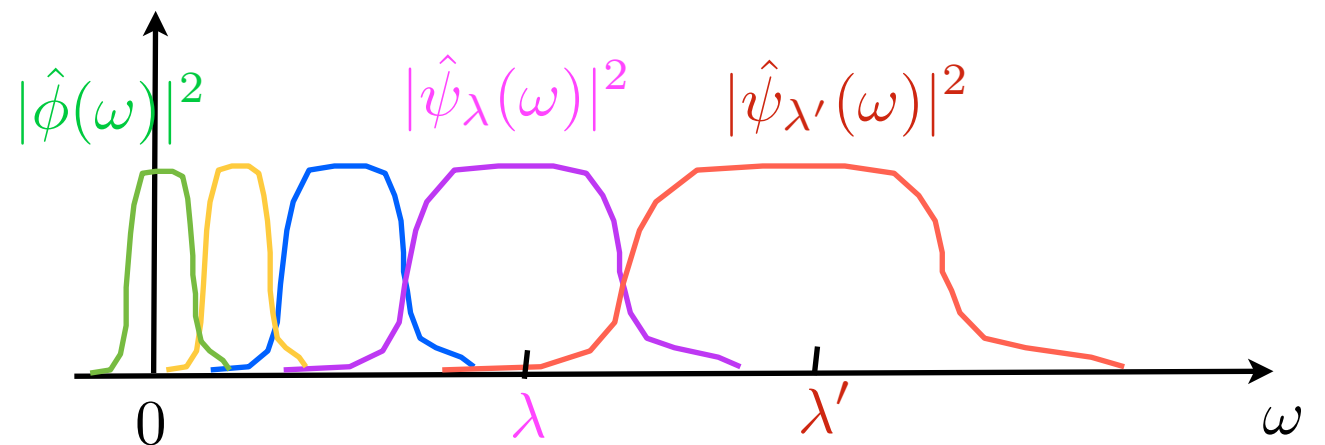
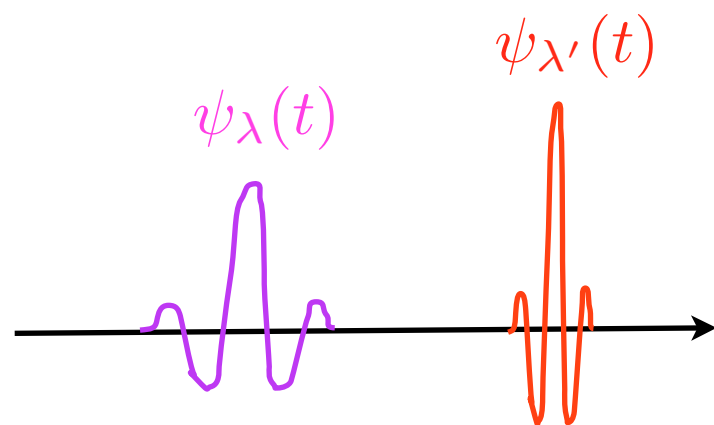
- Stability to deformations:

$$x_\tau(u) = x(u - \tau(u)) \Rightarrow \|\Phi x - \Phi x_\tau\| \leq C \|\nabla \tau\|_\infty \|x\|$$

Failure of Fourier and classic invariants

Wavelet Transform

- Dilated wavelets: $\psi_\lambda(t) = 2^{-j/Q} \psi(2^{-j/Q}t)$ with $\lambda = 2^{-j/Q}$



Q-constant band-pass filters $\hat{\psi}_\lambda$

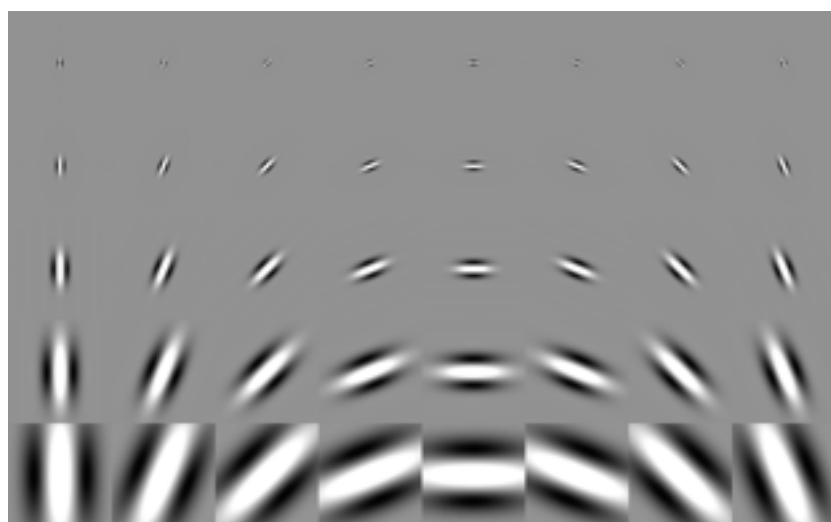
- Wavelet transform: $Wx = \begin{pmatrix} x \star \phi_{2^J}(t) \\ x \star \psi_\lambda(t) \end{pmatrix}_{\lambda \leq 2^J}$: average
: higher frequencies

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

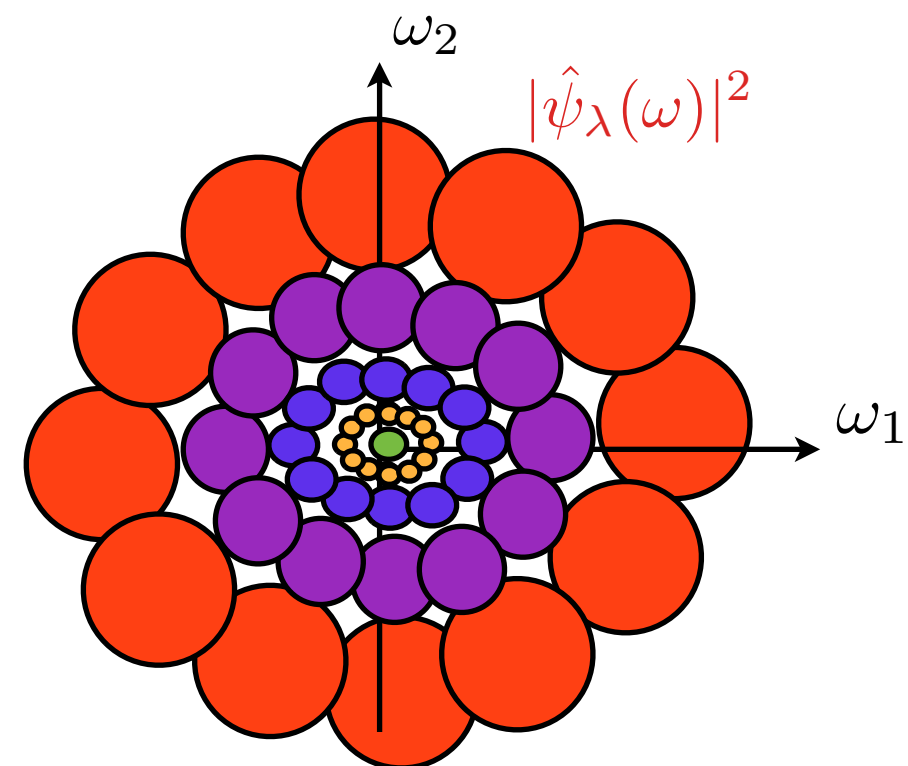
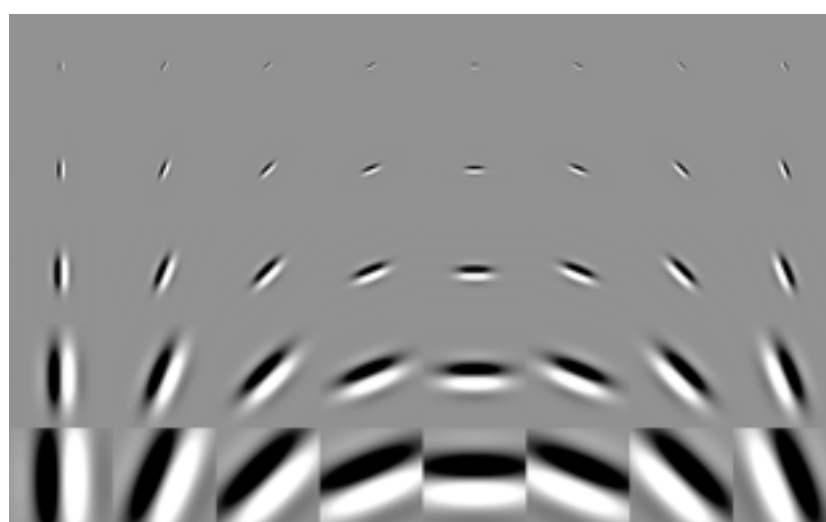
Scale separation with Wavelets

- Complex wavelet: $\psi(t) = g(t) \exp i\xi t$, $t = (t_1, t_2)$
rotated and dilated: $\psi_\lambda(t) = 2^{-j} \psi(2^{-j} r_\theta t)$ with $\lambda = (2^j, \theta)$

