

Deep Neural Network Mathematical Mysteries for High Dimensional Learning



Stéphane Mallat

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

High Dimensional Learning

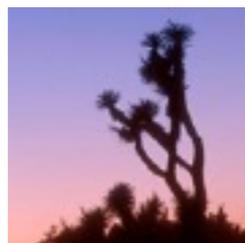
- High-dimensional $x = (x(1), \dots, x(d)) \in \mathbb{R}^d$:
- **Classification:** estimate a class label $f(x)$ given n sample values $\{x_i, y_i = f(x_i)\}_{i \leq n}$

Image Classification $d = 10^6$

Anchor



Joshua Tree



Beaver



Lotus



Water Lily



Huge variability
inside classes

Find invariants

High Dimensional Learning

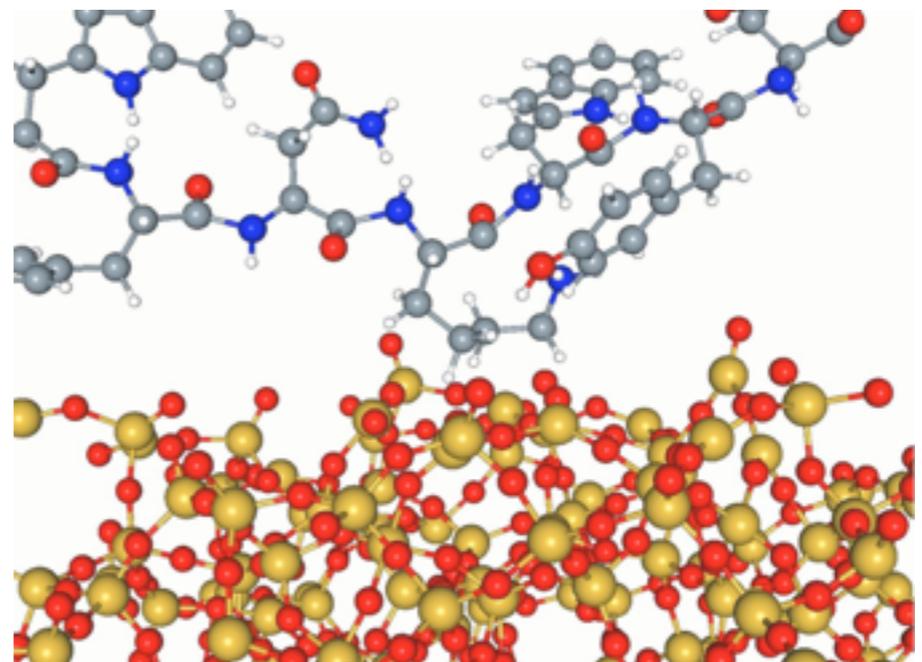
- High-dimensional $x = (x(1), \dots, x(d)) \in \mathbb{R}^d$:
- **Regression:** approximate a *functional* $f(x)$
given n sample values $\{x_i, y_i = f(x_i) \in \mathbb{R}\}_{i \leq n}$

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

Astronomy



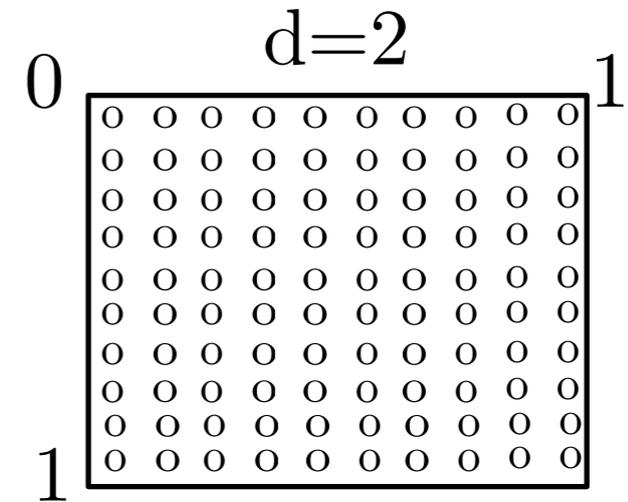
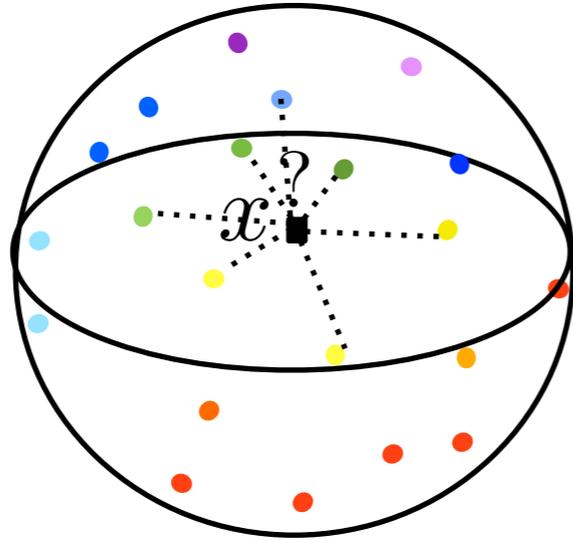
Quantum Chemistry



Importance of symmetries.

