

Colloquium du CERMICS



Mathematical Mysteries of Deep Neural Networks

Stéphane Mallat (Collège de France et École Normale Supérieure)

20 septembre 2019

Mathematical Mysteries of Deep Neural Networks



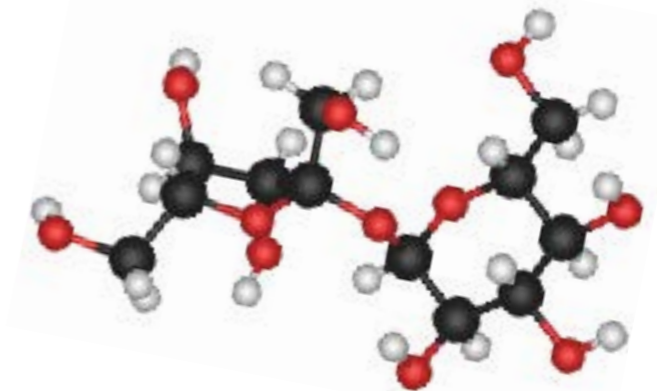
*Tomas Anglès, Roberto Leonarduzzi,
Stéphane Mallat, Louis Thiry,
John Zarka, Sixin Zhang*

Collège de France
École Normale Supérieure

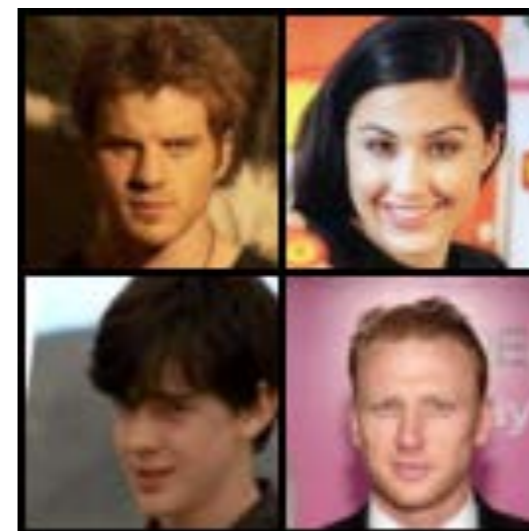
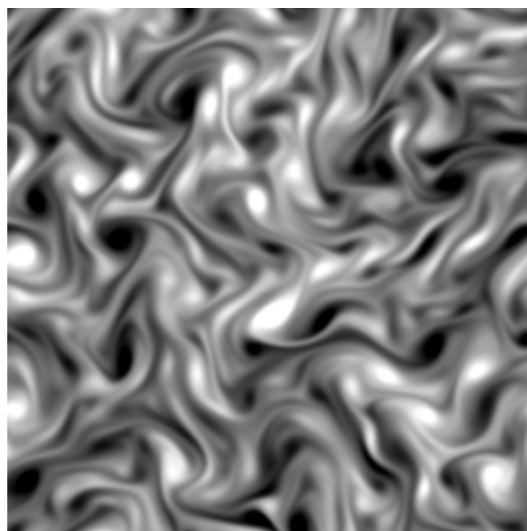
High-Dimensional Approximations

What regularity properties lead to low-dimensional approximations of $f(x)$ for a high-dimensional $x \in \mathbb{R}^d$ in physics and machine learning ?

- $f(x)$: class of an image x having $d = 10^6$ pixels or energy of a physical system in a state $x \in \mathbb{R}^d$

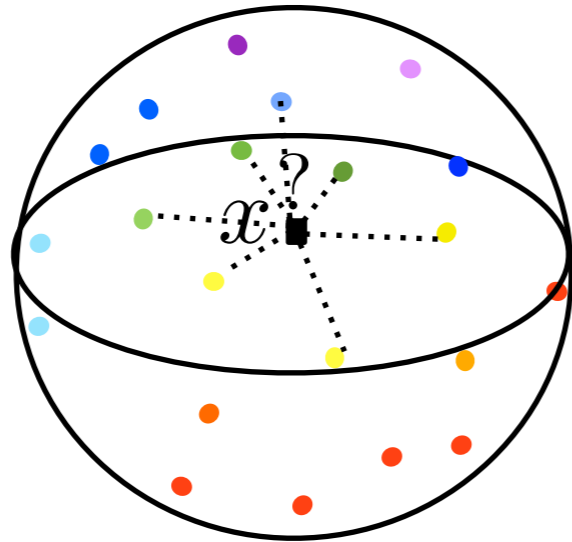


- $f(x) = p(x)$ a probability density.



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:

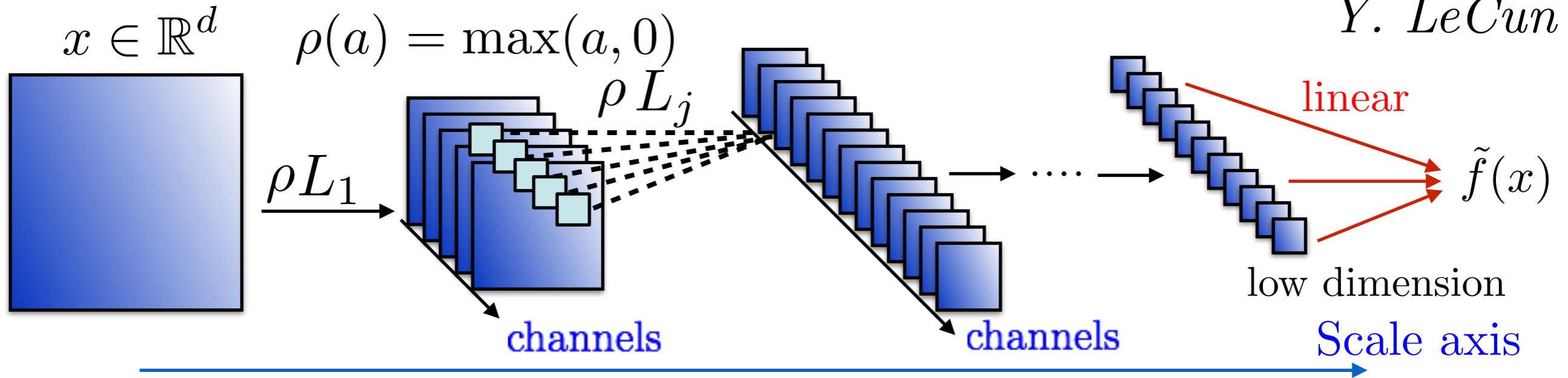


- Need $n \geq \epsilon^{-d}$ points to cover $[0, 1]^d$ at a Euclidean distance ϵ
Problem: $\|x - x_i\|$ is always large
- To estimate $f(x)$ when x is in a high-dimensional Ω requires *strong regularity* of f in Ω : what regularity ?

Deep Convolutional Network

- Deep convolutional neural network to predict $y = f(x)$:

Y. LeCun



L_j : spatial convolutions and linear combination of channels

Exceptional results for classification of *images, sounds, language, regressions in physics, signal and image generation...*

but not interpretable.

To create simpler interpretable networks:

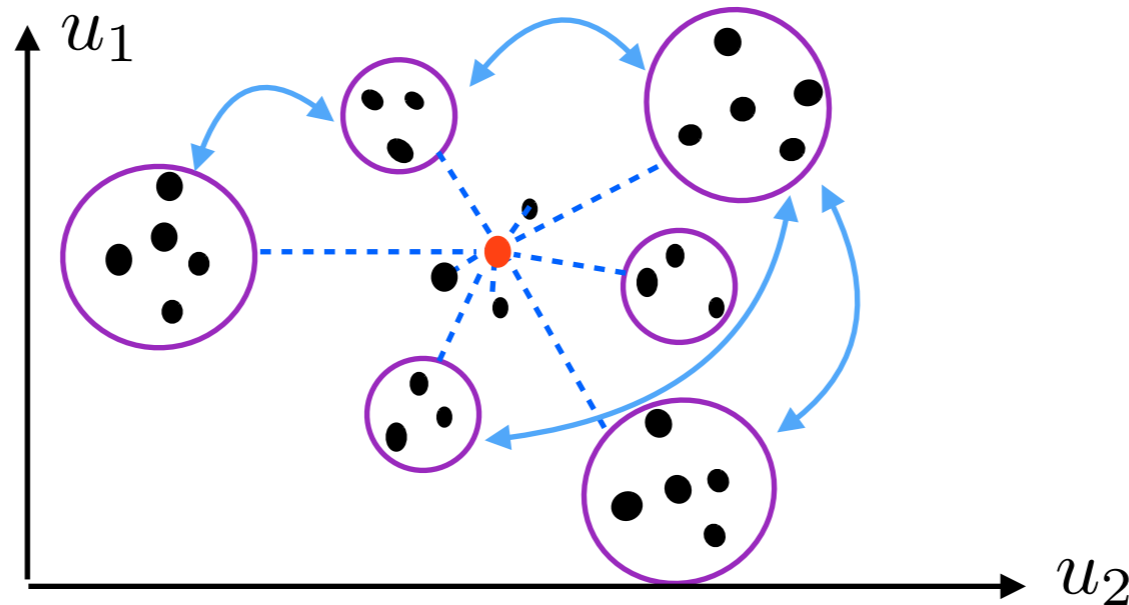
What underlying regularity is captured and how ?

3 ingredients: **Multiscale, Linearize group actions, Sparse**

- **Dimension reduction:**

Interactions de d bodies represented by $x(u)$: particles, pixels...

Interactions
across scales



Multiscale regroupement of interactions of d bodies into interactions of $O(\log d)$ groups.

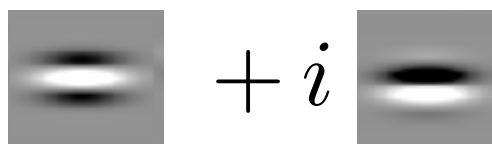
Scale separation \Rightarrow wavelet transforms.

How to capture scale interactions ?

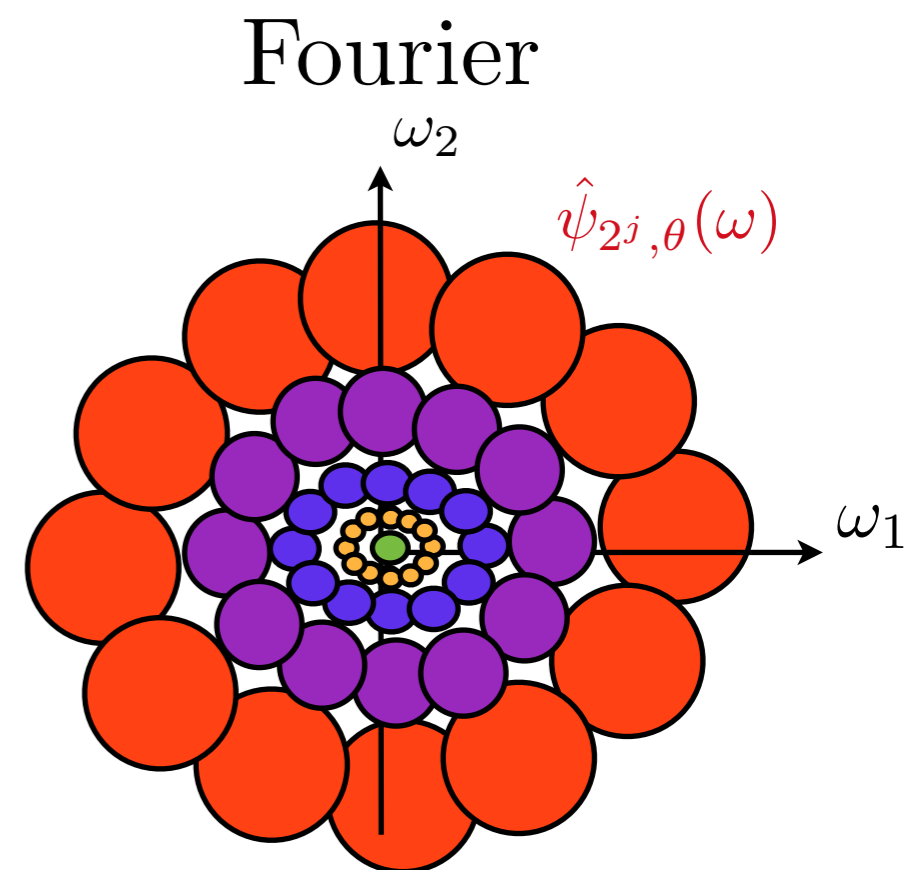
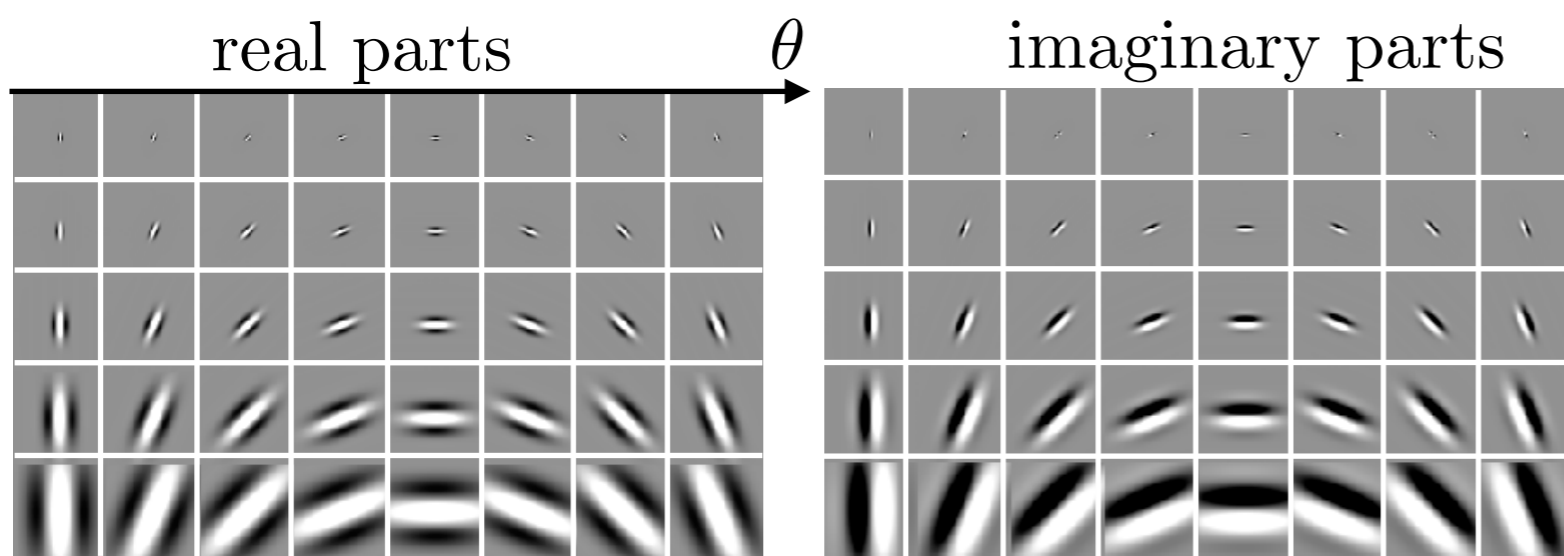
Critical
harmonic analysis
problems since 1970's