real parts



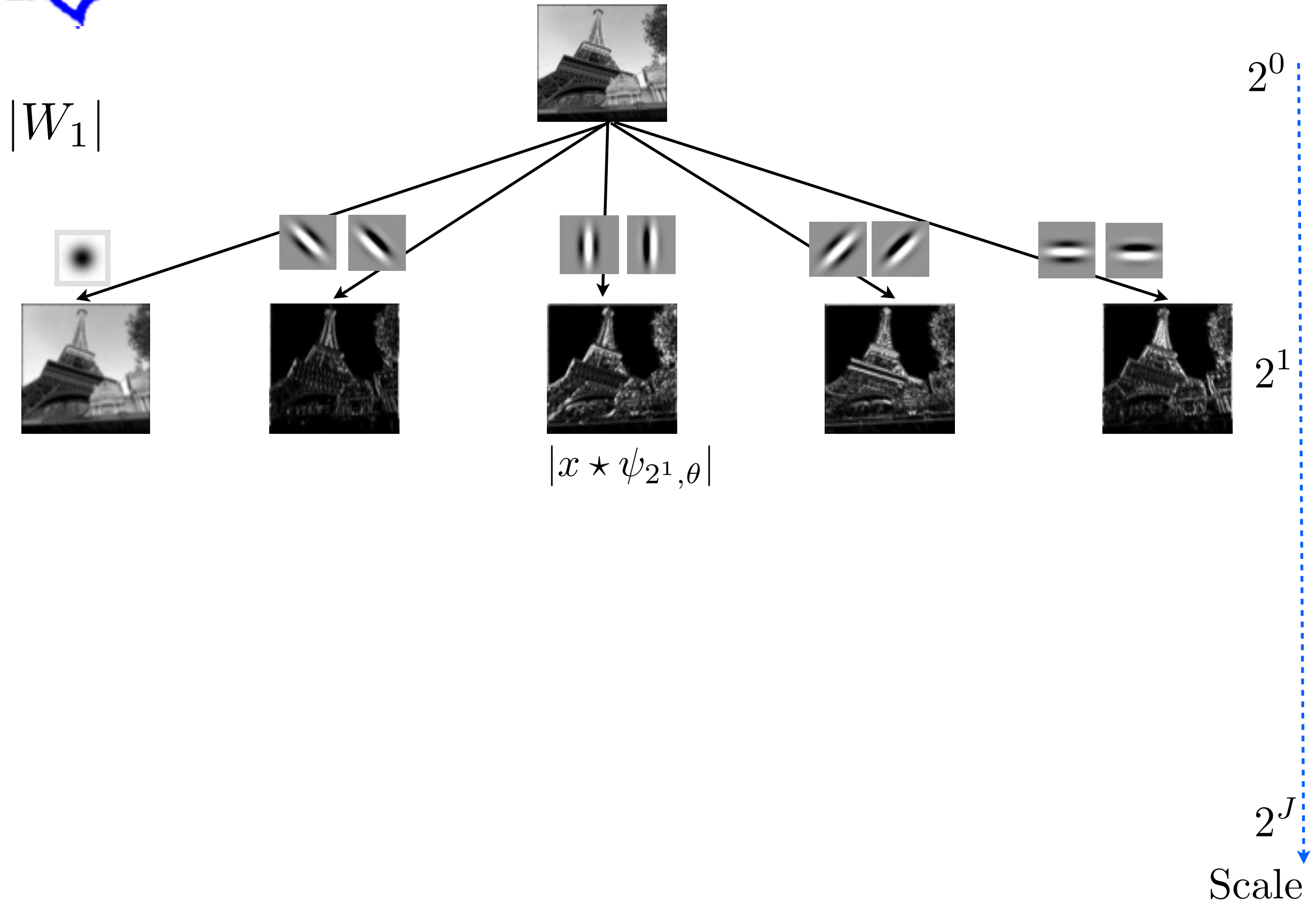
imaginary parts



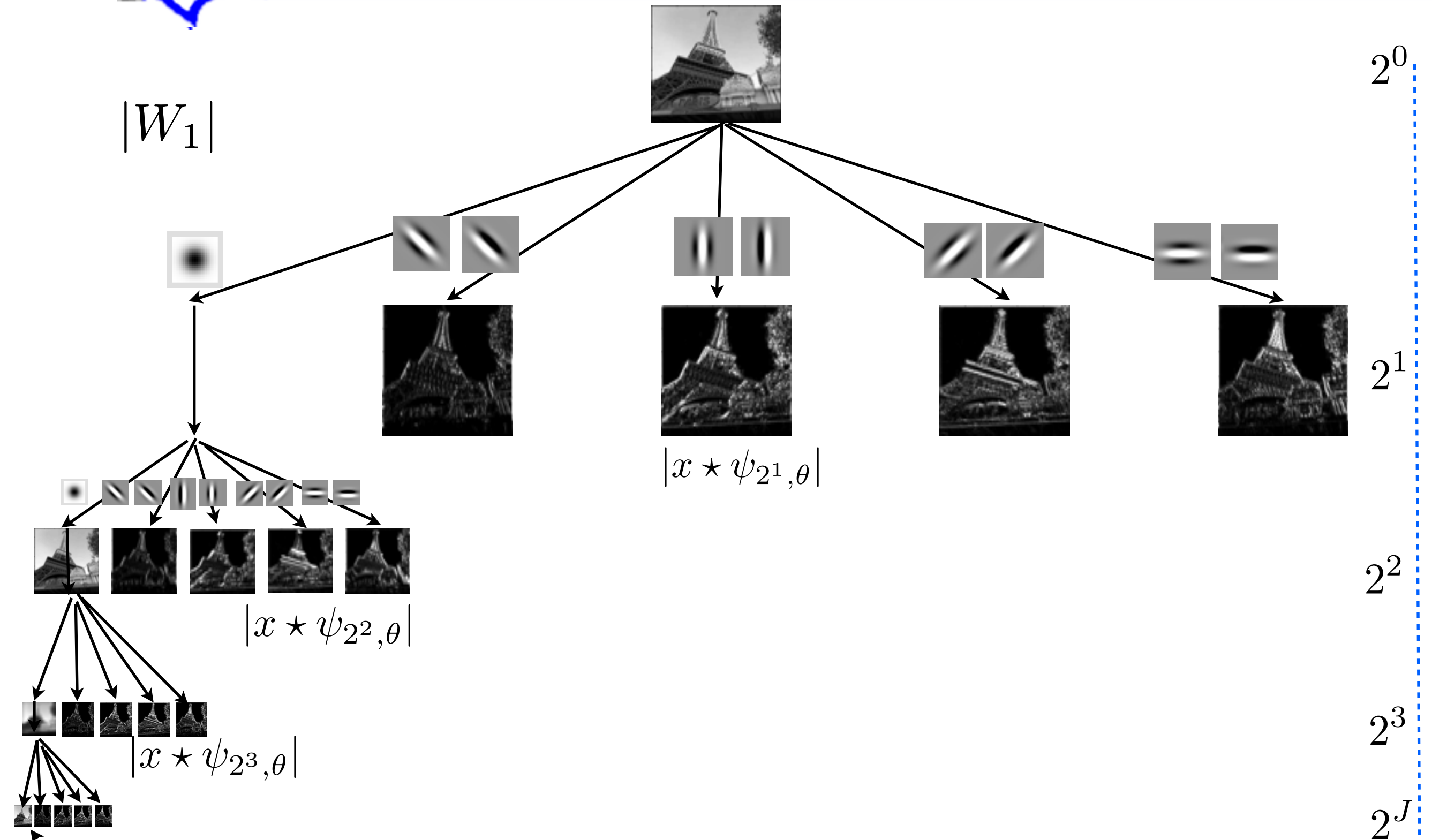
- Wavelet transform: $Wx = \begin{pmatrix} x \star \phi_{2^J}(t) \\ x \star \psi_\lambda(t) \end{pmatrix}_{\lambda \leq 2^J}$: average
: higher frequencies

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

Fast Wavelet Transform



Wavelet Transform



$x \star \phi_J$: locally invariant by translation

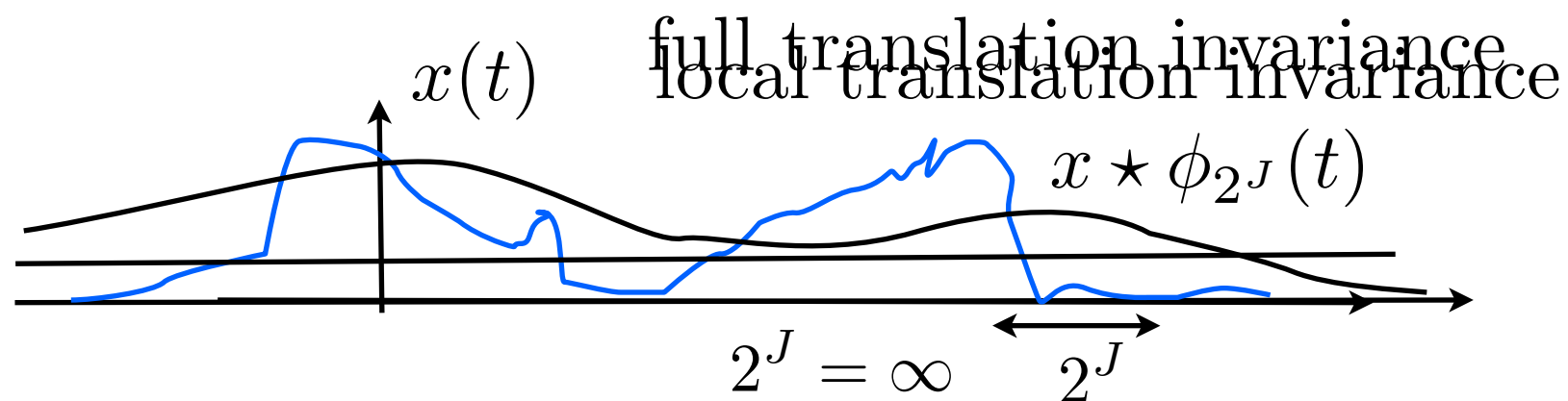
Depth: Scale

How to make everything invariant to translation ?

Wavelet Translation Invariance

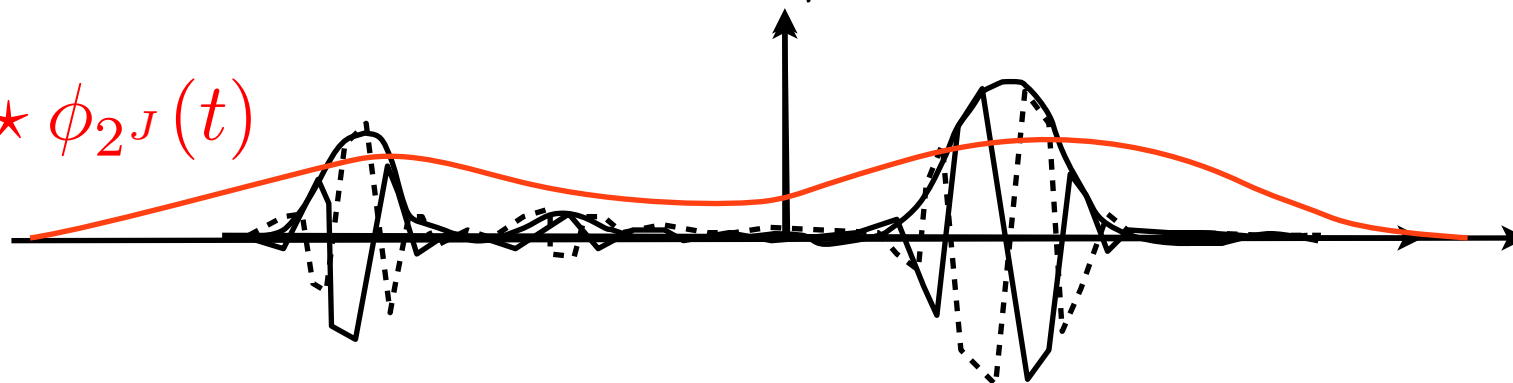
First wavelet transform

$$|W_1| x \equiv \left(\begin{array}{c} x \star \phi_{2^J} \\ x \star \phi_{2^J} \\ x \star \psi_{\lambda_1} \\ |x \star \psi_{\lambda_1}| \end{array} \right)_{\lambda_1}$$



Modulus improves invariance: $|x \star \psi_{\lambda_1}(t)| \neq \sqrt{|x \star \psi_{\lambda_1}^a(t)|^2 + |x \star \psi_{\lambda_1}^b(t)|^2}$ but covariant