- **Scale separation** with wavelets and interactions through phase
- Linear scale interaction models and **invariants** in:
 - Statistical physics for turbulence
 - Quantum chemistry and image classification
- Non-linear scale interactions models with **sparse dictionaries** in:
 - Classification of complex structures as in ImageNet
 - Generation of non-ergodic processes

Scale separation with Wavelets

- Wavelet filter $\psi(u)$: 

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



- Wavelet transform: invertible

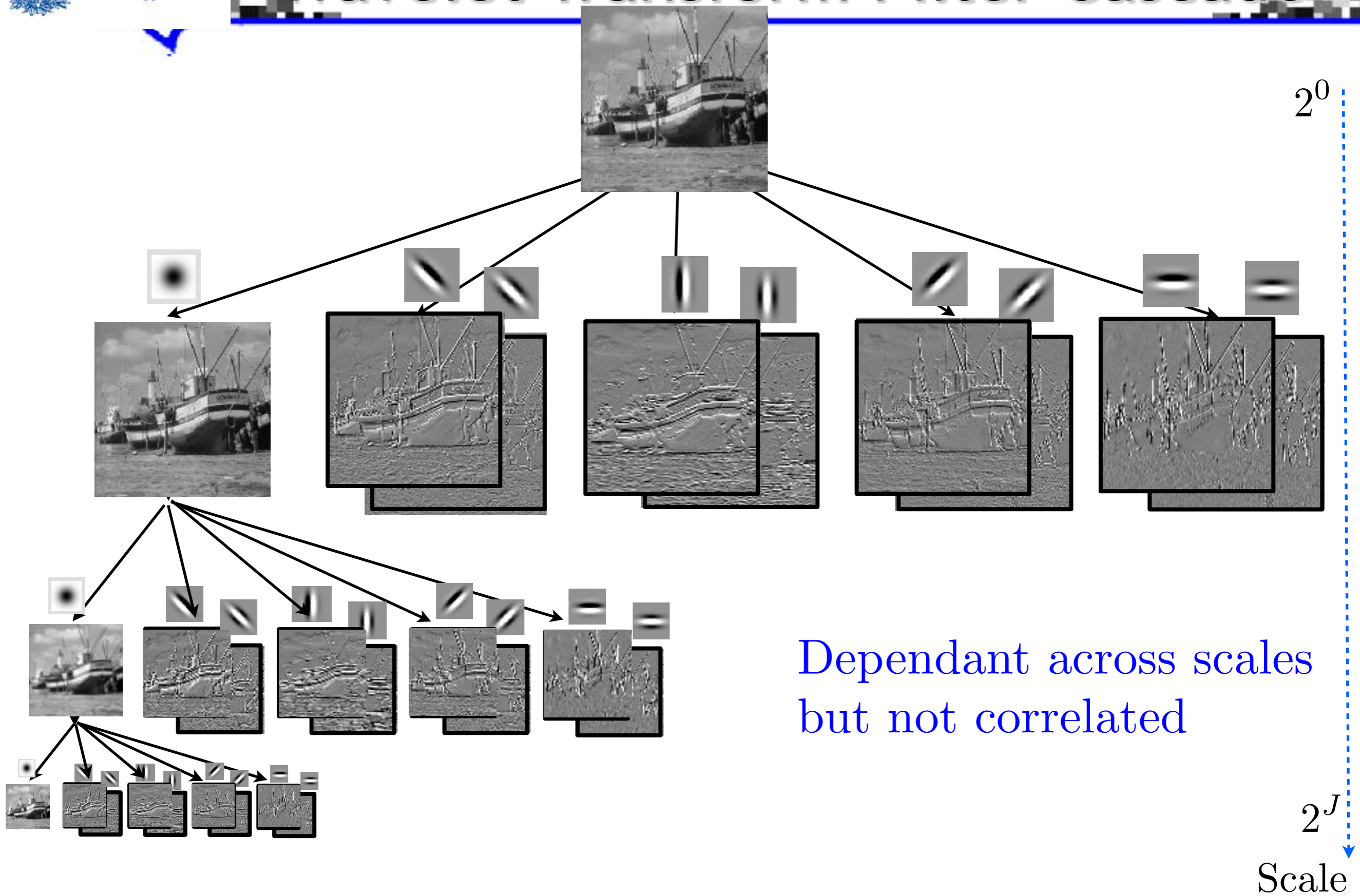
$$Wx = \left(x \star \psi_\lambda \right)_\lambda$$

$$\widehat{x \star \psi_\lambda}(\omega) = \hat{x}(\omega) \hat{\psi}_\lambda(\omega)$$

- Zero-mean and no correlations across scales: **problem!**

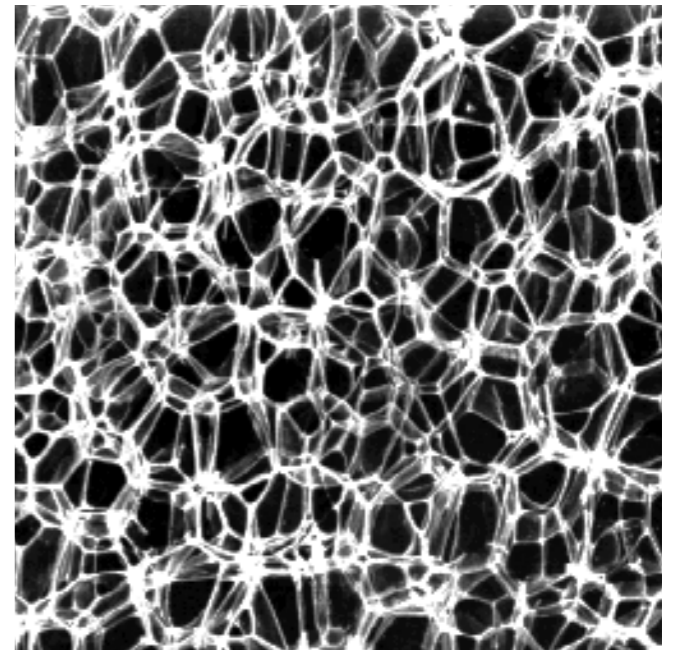
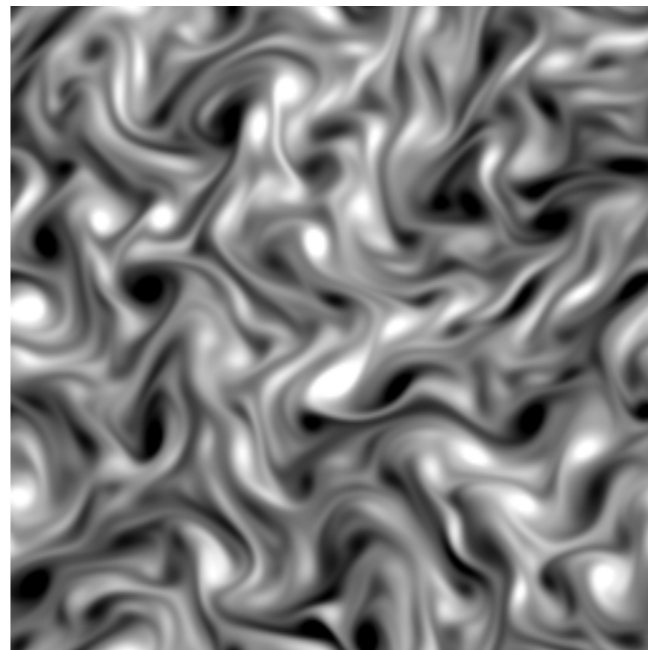
$$\sum_u x \star \psi_\lambda(u) x \star \psi_{\lambda'}^*(u) = \sum_\omega |\hat{x}(\omega)|^2 \psi_\lambda(\omega) \psi_\lambda(\omega)^* \approx 0 \text{ if } \lambda \neq \lambda'$$

Wavelet Transform Filter Cascade



What stochastic models
for turbulence ?

$$d = 6 \cdot 10^4$$



Prior: stationary $\Leftrightarrow p(x)$ is **invariant to translations**.

Maximum entropy distribution \tilde{p} conditioned by M moments

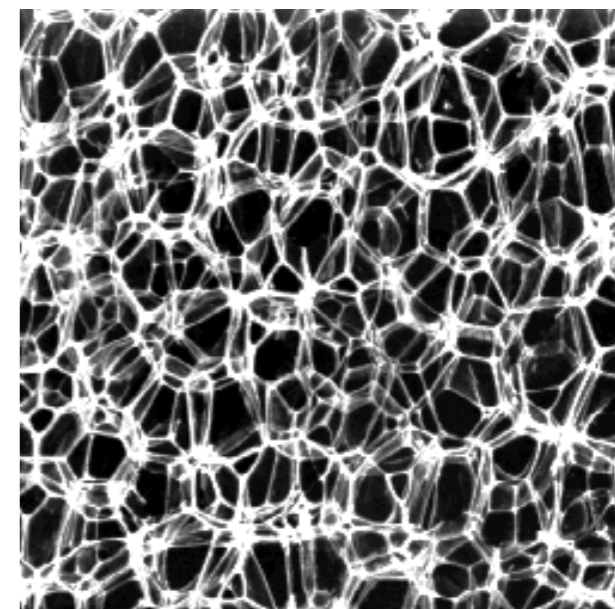
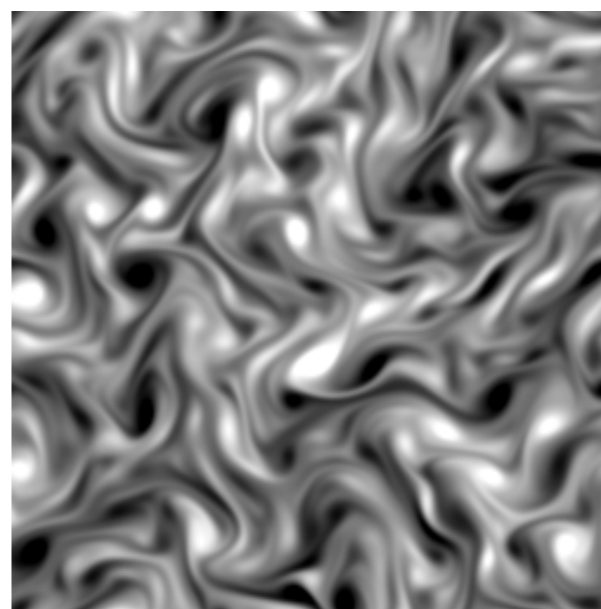
$$\mathbb{E}(\phi_m(x)) = \mu_m \quad \Rightarrow \quad \tilde{p}(x) = \mathcal{Z}^{-1} \exp \left(- \sum_{m=1}^M \beta_m \phi_m(x) \right)$$

With $M = d$ second order moments:

$$\phi_m(x) = \sum_u x(u)x(u - m) \quad \Rightarrow \quad \tilde{p}(x) \text{ is a Gaussian distribution}$$

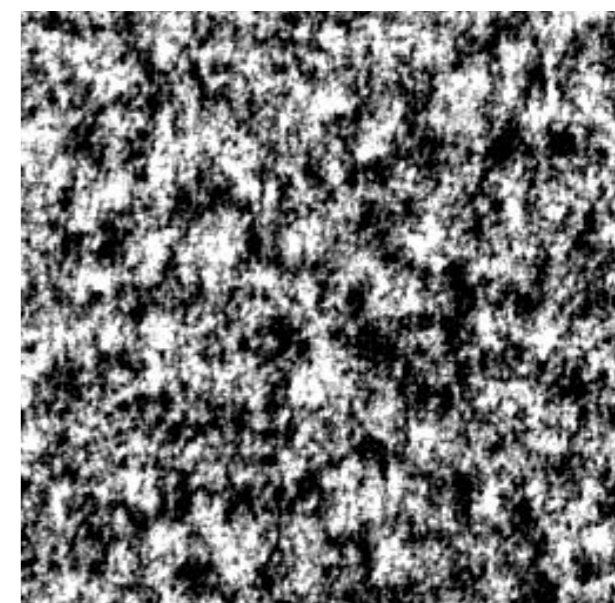
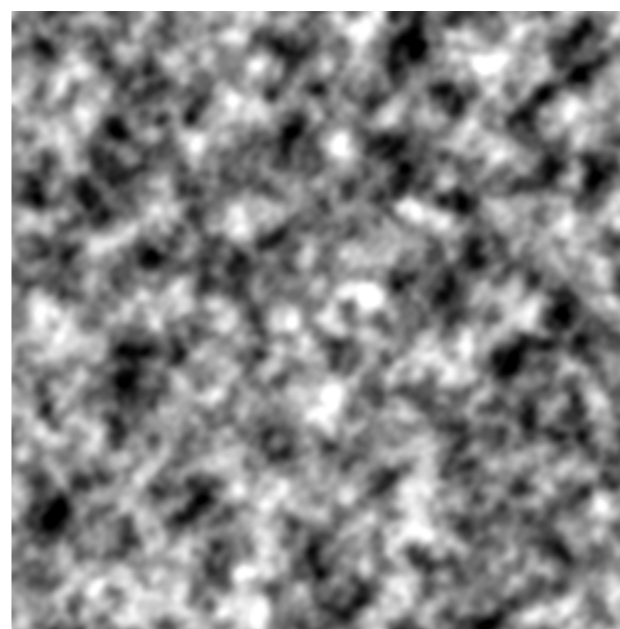
What stochastic models
for turbulence ?

$$d = 6 \cdot 10^4$$



$\tilde{p}(x)$ is a Gaussian distribution

\tilde{x}



No correlation is captured across scales and frequencies.

Random phases.

How to capture non-Gaussianity and long range interactions ?

Failure of high order moments. Deep net generations look better.

Rectifiers act on Phase

- Real wavelets of phase α : $\psi_{\alpha, \lambda} = \text{Real}(e^{-i\alpha} \psi_{\lambda})$

Rectifier: $\rho(a) = \max(a, 0)$

$$Ux(u, \alpha, \lambda) = \rho(x \star \text{Real}(e^{i\alpha} \psi_{\lambda})) = \rho(\text{Real}(e^{i\alpha} x \star \psi_{\lambda}))$$

$$x \star \psi_{\lambda} = |x \star \psi_{\lambda}| e^{i\varphi(x \star \psi_{\lambda})}$$

Homogeneous: $\rho(\alpha a) = \alpha \rho(a)$ if $\alpha > 0$

$$Ux(u, \alpha, \lambda) = |x \star \psi_{\lambda}| \rho(\cos(\alpha + \varphi(x \star \psi_{\lambda})))$$

A Relu computes phase harmonics:

Theorem : Fourier transform along the phase α :

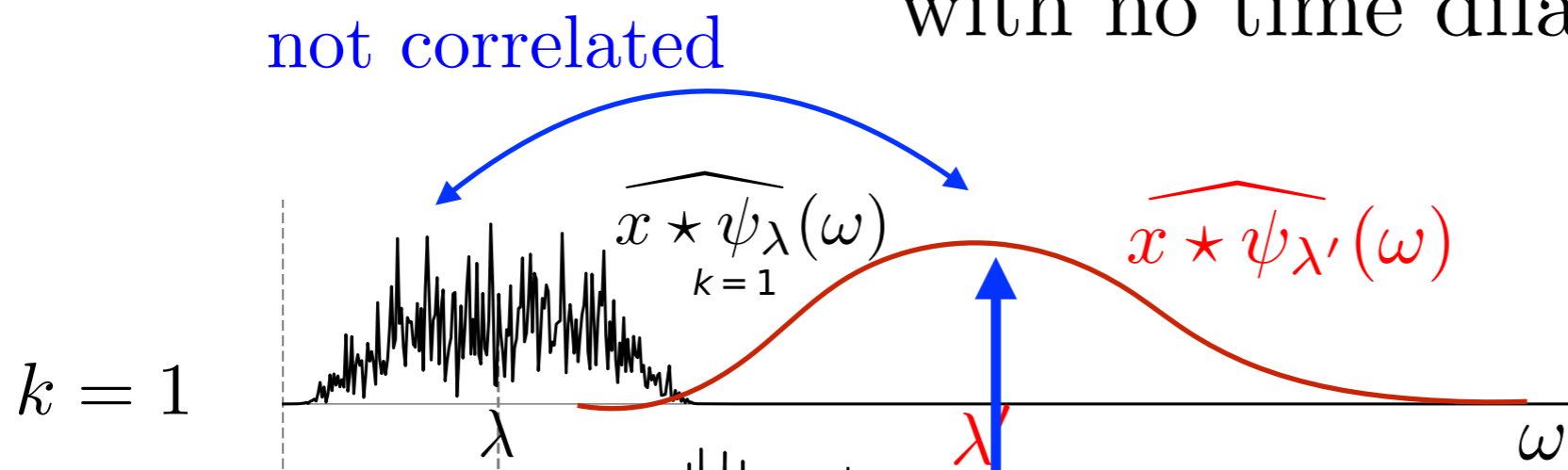
$$\widehat{U}x(u, k, \lambda) = \widehat{\gamma}(k) |x \star \psi_{\lambda}(u)| e^{ik \varphi(x \star \psi_{\lambda}(u))}$$

with $\gamma(\alpha) = \rho(\cos \alpha)$ for any homogeneous non-linearity ρ .

Frequency Transpositions

Phase harmonics: $|x \star \psi_\lambda(u)| e^{i k \varphi(x \star \psi_\lambda(u))}$

Performs a non-linear frequency dilation / transposition
with no time dilation



Phase
Harmonics

Correlated if $k\lambda \approx \lambda'$

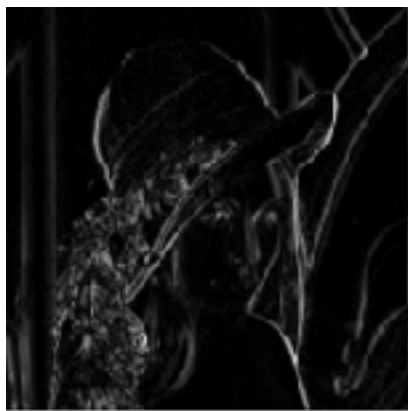
Rectified Wavelet Coefficients

Real wavelets: $\psi_{\alpha,\lambda} = \text{Real}(e^{-i\alpha} \psi_{\lambda})$ and $\rho(a) = \max(a, 0)$

$$\rho(x \star \psi_{\alpha,\lambda}) \xrightarrow[\text{along } \alpha]{\text{Fourier transform}} c_k |x \star \psi_{\lambda}| e^{ik \varphi(x \star \psi_{\lambda})}$$

Harmonics

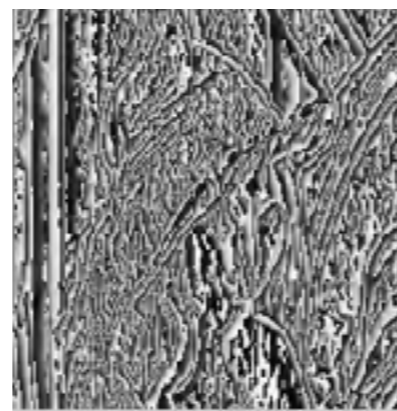
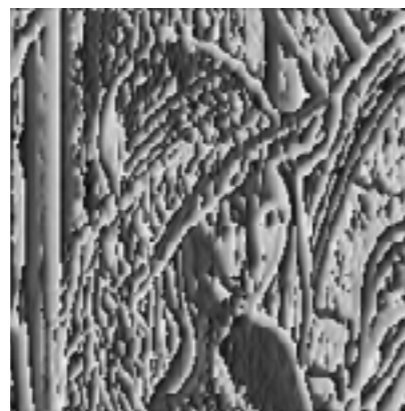
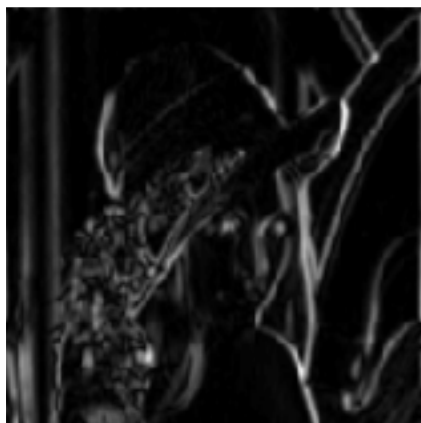
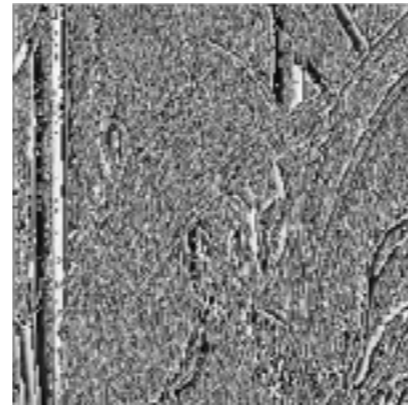
$|x \star \psi_{\lambda}|$