$$|x \star \psi_{\lambda_1}| \star \phi_{2^J}(t)$$

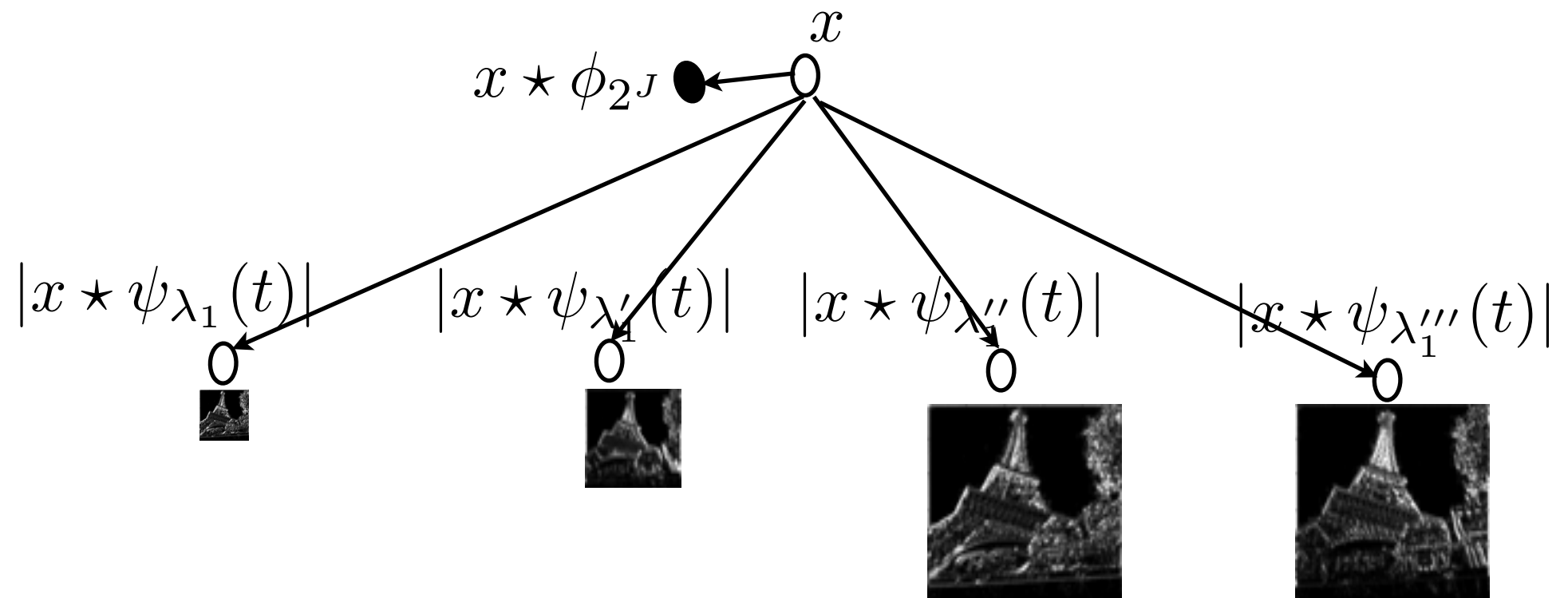


Second wavelet transform modulus

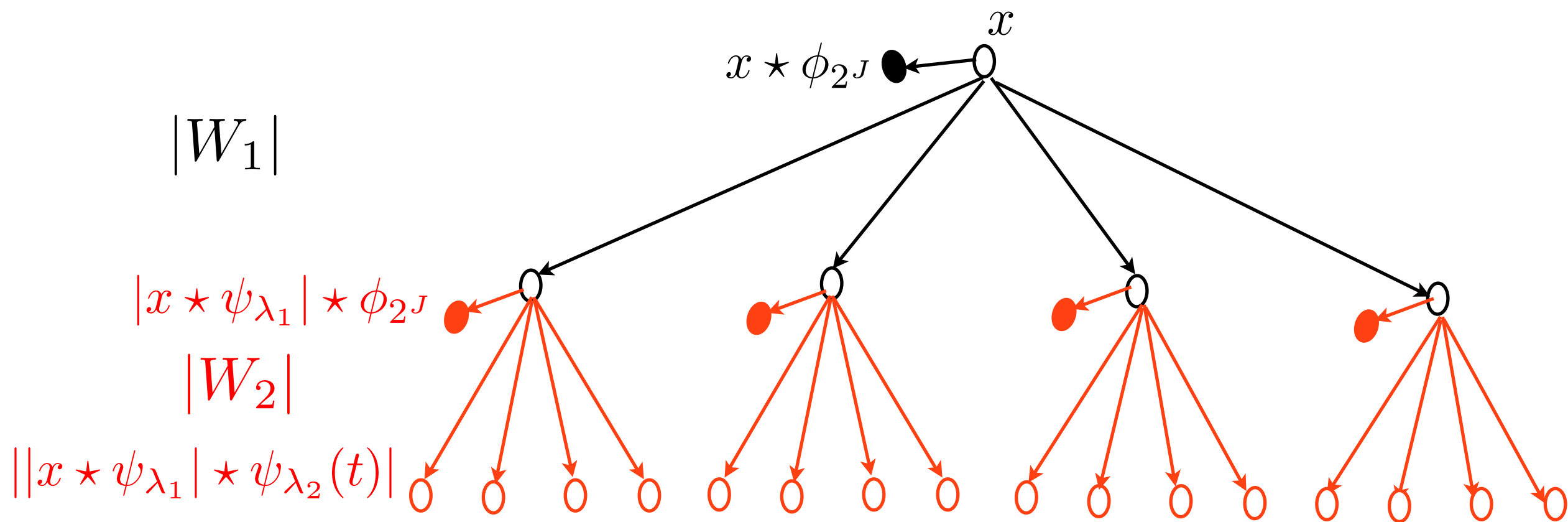
$$|W_2| |x \star \psi_{\lambda_1}| = \left(\begin{array}{c} |x \star \psi_{\lambda_1}| \star \phi_{2^J}(t) \\ ||x \star \psi_{\lambda_1}| \star \psi_{\lambda_2}(t)| \end{array} \right)_{\lambda_2}$$

Scattering Transform

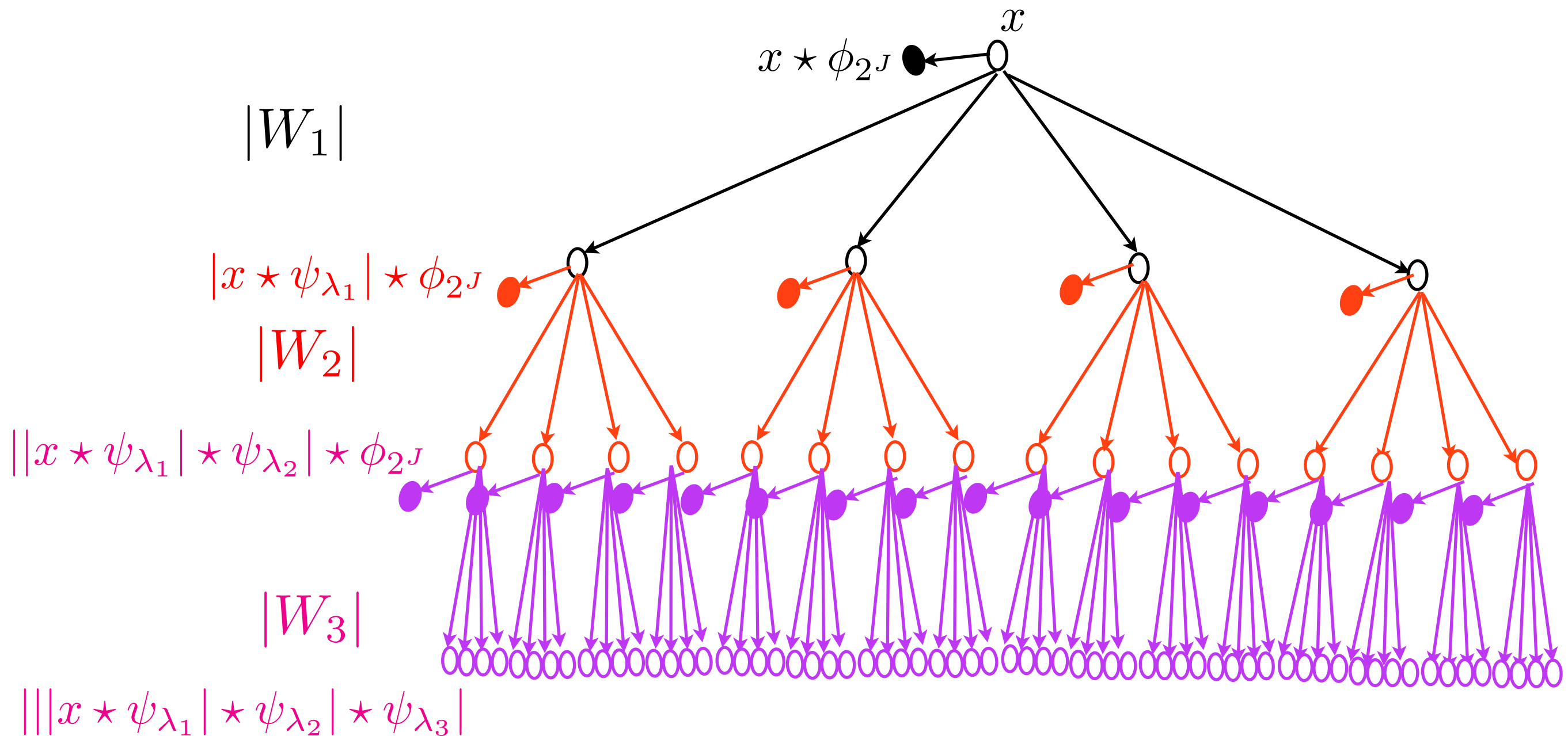
$|W_1|$



Scattering Transform



Scattering Neural Network



Scattering Properties

$$S_J x = \left(\begin{array}{c} x \star \phi_{2^J} \\ |x \star \psi_{\lambda_1}| \star \phi_{2^J} \\ ||x \star \psi_{\lambda_1}| \star \psi_{\lambda_2}| \star \phi_{2^J} \\ |||x \star \psi_{\lambda_2}| \star \psi_{\lambda_2}| \star \psi_{\lambda_3}| \star \phi_{2^J} \\ \dots \end{array} \right)_{\lambda_1, \lambda_2, \lambda_3, \dots} = \dots |W_3| |W_2| |W_1| x$$

W_k is unitary $\Rightarrow |W_k|$ is contractive

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:

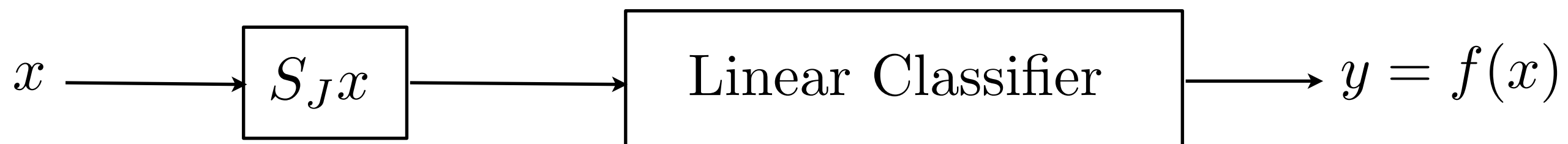
if $x_\tau(u) = x(u - \tau(u))$ then

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

Digit Classification: MNIST

Joan Bruna

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 5
4 8 1 9 0 1 8 8 9 4



Classification Errors

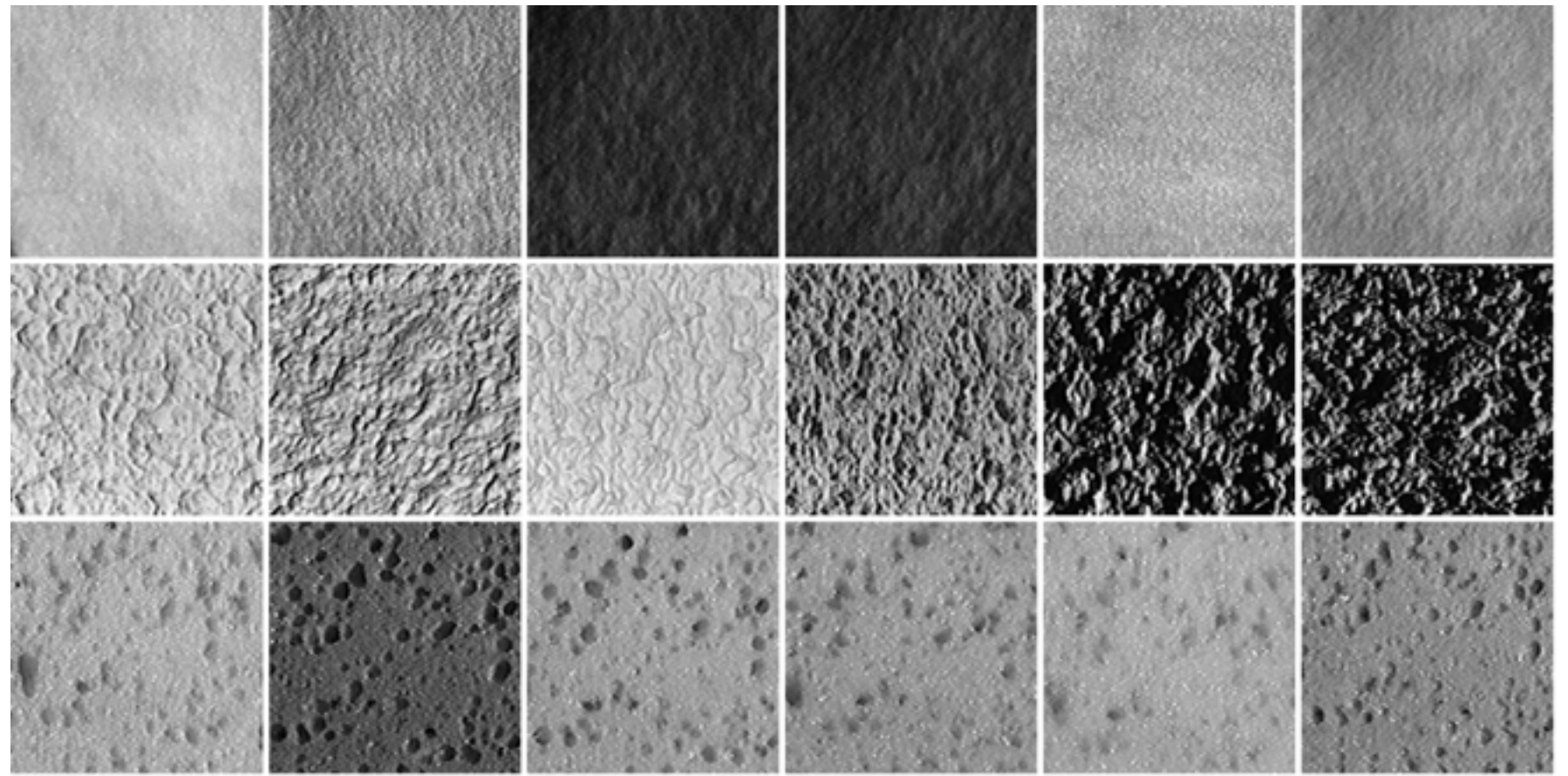
Training size	Conv. Net.	Scattering
50000	0.5%	0.4%

LeCun et. al.