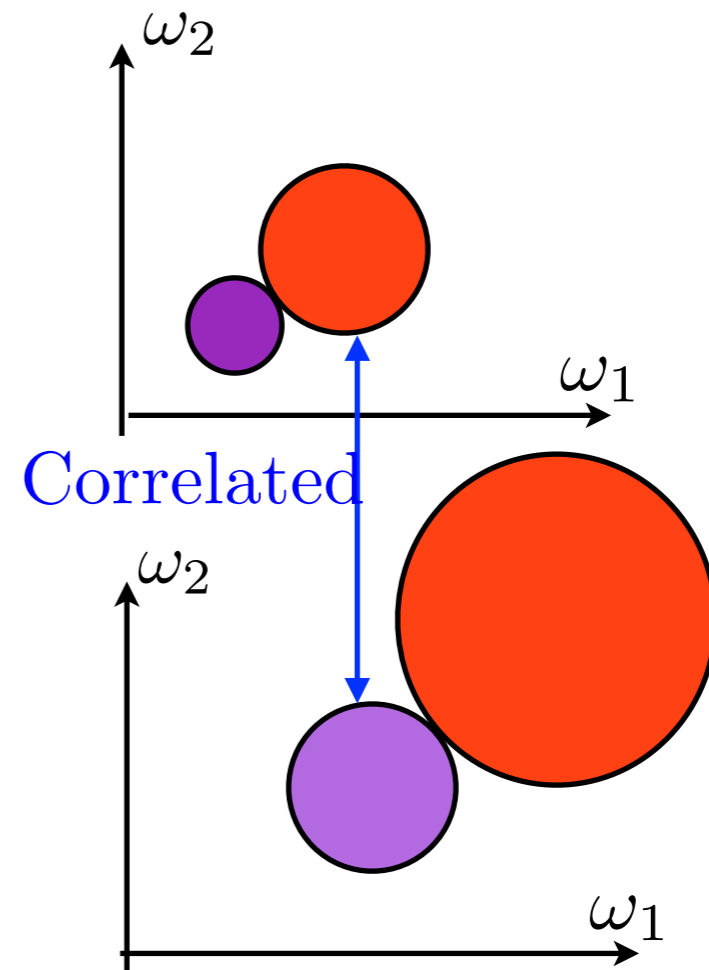
$\varphi(x \star \psi_{\lambda})$



$k \varphi(x \star \psi_{\lambda})$

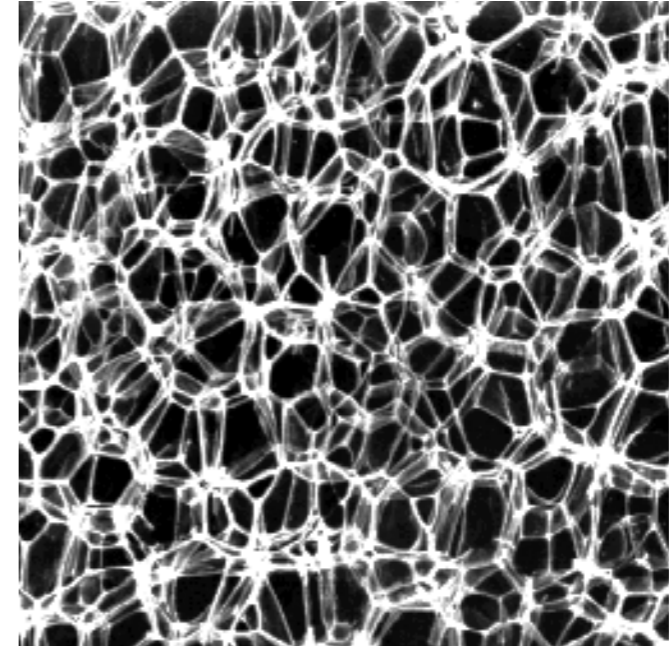
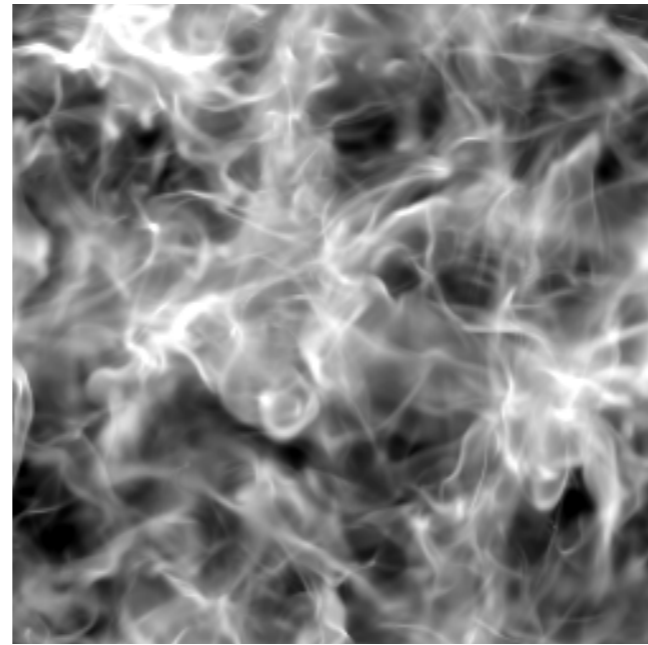
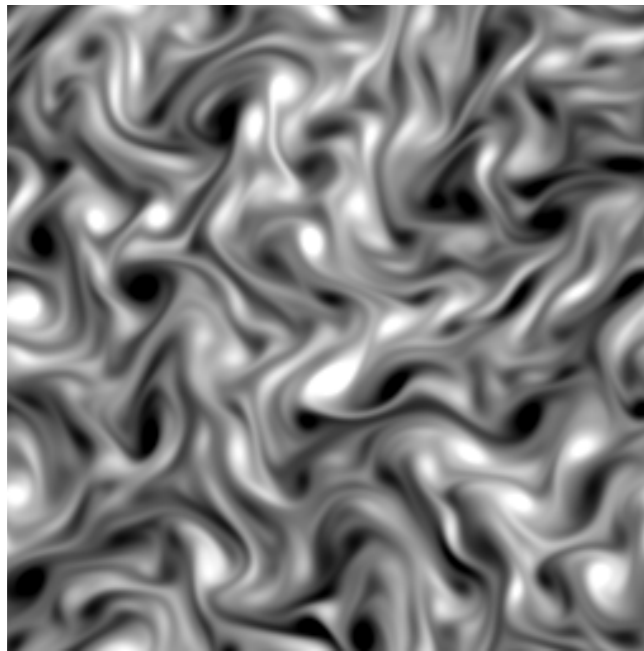


$k = 2$



Phase harmonics:
Frequency transpositions

x



Maximum entropy distribution conditioned by

$M = O(\log^2 d)$ wavelet harmonic correlations $\mathbb{E}(\phi_m(x))$

$$\phi_m(x) = \sum_u |x \star \psi_\lambda(u)| e^{ik\varphi(x \star \psi_\lambda(u))} |x \star \psi_{\lambda'}(u)| e^{-ik'\varphi(x \star \psi_{\lambda'}(u))}$$

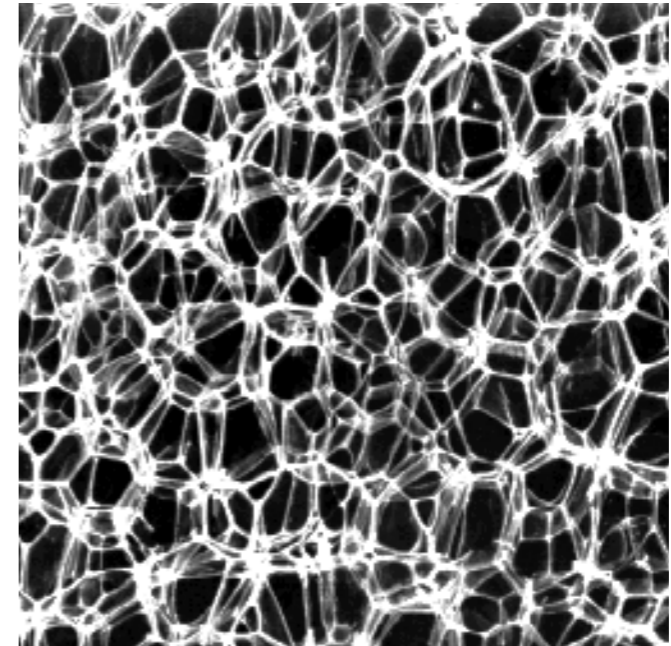
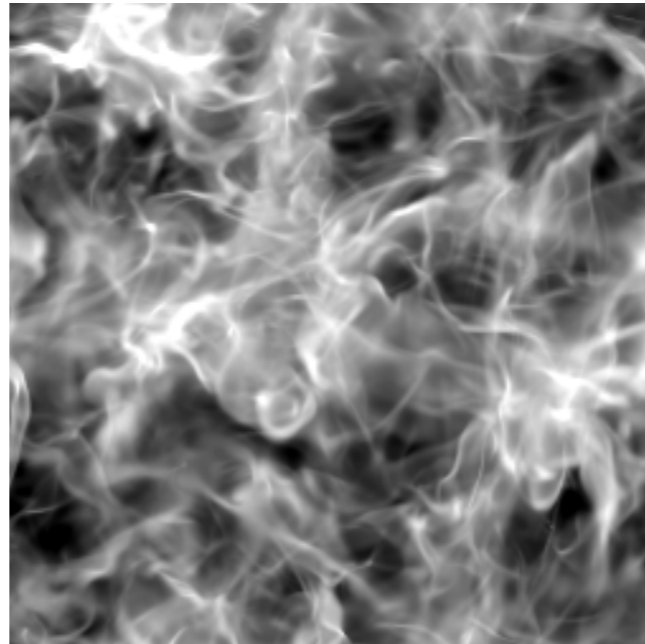
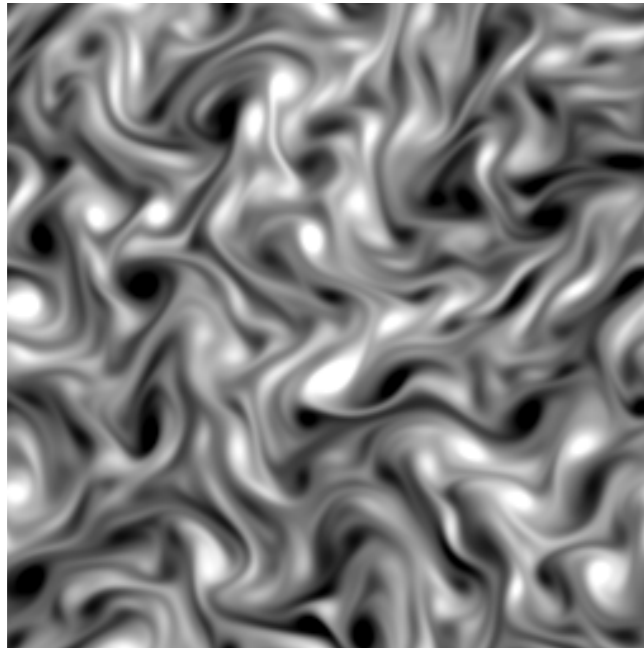
$$\tilde{p}(x) = \mathcal{Z}^{-1} \exp \left(- \sum_{m=1}^M \beta_m \phi_m(x) \right)$$

Ergodic Stationary Processes

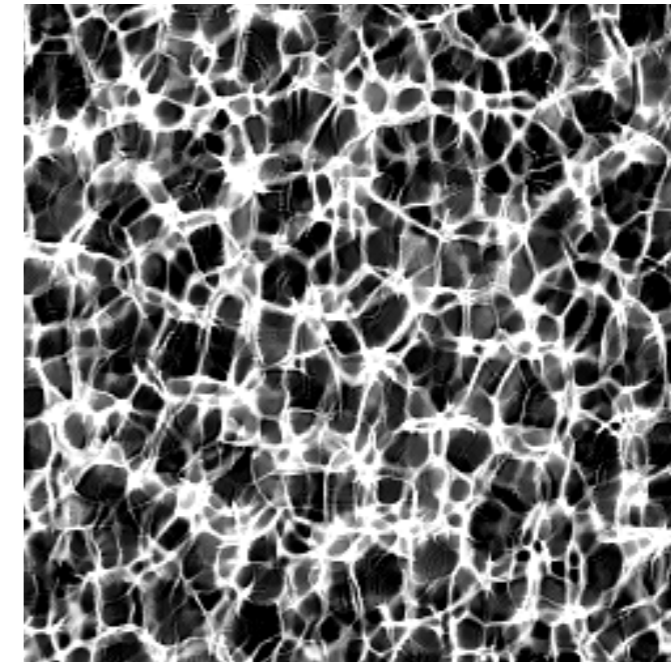
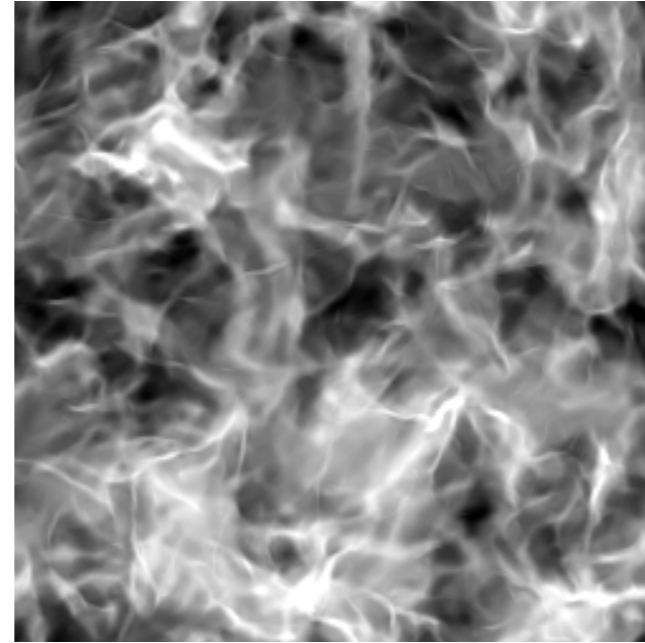
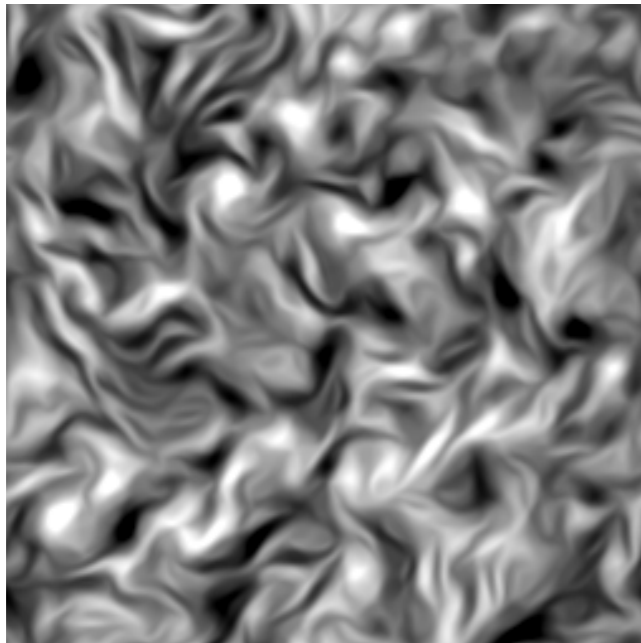
S. Zhang, J. Bruna, E. Allys, F. Levrier, F. Boulanger

$$d = 6 \cdot 10^4$$

x



\tilde{x}



$M = 3 \cdot 10^3$
number
of moments

Phase coherence is restored

How much physics are these models capturing ?

What about non-ergodic processes ?

Classification: invariance by translation by spatial averaging

$$\begin{pmatrix} x \star \phi_{2^J}(2^J n) \\ \rho(x \star \psi_{\alpha, \lambda}) \star \phi_J(2^J n) \end{pmatrix}_{\alpha, \lambda}$$

Recover the information loss with a second layer:

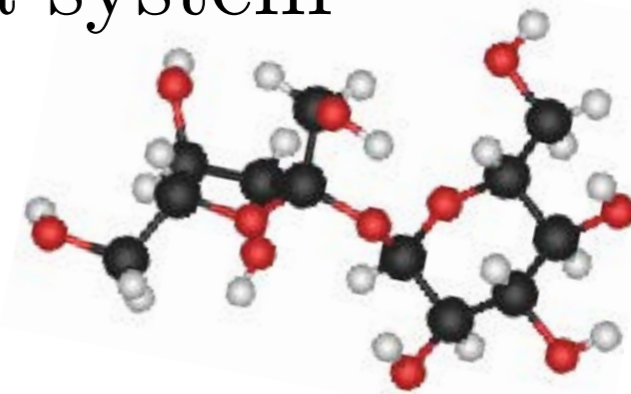
$$S_J x = U x \star \phi_J = \begin{pmatrix} x \star \phi_{2^J}(2^J n) \\ \rho(x \star \psi_{\alpha, \lambda}) \star \phi_J(2^J n) \\ \rho(\rho(x \star \psi_{\alpha, \lambda}) \star \psi_{j', \alpha'}) \star \phi_J(2^J n) \end{pmatrix}_{\alpha, \lambda, \alpha', \lambda'}$$

- Linearize small deformations

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

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

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



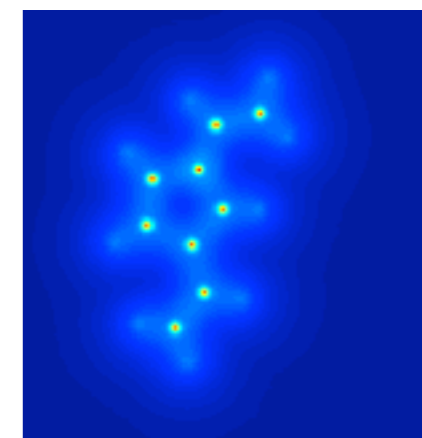
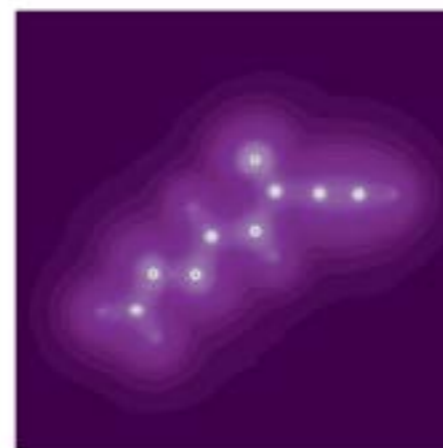
Symmetries:

$f(x)$ is invariant to translations and rotations,

multiscale interactions: chemical bounds, Van der Waal forces...

The energy depends upon the electronic density (Kohn-Sham)

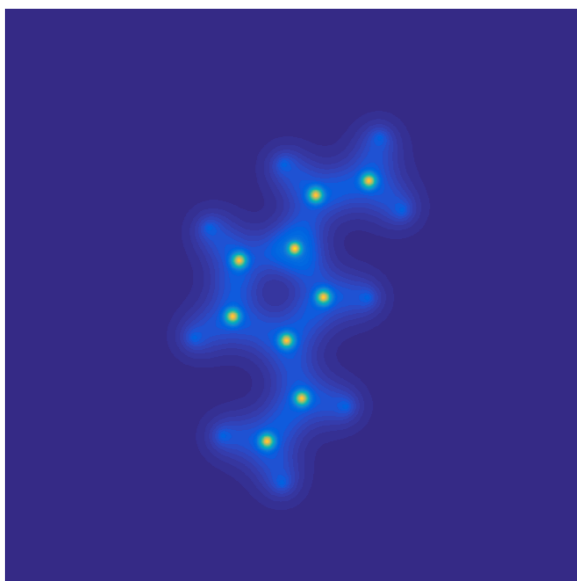
Ground state
electronic density
computed with Schroedinger



- We do not know the electronic density at equilibrium.
- The molecular state $\{z_k, r_k\}_{k \leq d}$ is represented by Diracs located at r_k weighted by charges z_k :

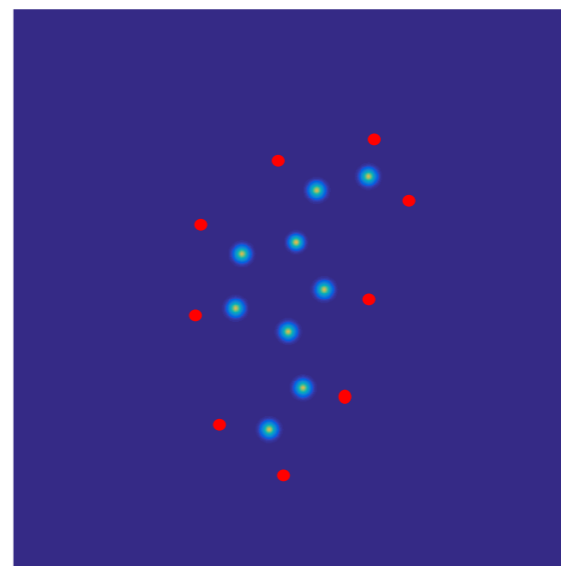
$$x(u) = \sum_{k=1}^d z_k \delta(u - r_k)$$

Electronic density



Dirac density

$x(u)$



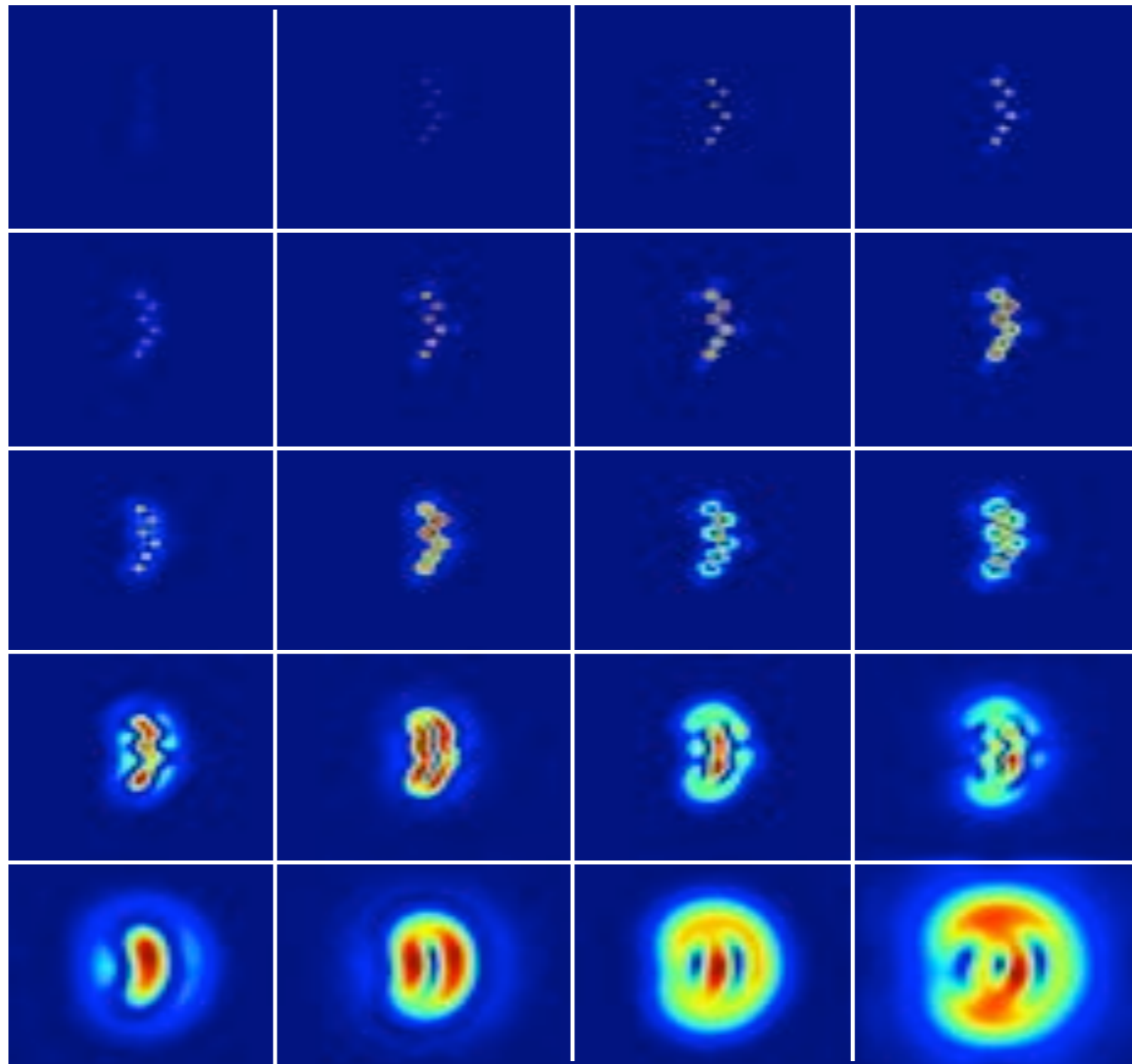
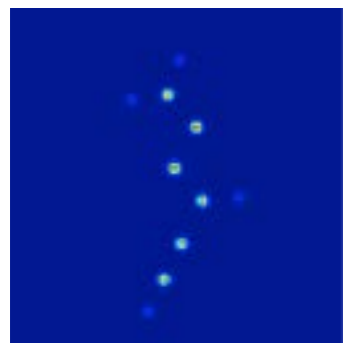
$$x = \sum_k z_k \delta(u - r_k) \Rightarrow \rho(x \star \psi_{2^j, \ell}(u)) = \rho\left(\sum_k z_k \psi_{2^j, \ell}(u - r_k)\right)$$