Classification of Textures

J. Bruna

CUREt database
61 classes



Classification Errors

$2^J = \text{image size}$

Training per class	Fourier Spectr.	Histogr. Features	Scattering
46	1%	1%	0.2 %

The scattering transform of a stationary process $X(t)$

$$\begin{array}{l}
S_J X = \begin{pmatrix} X \\ |X \star \psi_{\lambda_1}| \\ ||X \star \psi_{\lambda_1}| \star \psi_{\lambda_2}| \\ |||X \star \psi_{\lambda_2}| \star \psi_{\lambda_2}| \star \psi_{\lambda_3}| \\ \dots \end{pmatrix} \star \phi_{2^J} : \text{Gaussian for } 2^J \text{ large} \\
\downarrow J \rightarrow \infty \\
\mathbb{E}(S X) = \begin{pmatrix} \mathbb{E}(X) \\ \mathbb{E}(|X \star \psi_{\lambda_1}|) \\ \mathbb{E}(||X \star \psi_{\lambda_1}| \star \psi_{\lambda_2}|) \\ \mathbb{E}(|||X \star \psi_{\lambda_2}| \star \psi_{\lambda_2}| \star \psi_{\lambda_3}|) \\ \dots \end{pmatrix}_{\lambda_1, \lambda_2, \lambda_3, \dots} \\
\text{if } X \text{ is ergodic}
\end{array}$$

Representation of Random Processes

$$\mathbb{E}(SX) = \left(\begin{array}{ccc} \mathbb{E}(X) & = & \mathbb{E}(U_0 X) \\ \mathbb{E}(|X \star \psi_{\lambda_1}|) & = & \mathbb{E}(U_1 X) \\ \mathbb{E}(|X \star \psi_{\lambda_1} \star \psi_{\lambda_2}|) & = & \mathbb{E}(U_2 X) \\ \mathbb{E}(|X \star \psi_{\lambda_2} \star \psi_{\lambda_2} \star \psi_{\lambda_3}|) & = & \mathbb{E}(U_3 X) \\ \dots & & \end{array} \right)_{\lambda_1, \lambda_2, \lambda_3, \dots}$$

Theorem (Boltzmann) The distribution $p(x)$ which satisfies

$$\int_{\mathbb{R}^N} U_m x p(x) dx = E(U_m X)$$

with a maximum entropy $H_{\max} = - \int p(x) \log p(x) dx$ is

$$p(x) = \frac{1}{Z} \exp \left(\sum_{m=1}^{\infty} \lambda_m \cdot U_m x \right)$$

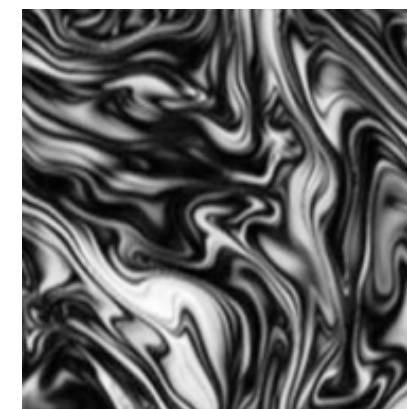
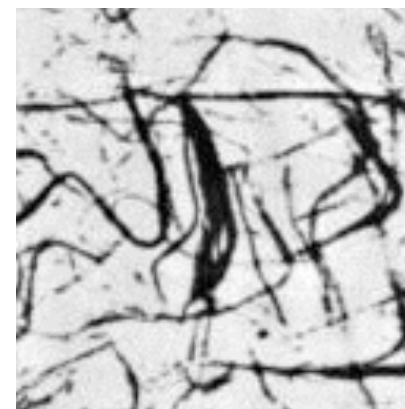
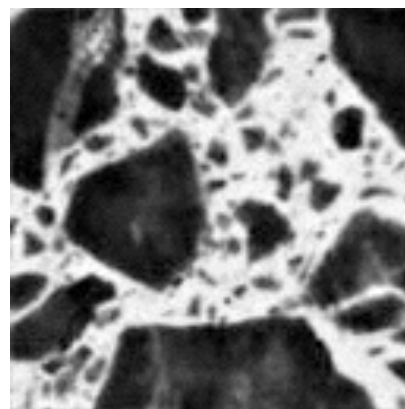
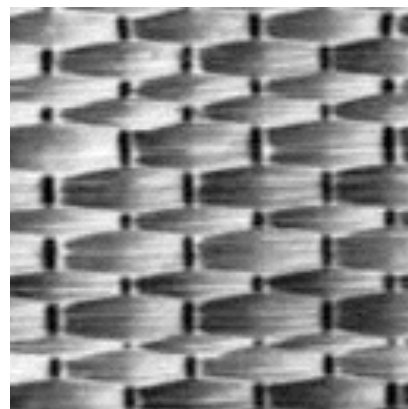
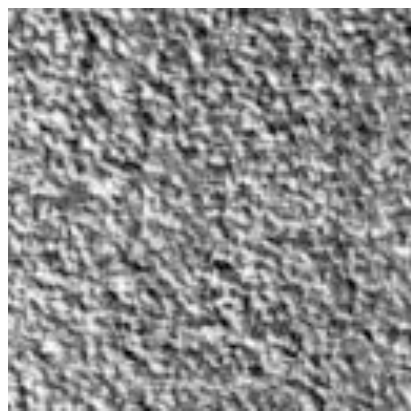
$$H_{\max} \geq H(X) \text{ (entropie of X)}$$

Little loss of information: $H_{\max} \approx H(X)$

Ergodic Texture Reconstructions

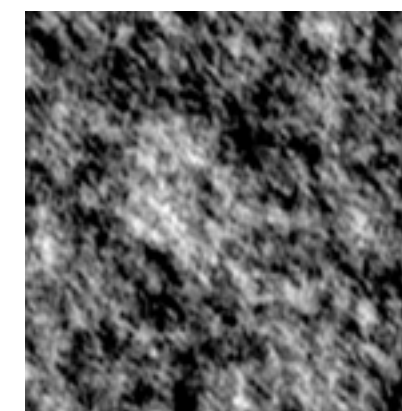
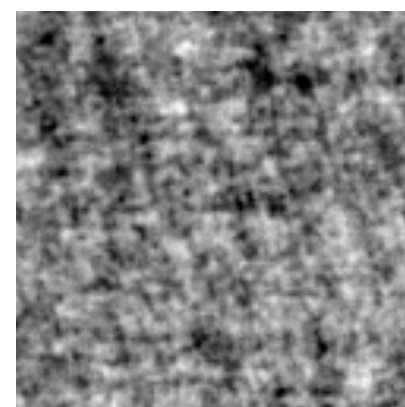
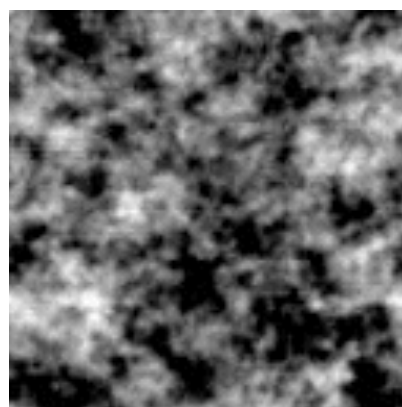
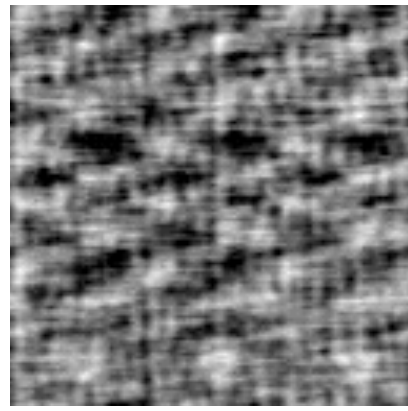
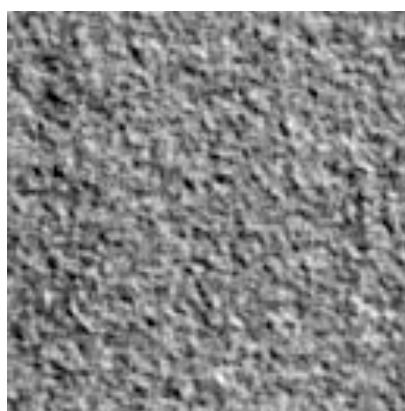
Joan Bruna

Original Textures



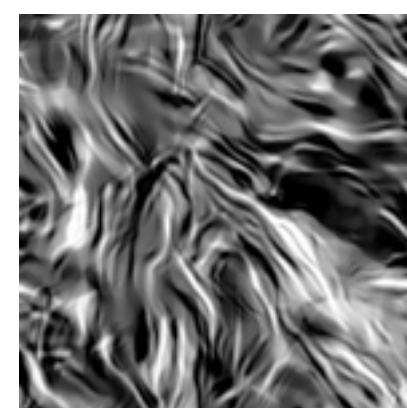
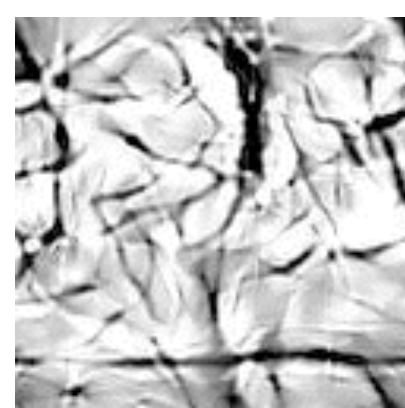
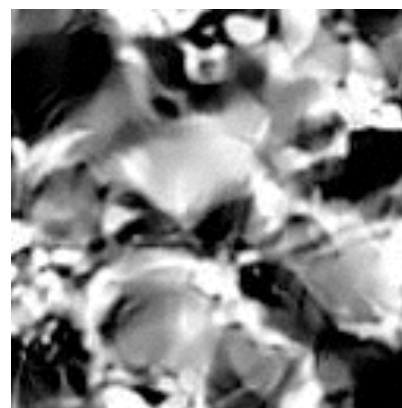
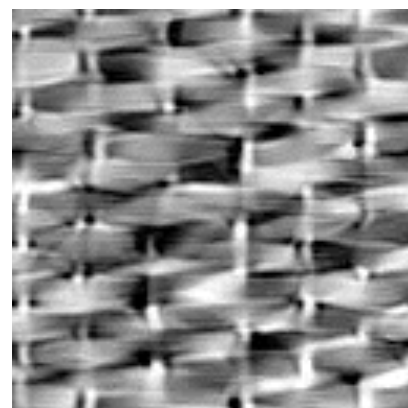
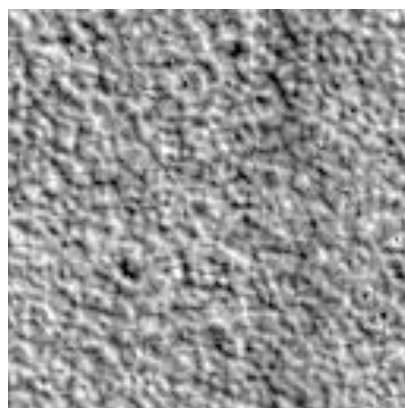
2D Turbulence

Gaussian process model with same second order moments



Second order Gaussian Scattering: $O(\log N^2)$ moments

$$\mathbb{E}(|x \star \psi_{\lambda_1}|) \quad , \quad \mathbb{E}(|x \star \psi_{\lambda_1}| \star \psi_{\lambda_2}|)$$