$\ell = 0$ $\ell = 1$ $\ell = 2$ $\ell = 3$

$j = 0$

$j = 1$

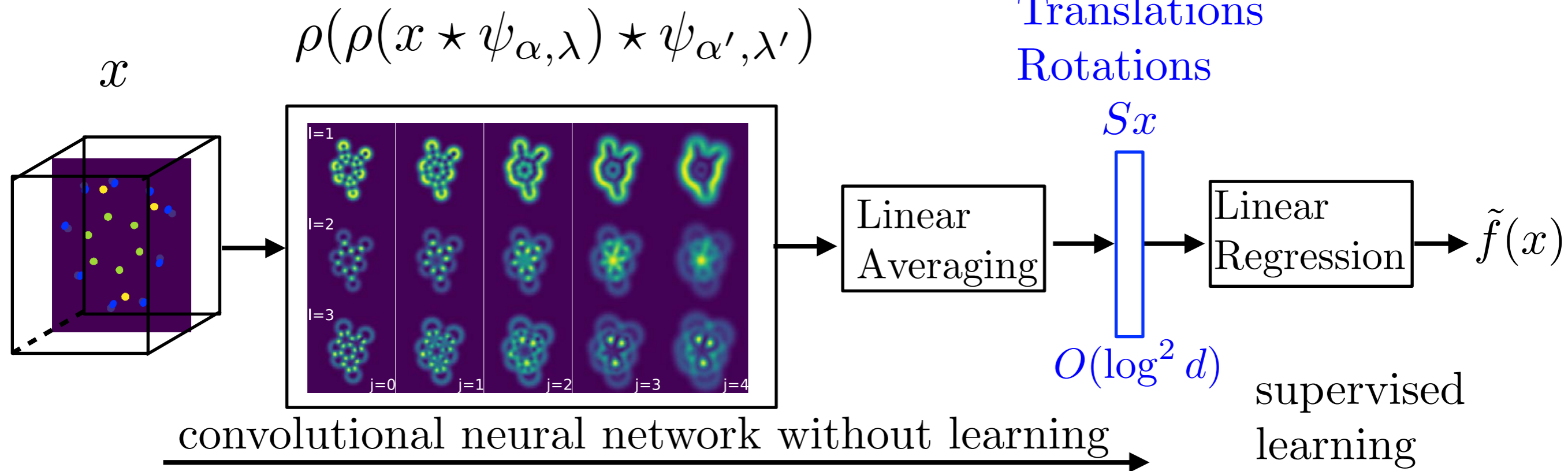
$j = 5$



Scattering Energy Regression

M. Eickenberg, G. Excarchakis M. Hirn, N. Poilvert, L. Thiry

Invariants to
Translations
Rotations



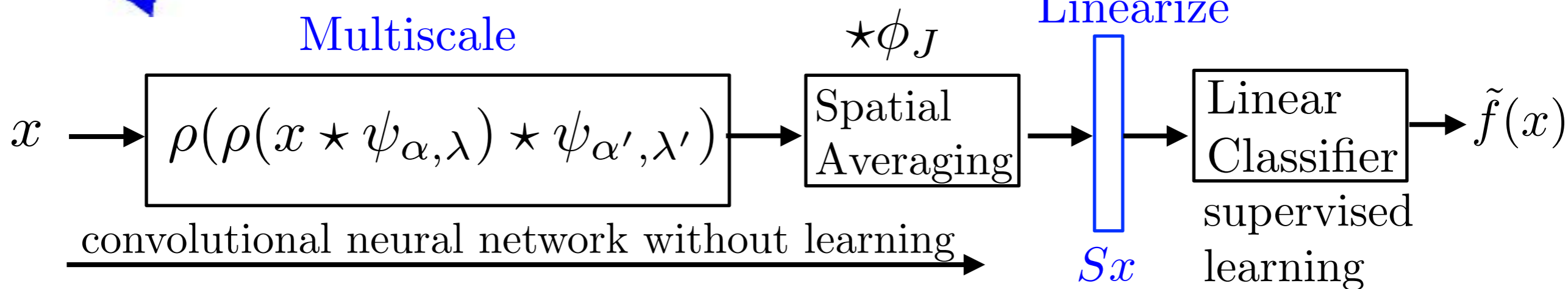
QM9: Data basis of 130.000 organic molecules with C, H, O, N, F with DFT atomisation energies

Regression error ~ 0.5 kcal/mol \sim Deep Nets.

But small molecules with at most 29 atoms and 9 heavy ones

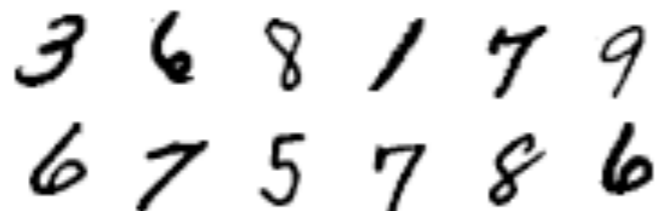


Multiscale



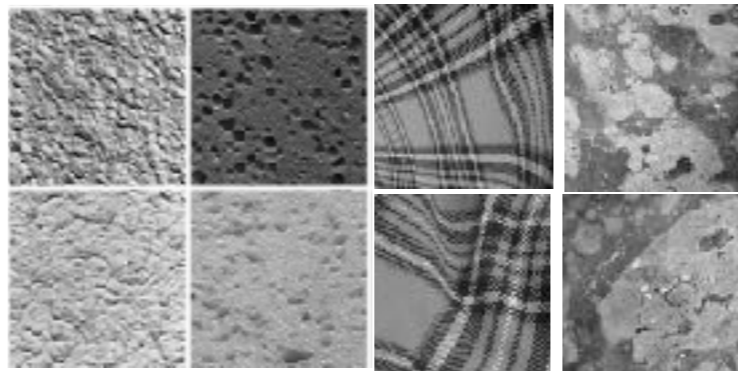
Errors: Scattering | Deep Nets.

Digits
10 classes



0.5%

Textures
60 classes



0.5%

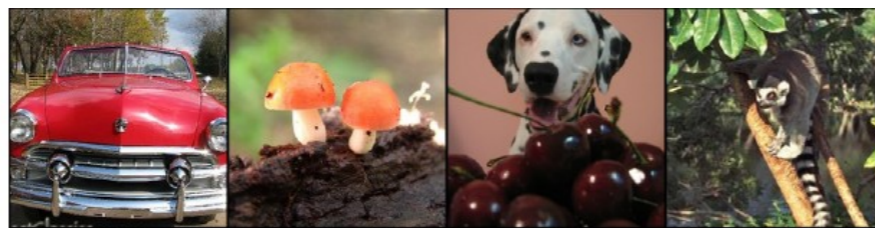
ImageNet
 10^3 classes



mite container ship motor scooter leopard

60%

AlexNet
20% : 2012



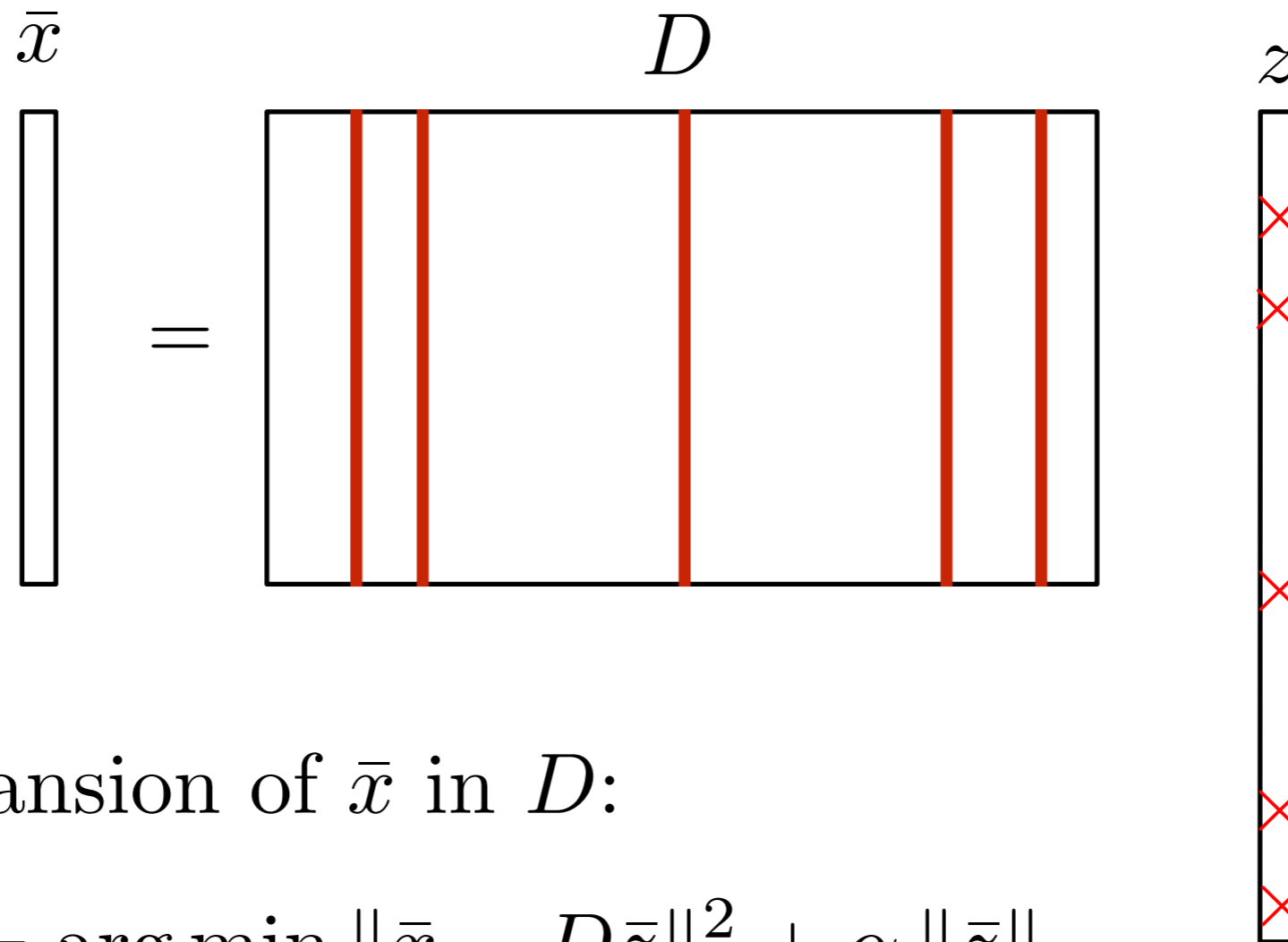
grille mushroom cherry Madagascar cat

What is learned ?

Sparse Dictionary Representation

- Need to learn "sparse informative patterns"

Pattern representations with sparse dictionary expansions:



Sparse l^1 expansion of \bar{x} in D :

$$z = \arg \min_{\bar{z}} \|\bar{x} - D\bar{z}\|_2^2 + \alpha \|\bar{z}\|_1$$

- How to minimise this convex cost ?

Homotopy ISTA Network

- Homotopy algorithms decrease the multiplier α_k :

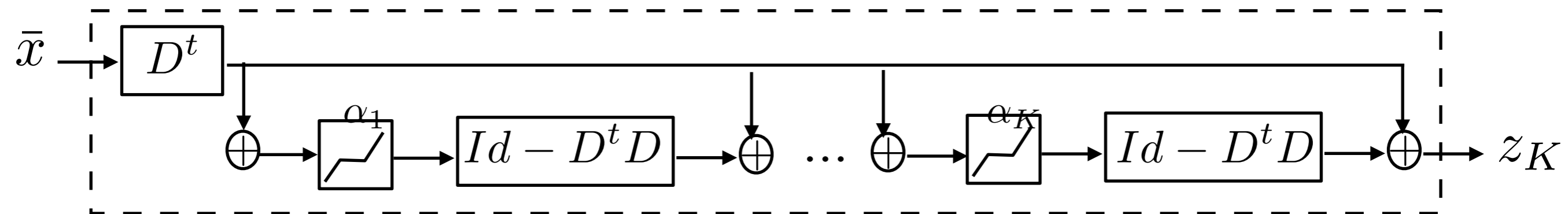
$$z = \arg \min_{\bar{z}} \|\bar{x} - D\bar{z}\|_2^2 + \alpha_k \|z\|_1$$

with an iterated soft-threshold decreasing thresholds:

$$z_{k+1} = T_{\alpha_k} (D^t \bar{x} + (I - D^t D) z_k) \xrightarrow[k \rightarrow \infty]{\alpha_k \sim \gamma^{-k}} z$$

where $T_\alpha(a) = \text{sign}(a) \max(|a| - \alpha, 0)$ is a soft-thresholding.

- Implemented with a convolutional network of depth K :

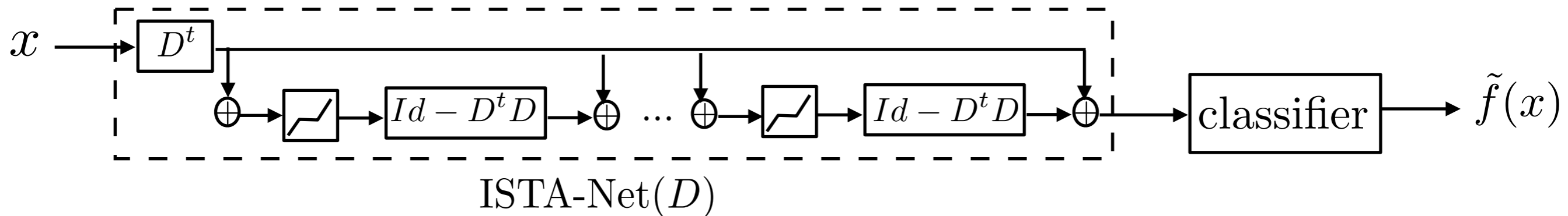


ISTA-Net(D)

with a convolutional dictionary D

- Deep network with sparse coding and classification:

l^1 sparse coding in D



- Optimize the dictionary D and the classifier to minimize the classification loss over a supervised data basis $\{x_i, y_i\}_i$:

$$\text{Loss}(D) = \sum_i \text{loss}(y_i, \tilde{f}(x_i))$$

- Stochastic gradient descent

ImageNet Classification

J. Zarka, L. Thiry, T. Angles

ImageNet
 10^3 classes



Learned sparse coding



Multiscale Scattering

Linearize

Learned sparse coding

18% error



Interpretable network : patterns stored in D

AlexNet: 20% error

Why such an error reduction ? Linearize classes in separate spaces

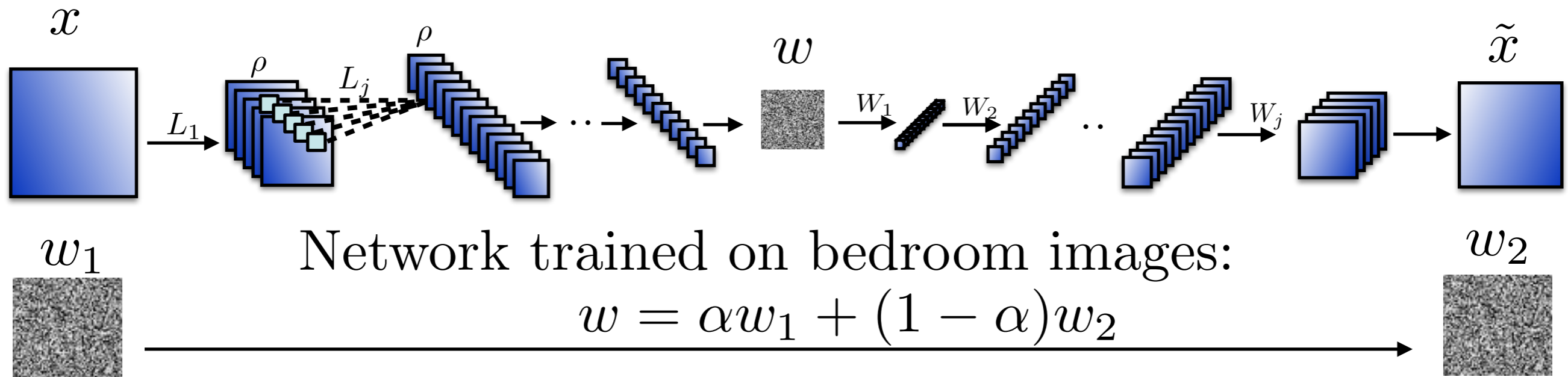
Non Ergodic Processes

autoencoder: trained on n examples $\{x_i\}_{i \leq n}$

Encoder

Gaussian white

Decoder



Network trained on bedroom images:

$$w = \alpha w_1 + (1 - \alpha) w_2$$



Network trained on faces of celebrities:



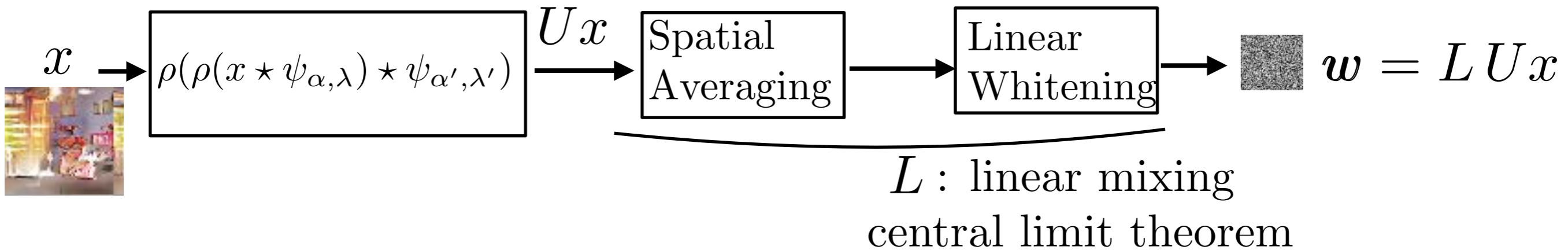
Tomas Angles

Encoder

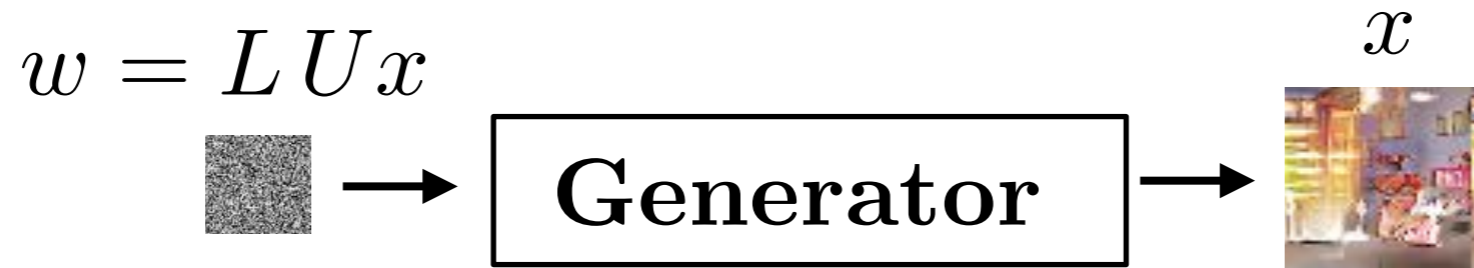
Multiscale

Nearly
Gaussian

Low-dimensional
Gauss. white noise



Inversion:



U has a linear inverse $U^{-1} : \rho(a) + \rho(-a) = a$

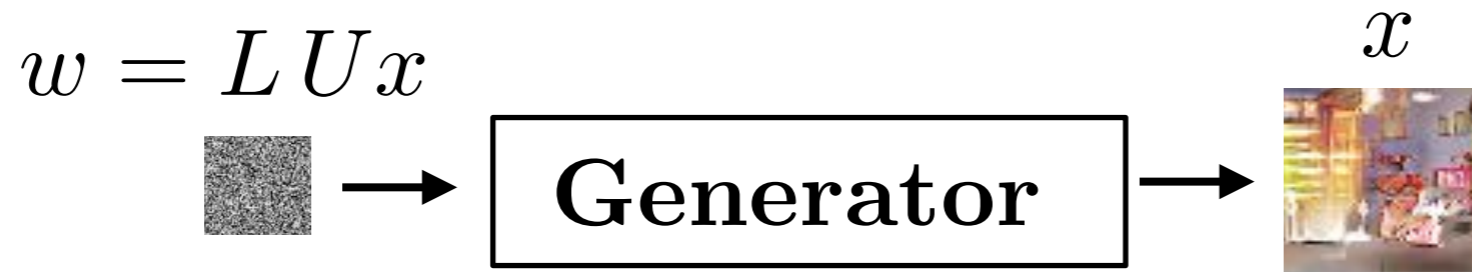
L is non-invertible linear projector

Regularization: inversion in a dictionary D where Ux is sparse

Compute z such that $Ux = Dz$ where z is sparse

Non-linear multiscale model

Inversion:

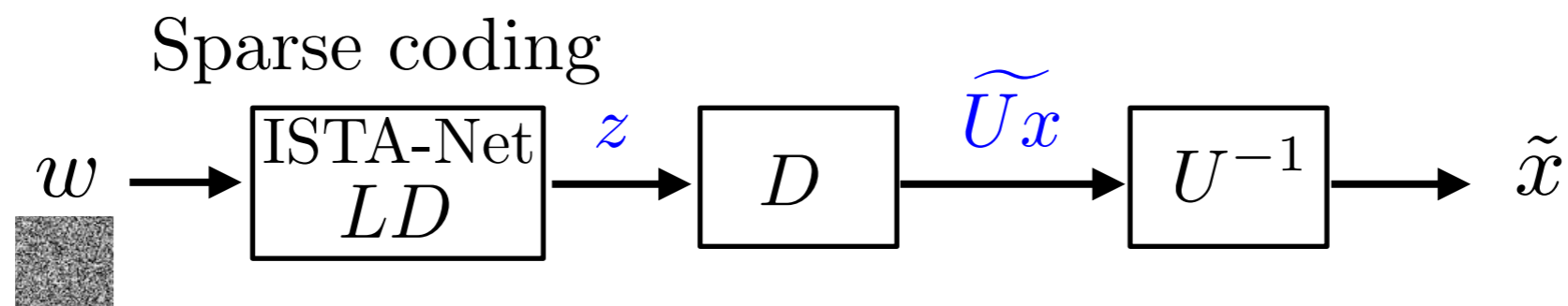


U has a linear inverse U^{-1} , L is non-invertible linear projector

Inversion in a dictionary D where Ux is sparse:

$$Ux = Dz \Rightarrow w = LUx = LDz$$

compute sparse z from w in LD



- How to optimise the dictionary D ?

Learn the dictionary D by minimizing $\sum_i \|x_i - \tilde{x}_i\|^2$

with a stochastic gradient descent on a training set $\{x_i\}_i$

Training Reconstruction

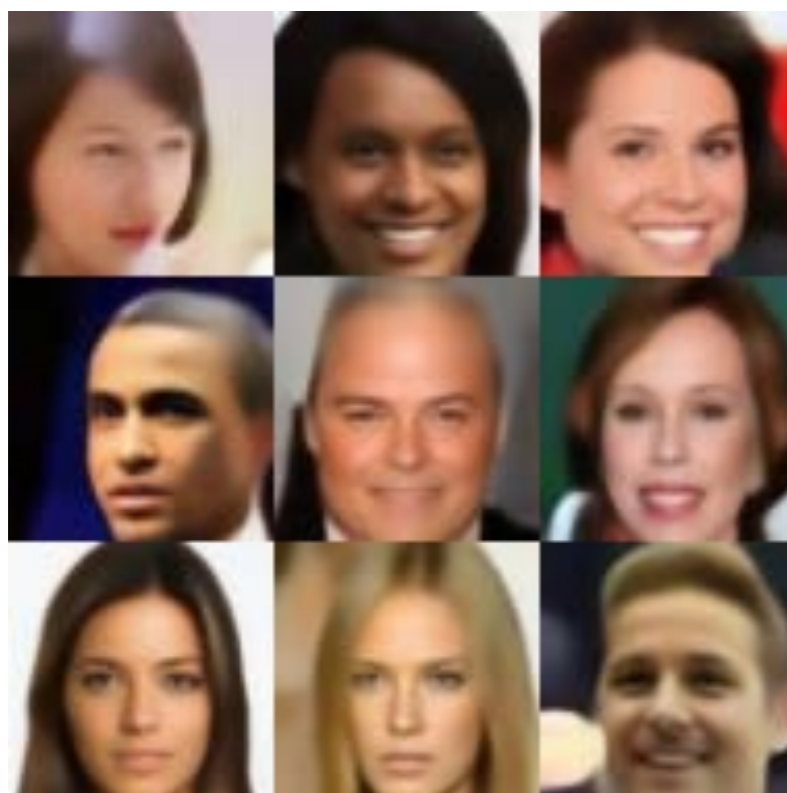
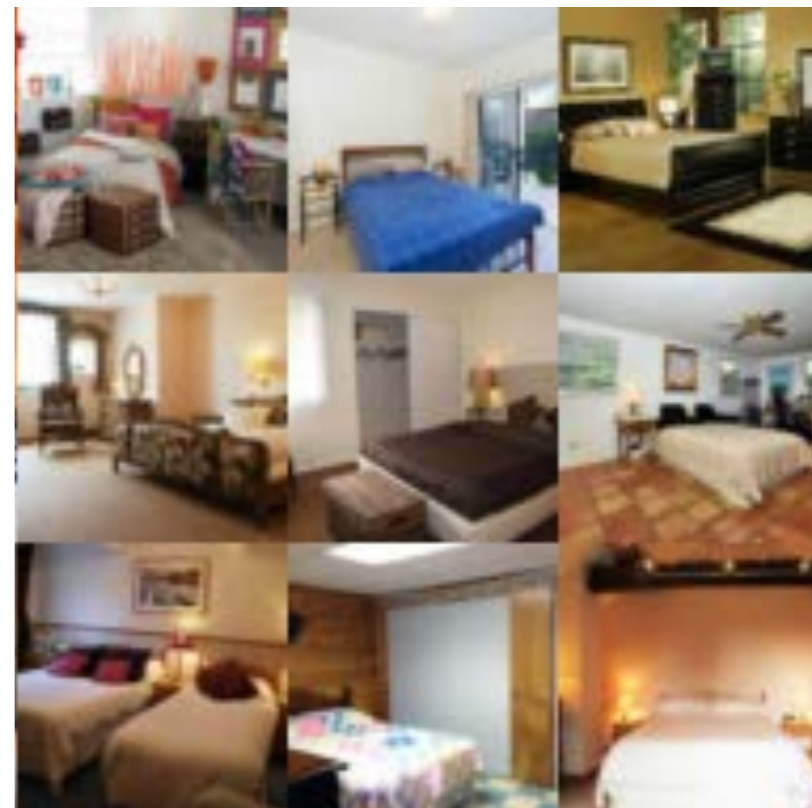
Celebrities Data Basis