Representation of Audio Textures

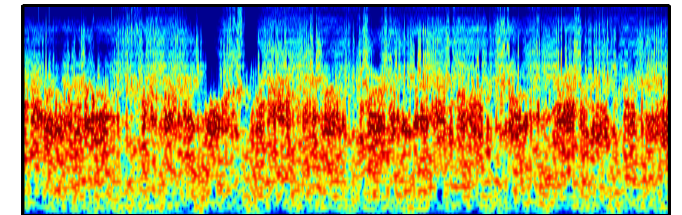
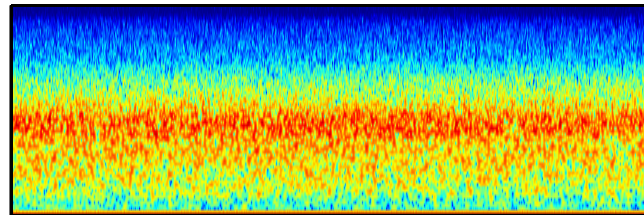
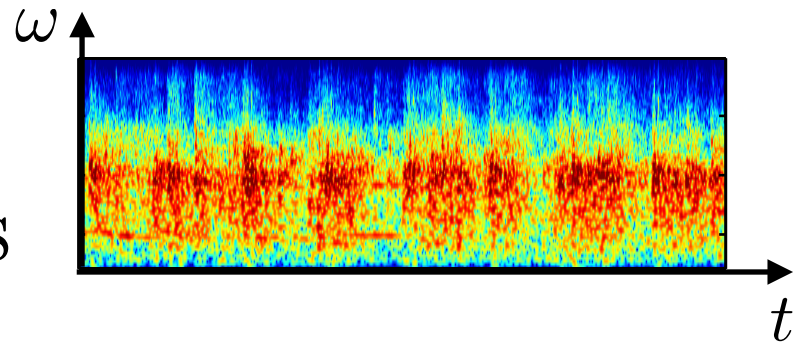
Joan Bruna

Original

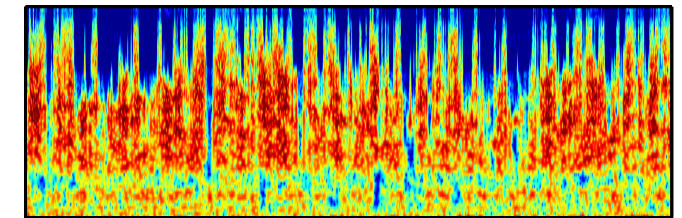
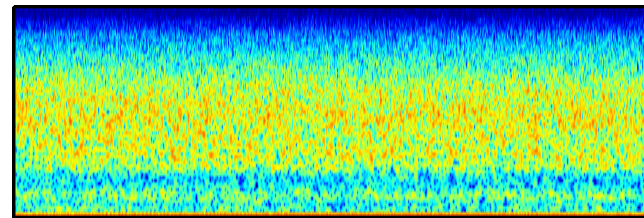
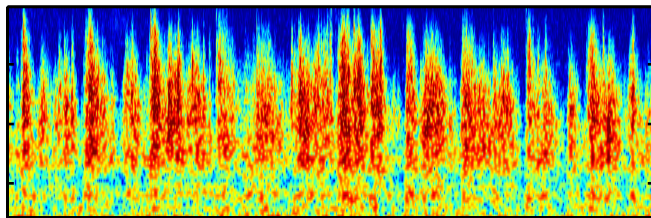
Gaussian
in time

Gaussian
in scattering

Applauds



Paper

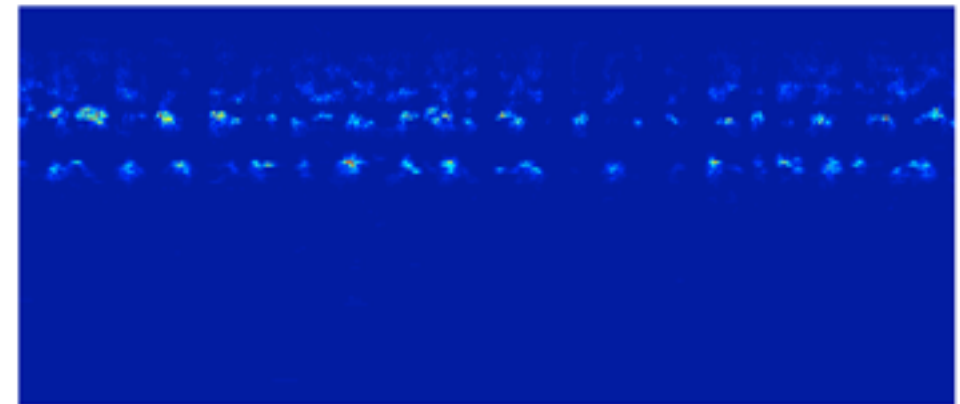
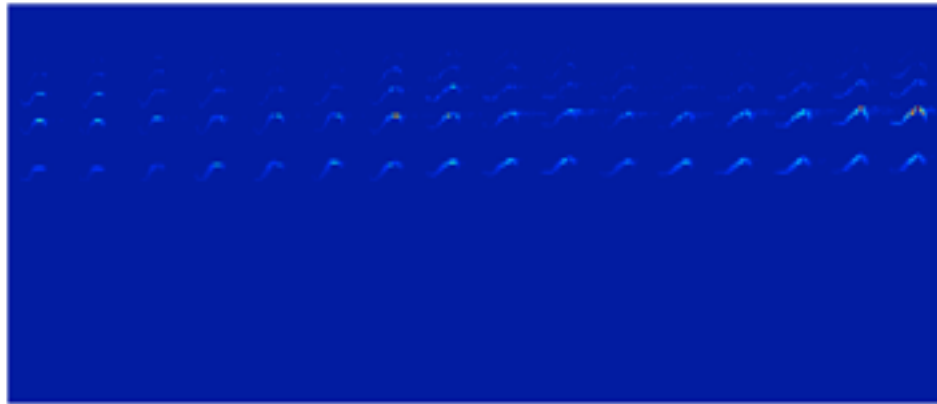


Cocktail Party

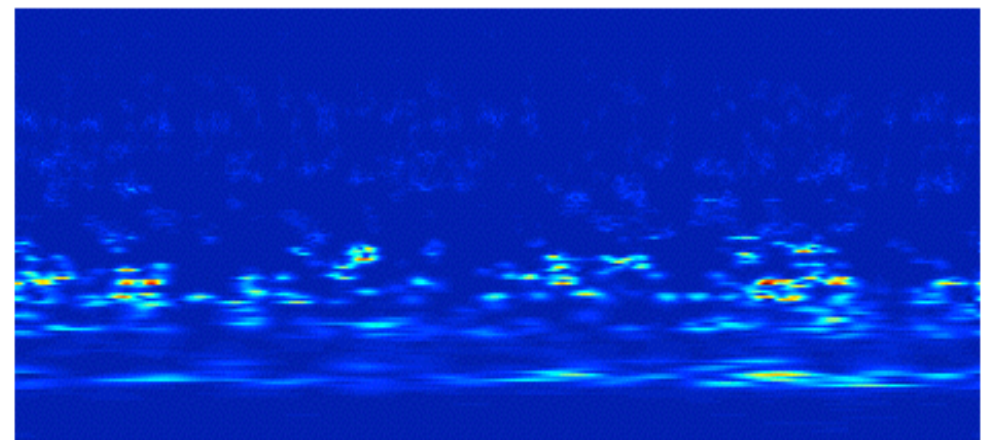
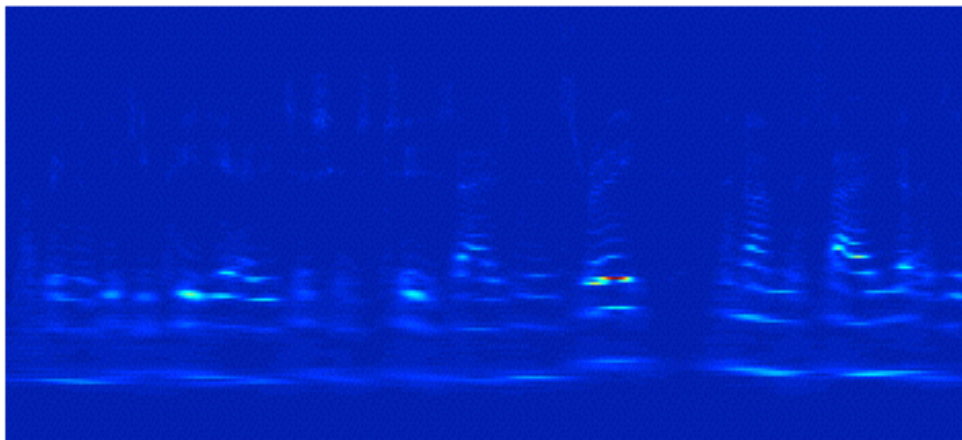
Failures: Harmonic Sounds *V. Lostanlen*

Need to express frequency channel interactions: time-frequency image

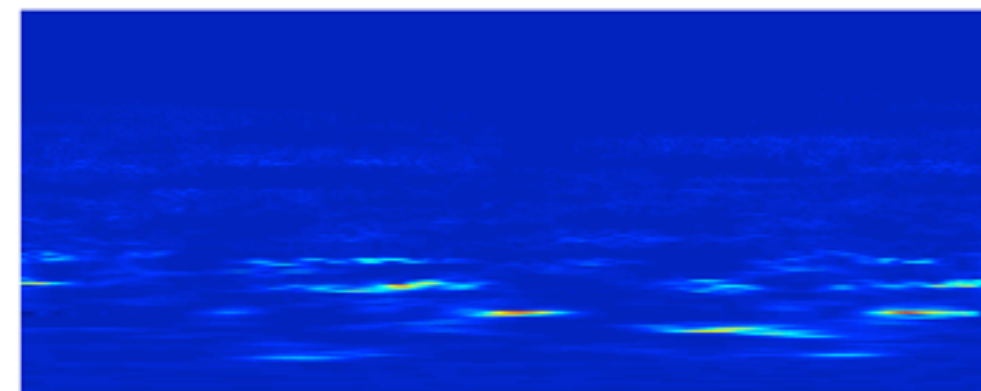
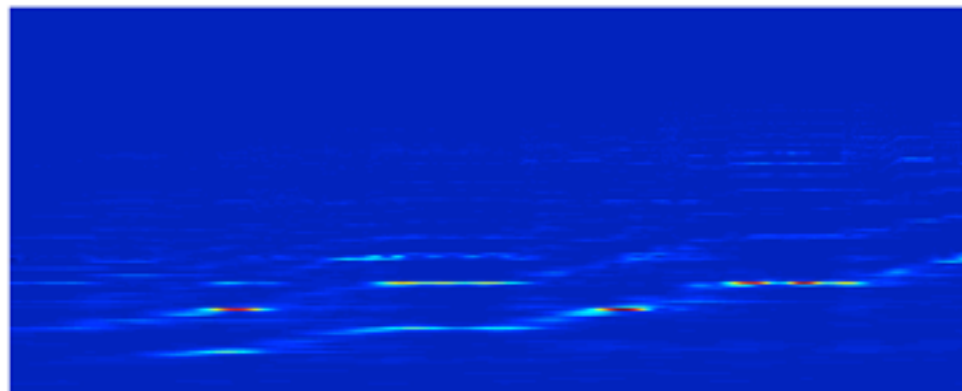
Bird