Tomas Angles

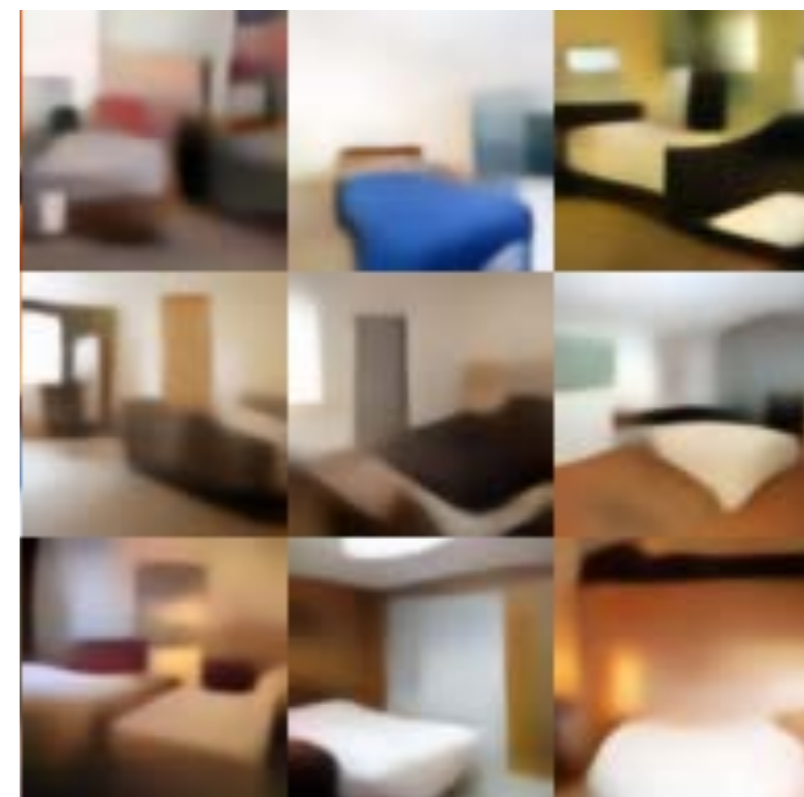
Bedrooms



x_i



\tilde{x}_i



Generative Interpolations

Celebrities

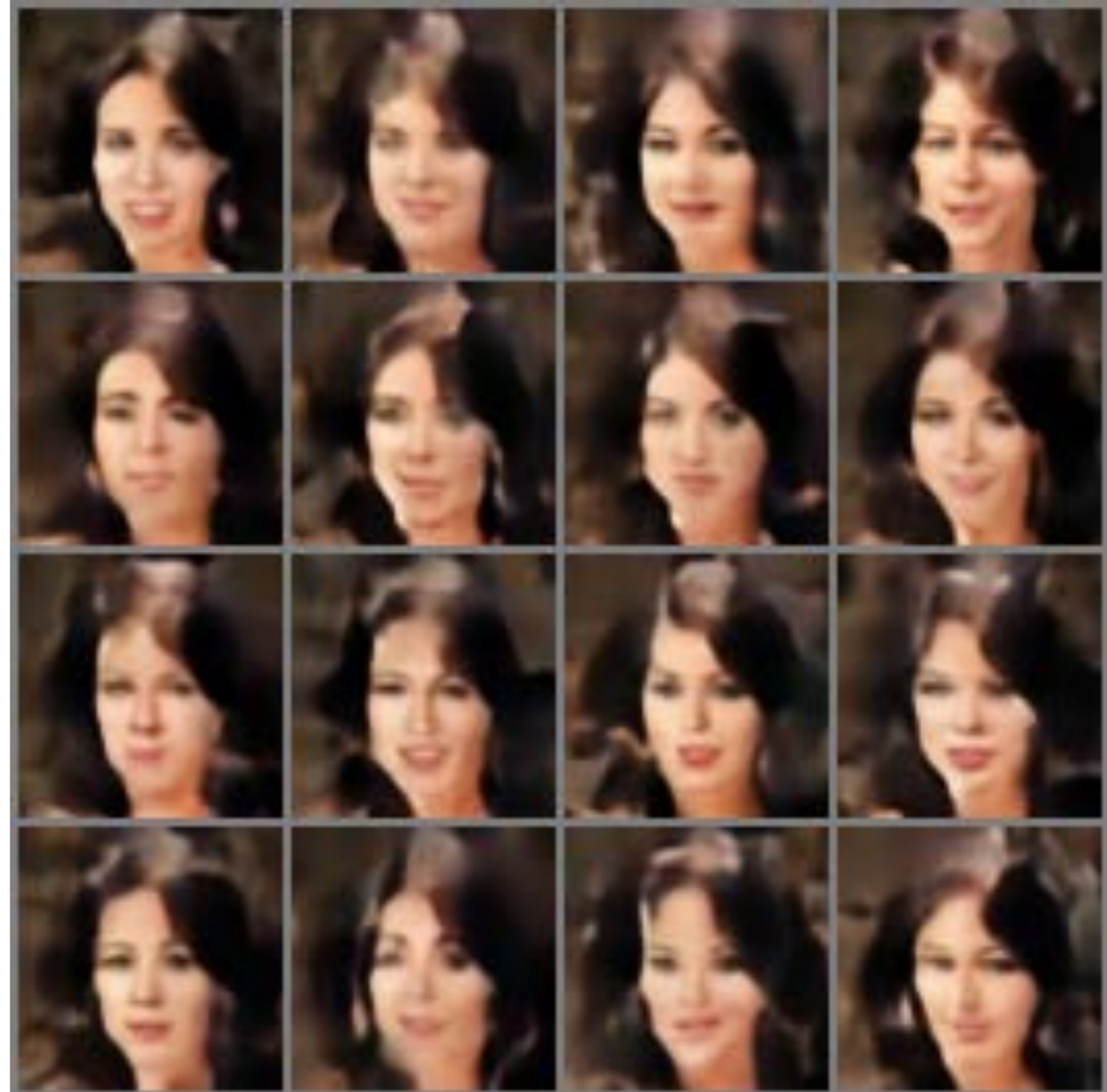
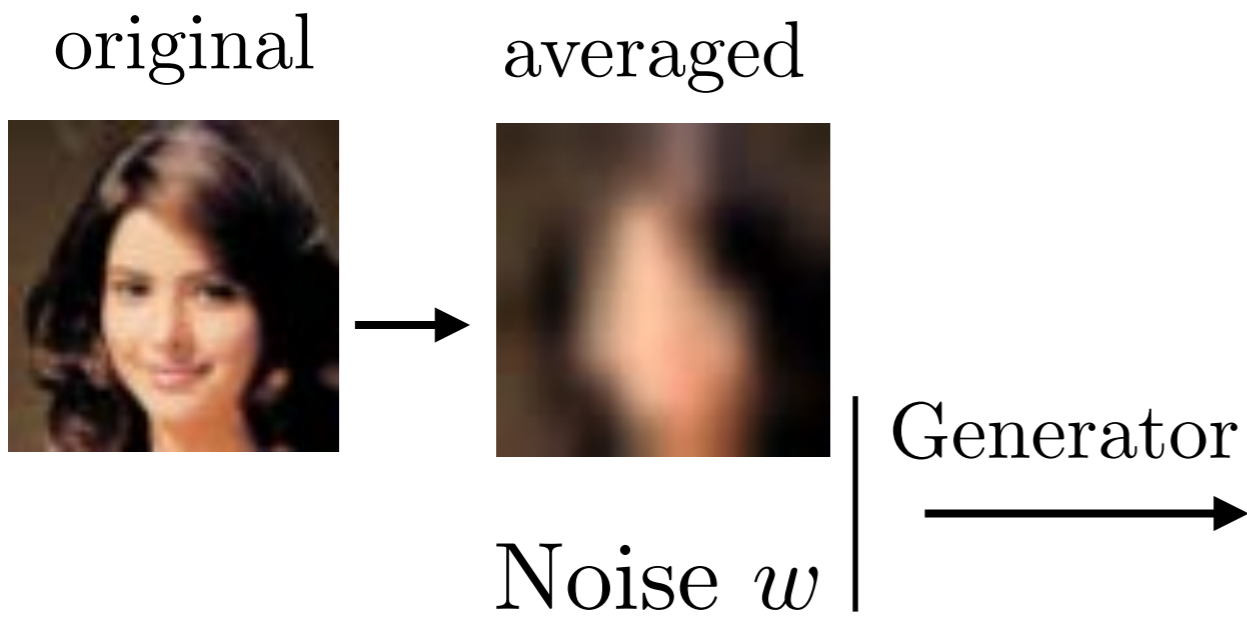
Tomás Angles

$$w_1 \xrightarrow{\quad\quad\quad} w_2$$

$$w = \alpha w_1 + (1 - \alpha)w_2$$

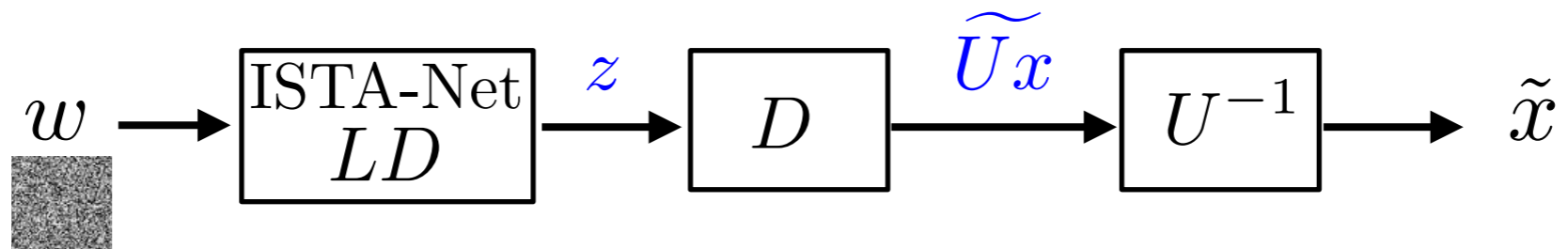


Syntheses with different input noises



Random Generations from Noise

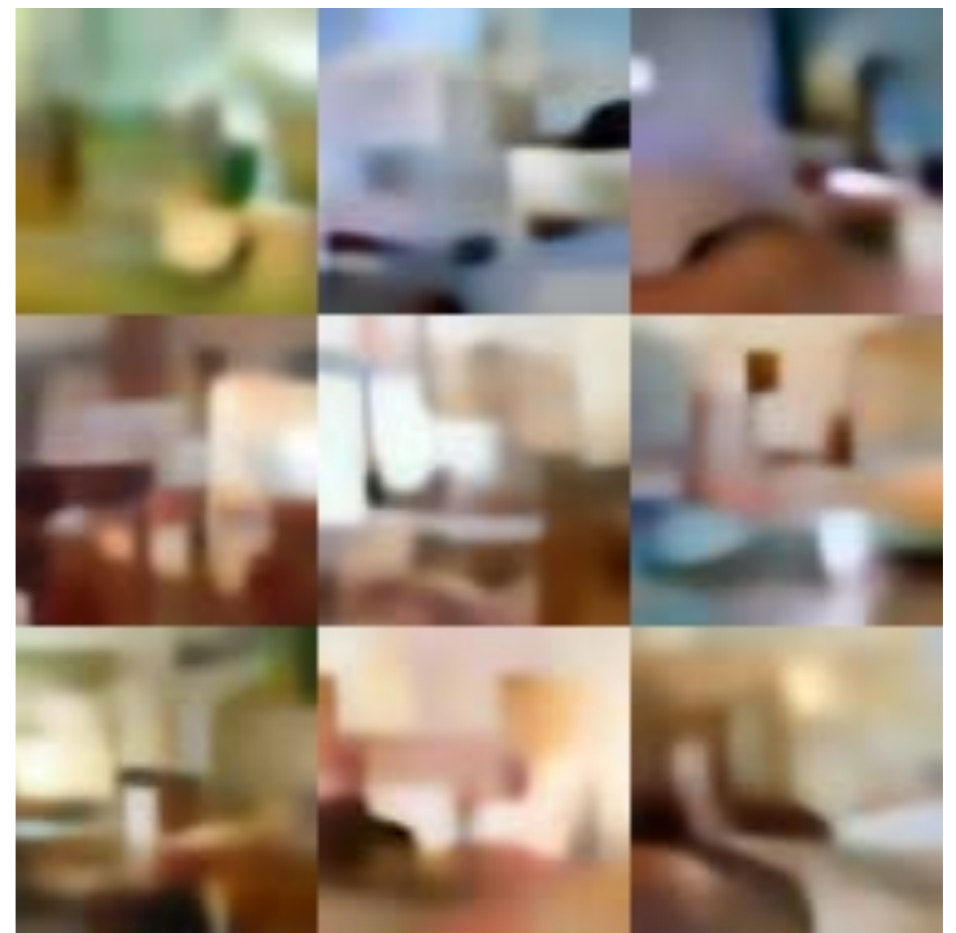
Tomás Angles



Celebrities



Bedrooms



Conclusion

- Deep neural networks are complex computational machines whose flexibility can be compared with Turing machines.
- A ReLU on multiscale filters can produce scale interactions: creates phase harmonics, it may also be used to compute sparse representations, or piecewise linear approximations.
- One can define structured networks which are interpretable: similar to a structured program, with state of the art results.
- Still need functional analysis models and approximation theorems with decay rates.