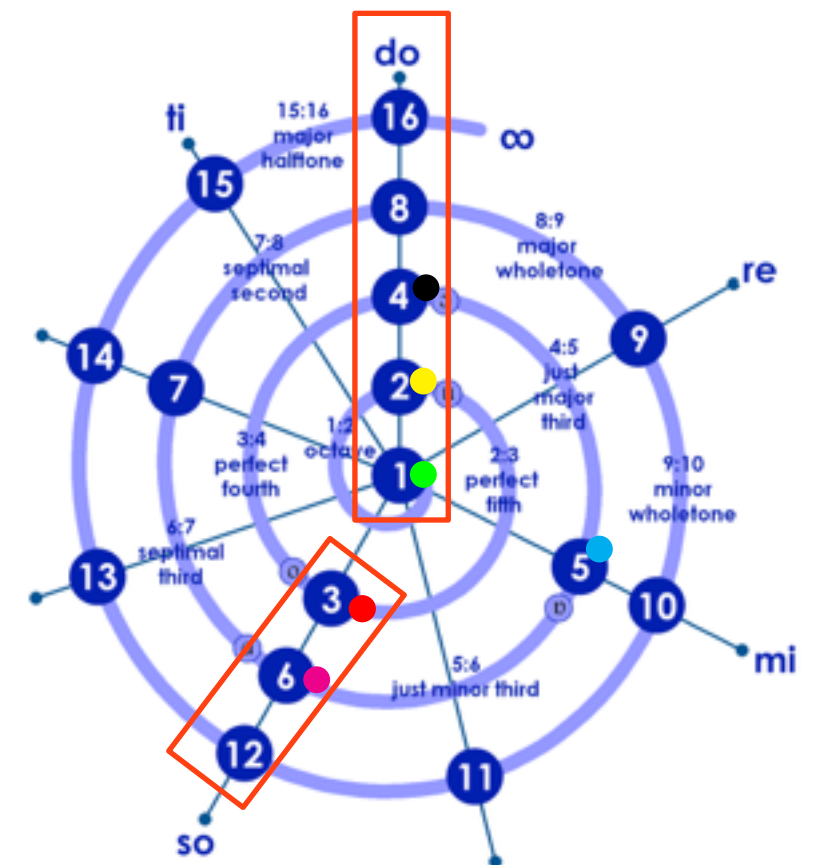
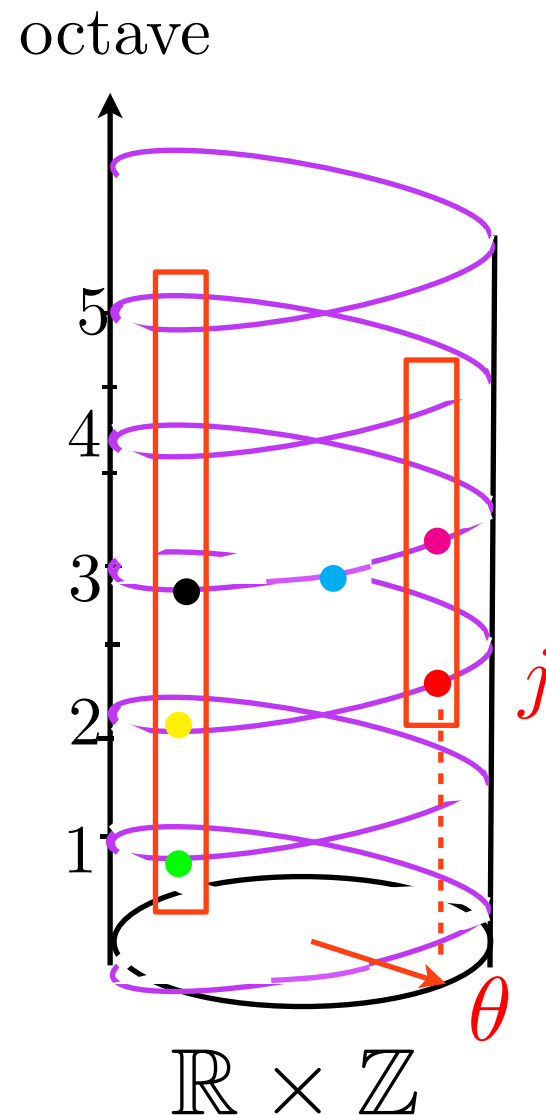
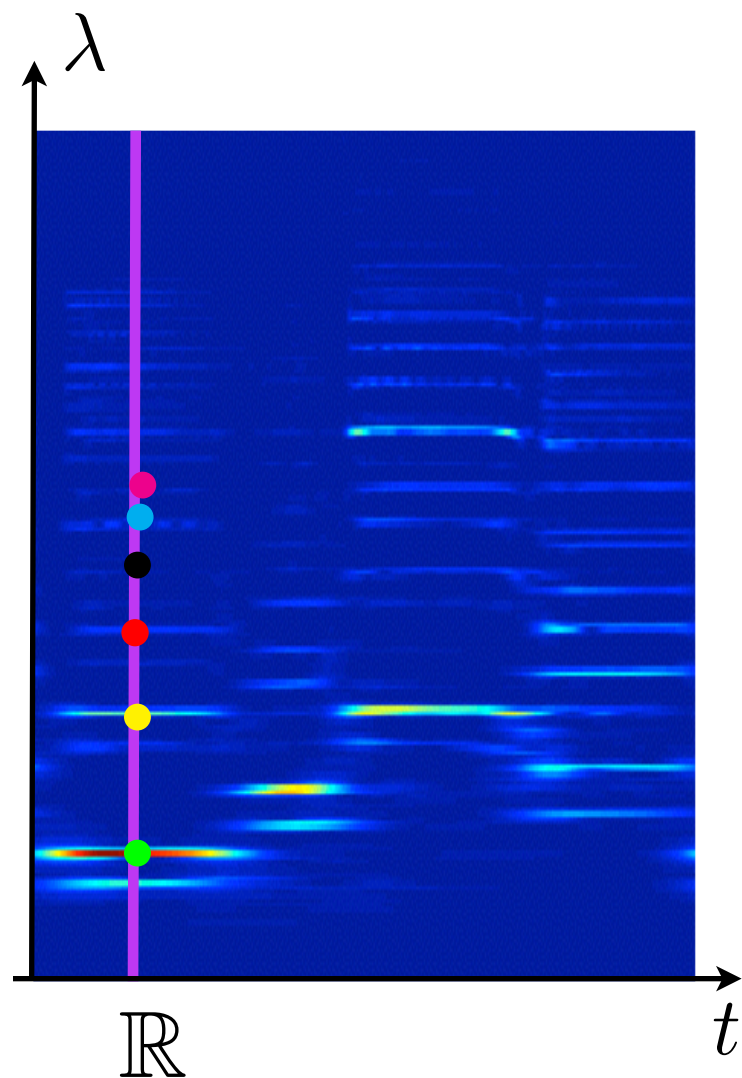
Speech



Cello

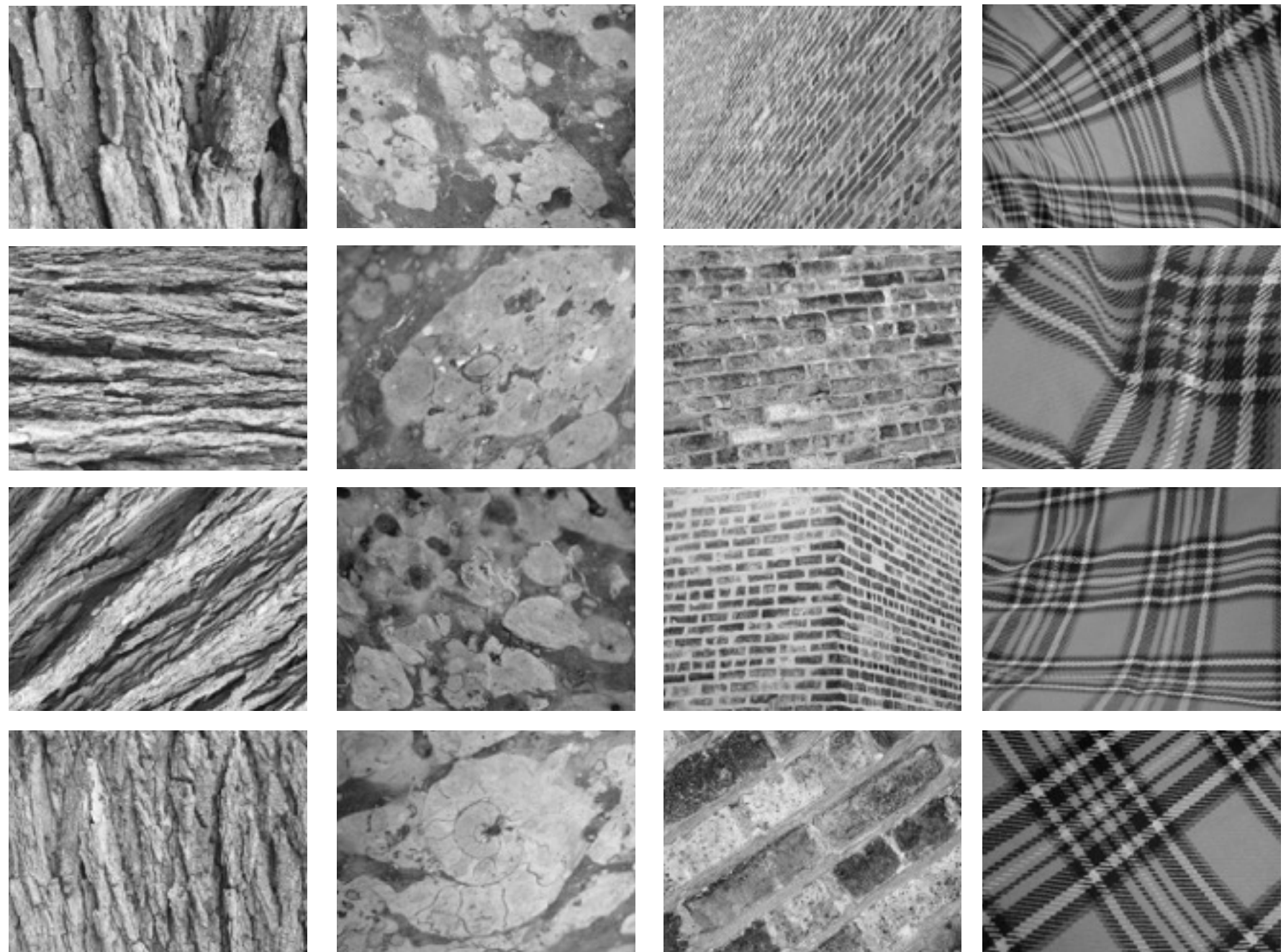


Need to capture frequency variability and structures.



- Alignment of harmonics in two main groups.
More regular variations along (θ, j) than λ

UIUC database:
25 classes



Scattering classification errors

Training	Scat. Translation
20	20 %

Extension to Rigid Mouvements

Laurent Sifre

Need to capture the variability of spatial directions.

- Group of rigid displacements: translations and rotations
- Action on wavelet coefficients:

rotation & translation

rotation & translation , angle translation

$$\begin{array}{ccccc}
 x(r_\alpha(u) x(u)) & \longrightarrow & \boxed{|W_1|} & \longrightarrow & x_j(r_\alpha(u) x(u)) = |x, \psi_{2^j \alpha_\theta}(u)| \\
 & & \downarrow & & \\
 & & \int x(u) du & &
 \end{array}$$

Extension to Rigid Mouvements

Laurent Sifre

- To build invariants: second wavelet transform on $\mathbf{L}^2(G)$:
convolutions of $x_j(u, \theta)$ with wavelets $\psi_{\lambda_2}(u, \theta)$

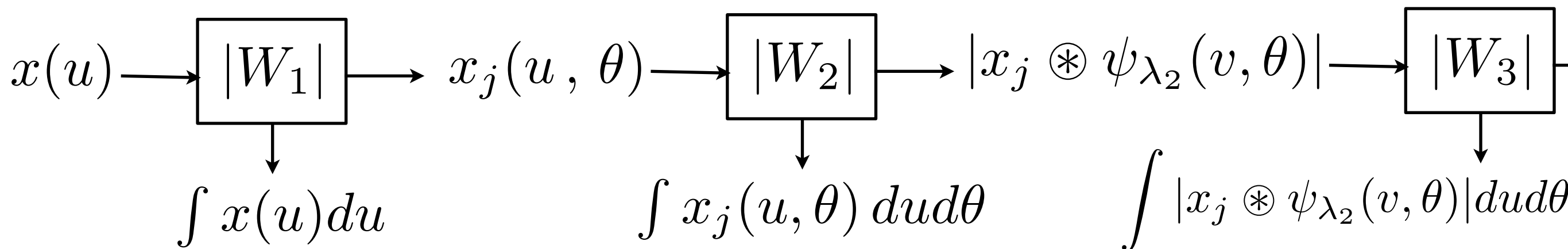
$$x_j \circledast \psi_{\lambda_2}(u, \theta) = \int_{\mathbb{R}^2} \int_0^{2\pi} x_j(v, \alpha) \psi_{\lambda_2}(u - v, \theta - \alpha) dv d\alpha$$

- Scattering on rigid mouvements:

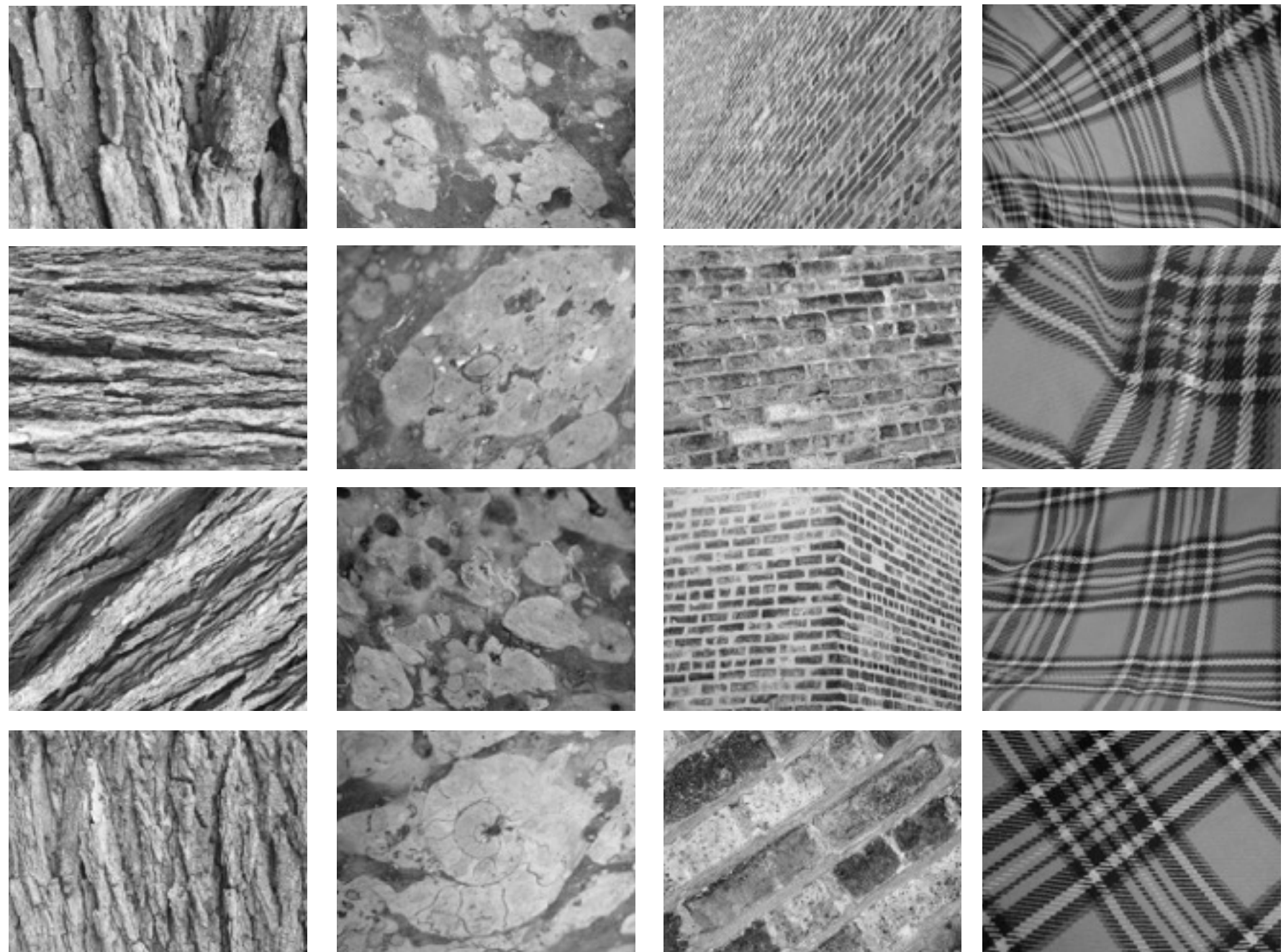
Wavelets on Translations

Wavelets on Rigid Mvt.

Wavelets on Rigid Mvt.



UIUC database:
25 classes



Scattering classification errors

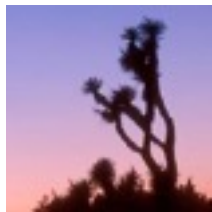
Training	Scat. Translation	Scat. Rigid Mouvt.
20	20 %	0.6%

Complex Image Classification

CalTech 101 data-basis:

Edouard Oyallon

Arbre de Joshua



Ancre



Metronome



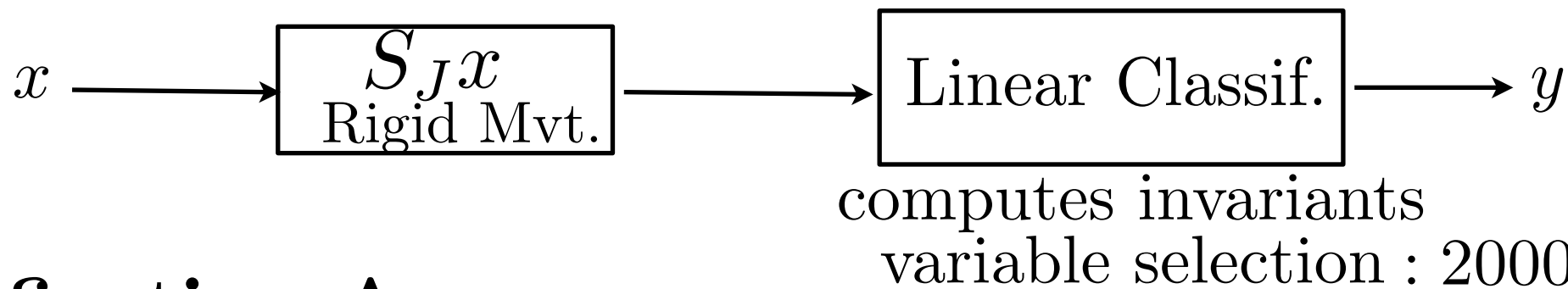
Castore



Nénuphare



Bateau



Classification Accuracy

Data Basis	Deep-Net	Scat.-2
CalTech-101	85%	80%
CIFAR-10	90%	80%

State of the art
Unsupervised

*N. Poilvert
Matthew Hirn*

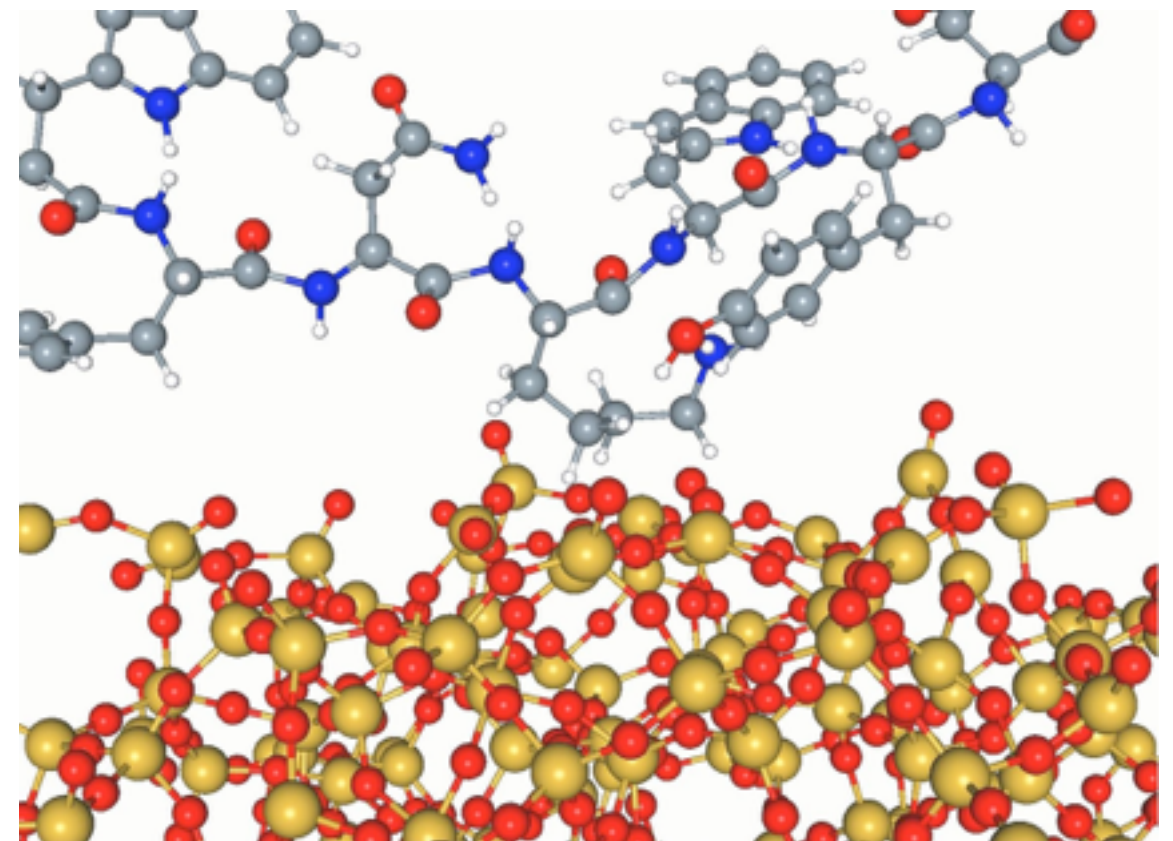
- Energy of d interacting bodies:

Can we learn the interaction energy $f(x)$ of a system
with $x = \left\{ \text{positions, values} \right\}$?

Astronomy

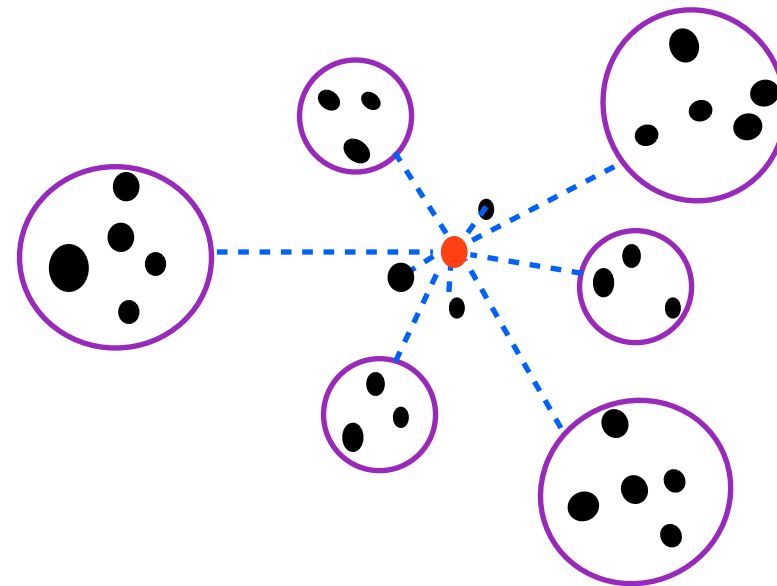


Quantum Chemistry



Multiscale Interactions

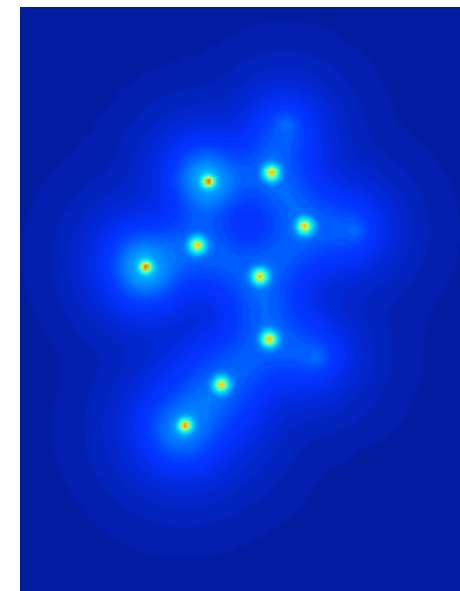
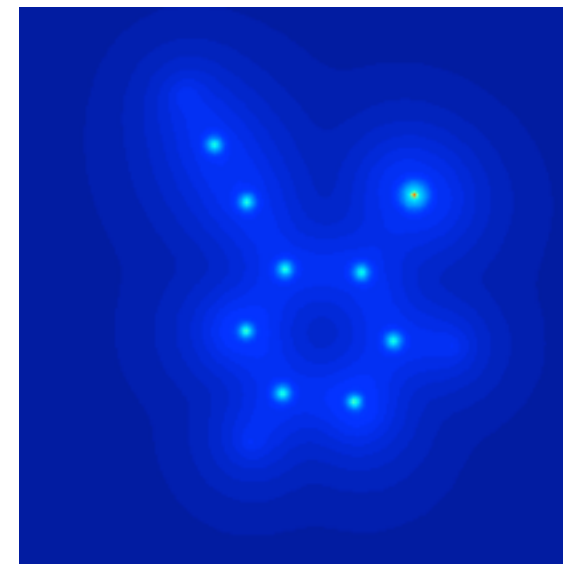
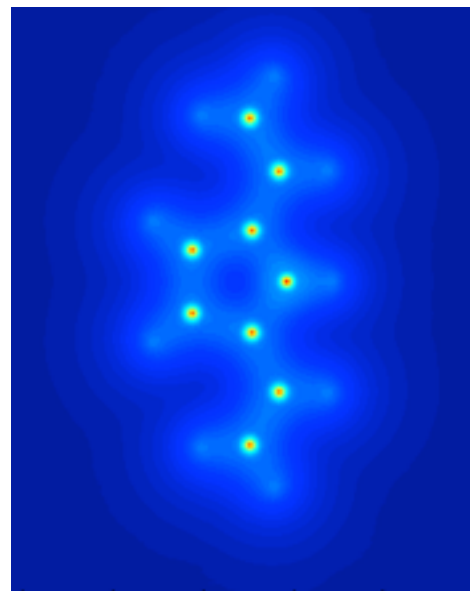
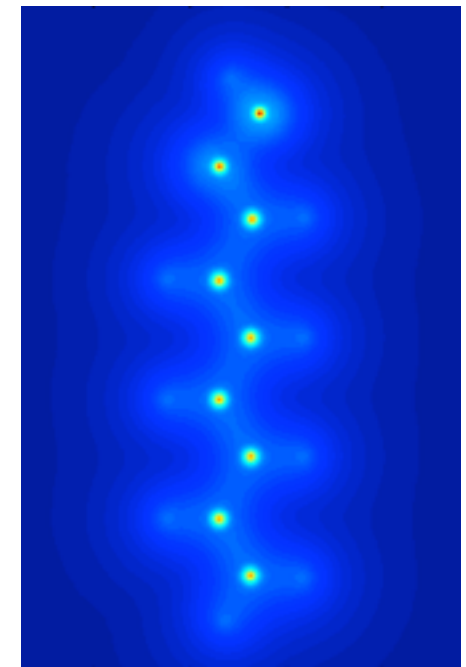
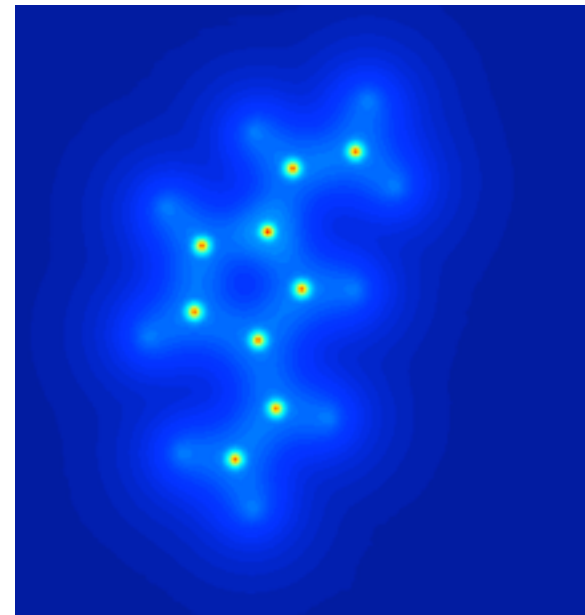
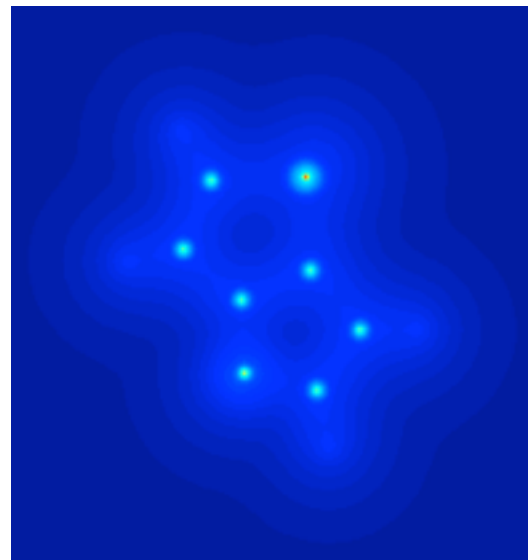
- A system of d particles involves d^2 interactions
- Multiscale separation into $O(\log^2 d)$ interactions



Quantum Chemistry

Electronic density $\rho_x(u)$: computed by solving Schrodinger

Organic molecules
with
Hydrogene, Carbon
Nitrogen, Oxygen
Sulfur, Chlorine



Kohn-Sham model:

$$E(\rho) = T(\rho) + \int \rho(u) V(u) + \frac{1}{2} \int \frac{\rho(u)\rho(v)}{|u-v|} dudv + E_{xc}(\rho)$$

\downarrow
 Molecular
energy

\downarrow
 Kinetic
energy

\downarrow
 electron-nuclei
attraction

\downarrow
 electron-electron
Coulomb repulsion

\downarrow
 Exchange
correlat. energy

At equilibrium:

$$f(x) = E(\rho_x) = \min_{\rho} E(\rho)$$

- $f(x)$ is invariant to isometries and is deformation stable

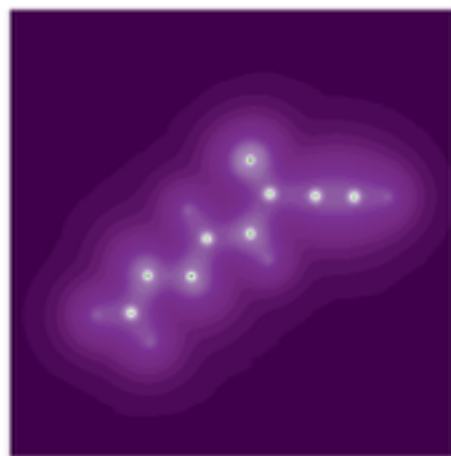
Atomization Density

- We do not know the electronic density ρ_x at equilibrium.

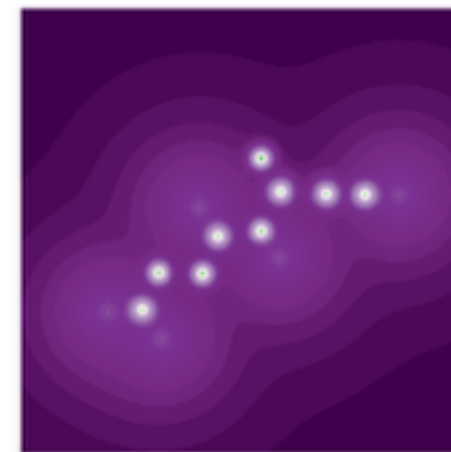
approximated by the sum of the densities of all atoms:

$$\tilde{\rho}_x(u) = \sum_{k=1}^d \rho_{z_k}(u - r_k)$$

Electronic density $\rho_x(u)$



Approximate density $\tilde{\rho}_x(u)$



- Sparse regression computed over a representation invariant to action of isometries in \mathbb{R}^3 :

$$\Phi x = \{\phi_n(\tilde{\rho}_x)\}_n : \left| \begin{array}{l} \text{Fourier modulus coefficients and squared} \\ \text{or} \\ \text{scattering coefficients and squared} \end{array} \right.$$

Partial Least Square regression on the training set:

$$f_M(x) = \sum_{k=1}^M w_k \phi_{n_k}(\tilde{\rho}_x)$$

M: number of variables

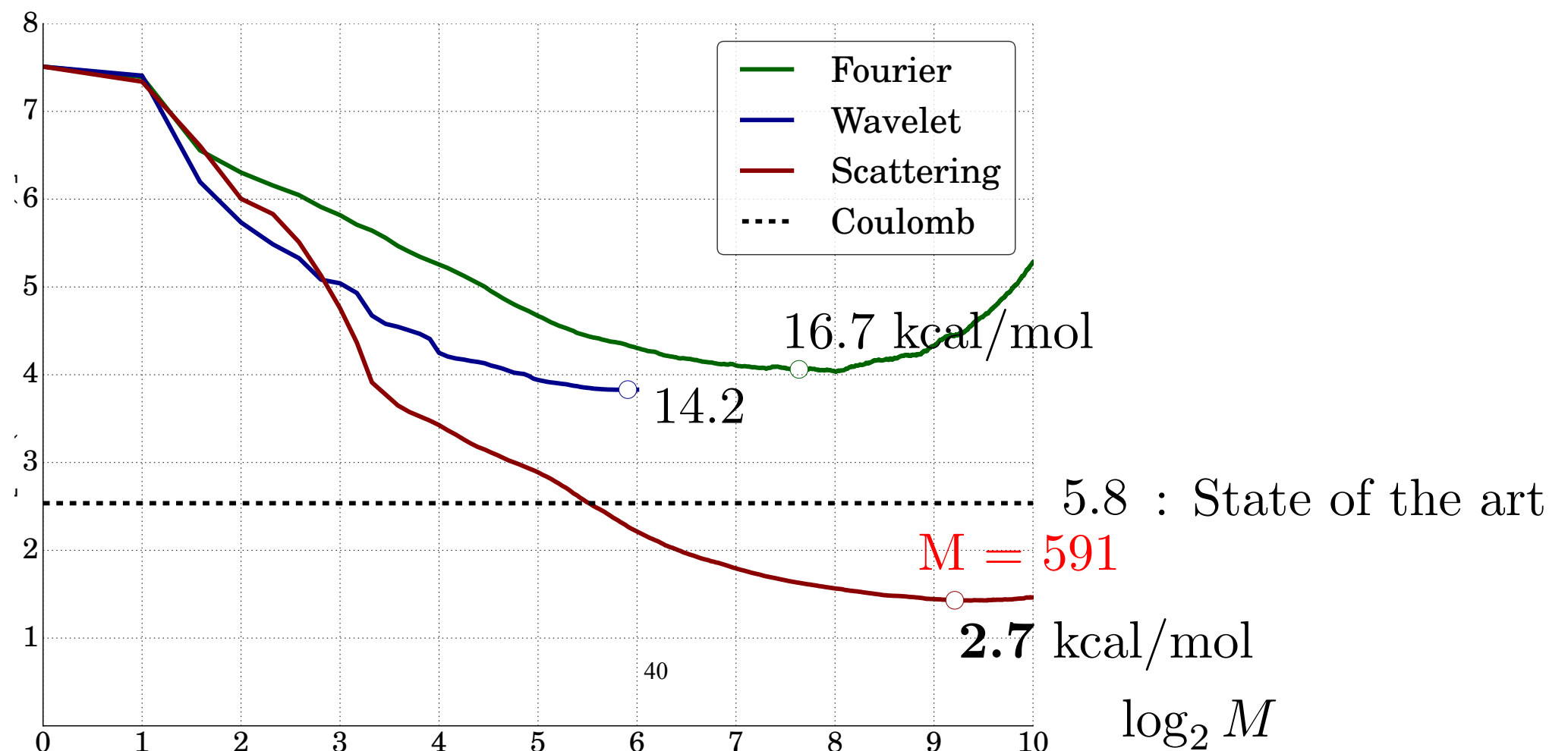
Scattering Regression

Data basis $\{x_i, f(x_i) = E(\rho_{x_i})\}_{i \leq N}$ of 4357 planar molecules

$$\text{Regression: } f_M(x) = \sum_{m=1}^M w_m \phi_{k_m}(\tilde{\rho}_x)$$

Testing error

$$2^{-1} \log_2 \mathbb{E} |f_M(x) - y(x)|^2$$



Conclusion

- A major challenge of data analysis is to find Euclidean embeddings of metrics \Leftrightarrow build Gaussian models
- Continuity to action of diffeomorphisms \Rightarrow wavelets
- Known geometry \Rightarrow no need to learn.
Unknown geometry: learn wavelets on appropriate groups.
- Can learn physics from prior on geometry and invariants.
- Applications to images, audio and natural languages

www.di.ens.fr/data/scattering