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:

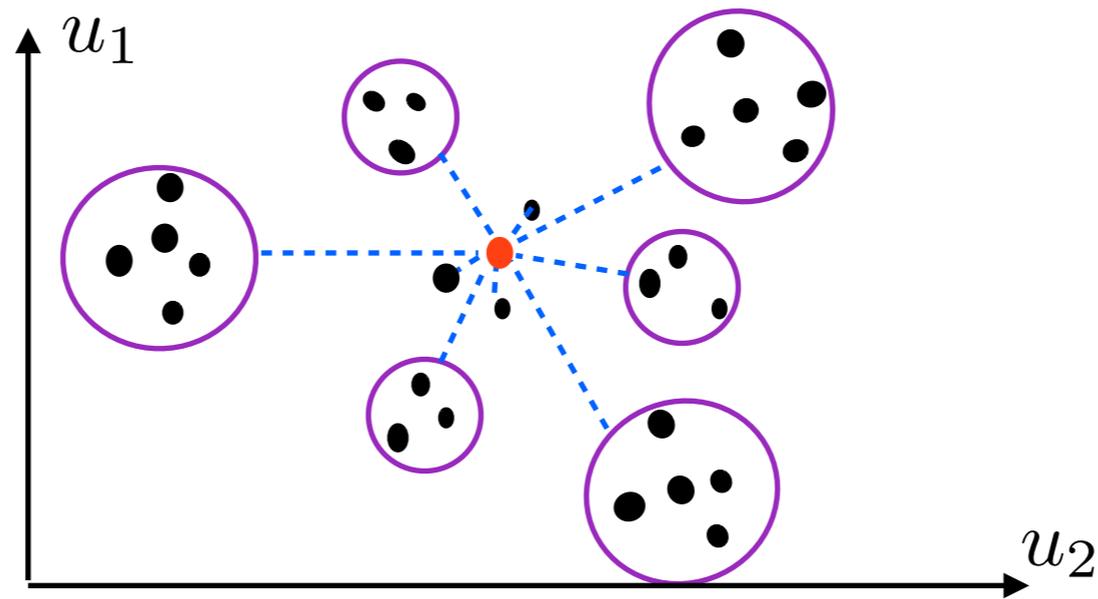


- Need ϵ^{-d} points to cover $[0, 1]^d$ at a Euclidean distance ϵ

Problem: $\|x - x_i\|$ is always large



- Variables $x(u)$ indexed by a low-dimensional u : time/space... pixels in images, particles in physics, words in text...
- Multiscale interactions of d variables:



From d^2 interactions to $O(\log^2 d)$ multiscale interactions.

- Multiscale analysis: wavelets on groups of symmetries.
hierarchical architecture.

- 1 Hidden Layer Network, Approximation theory and Curse
- Kernel learning
- Dimension reduction with change of variables
- Deep Neural networks and symmetry groups
- Wavelet Scattering transforms
- Applications and many open questions

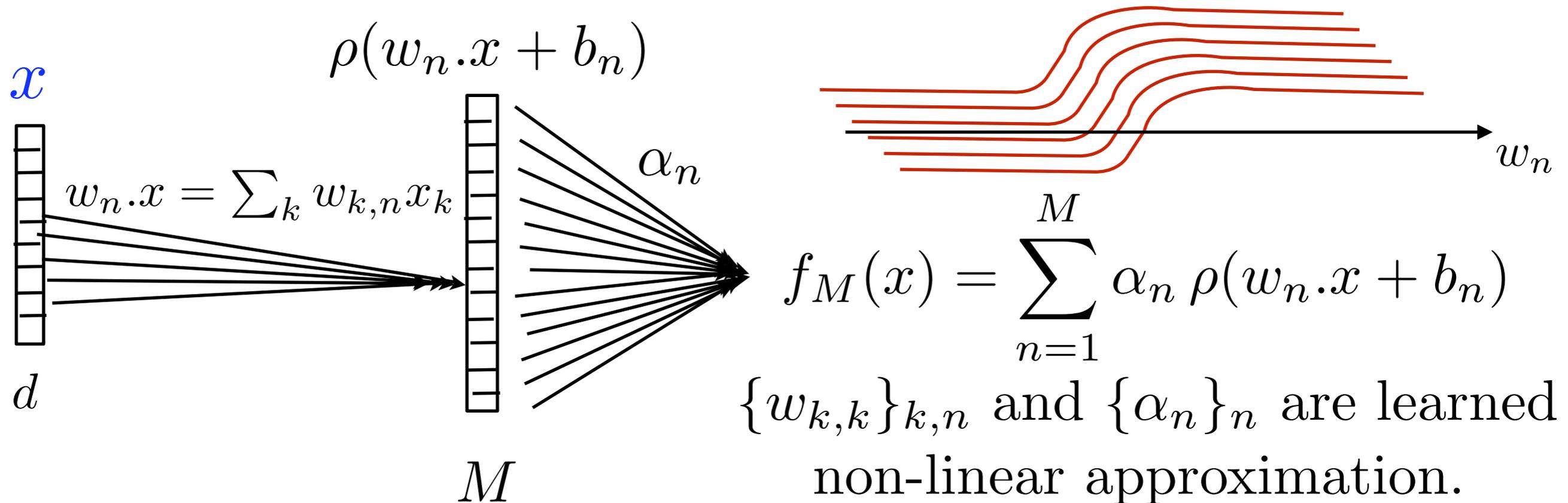
Understanding Deep Convolutional Networks, arXiv 2016.

Learning as an Approximation

- To estimate $f(x)$ from a sampling $\{x_i, y_i = f(x_i)\}_{i \leq M}$ we must build an M -parameter approximation f_M of f .
- Precise sparse approximation requires some "regularity".
- For binary classification $f(x) = \begin{cases} 1 & \text{if } x \in \Omega \\ -1 & \text{if } x \notin \Omega \end{cases}$
$$f(x) = \text{sign}(\tilde{f}(x))$$
where \tilde{f} is potentially regular.
- What type of regularity ? How to compute f_M ?

1 Hidden Layer Neural Networks

One-hidden layer neural network: ridge functions $\rho(x \cdot w_n + b_n)$



Cybenko, Hornik, Stinchcombe, White

Theorem: For "reasonable" bounded $\rho(u)$

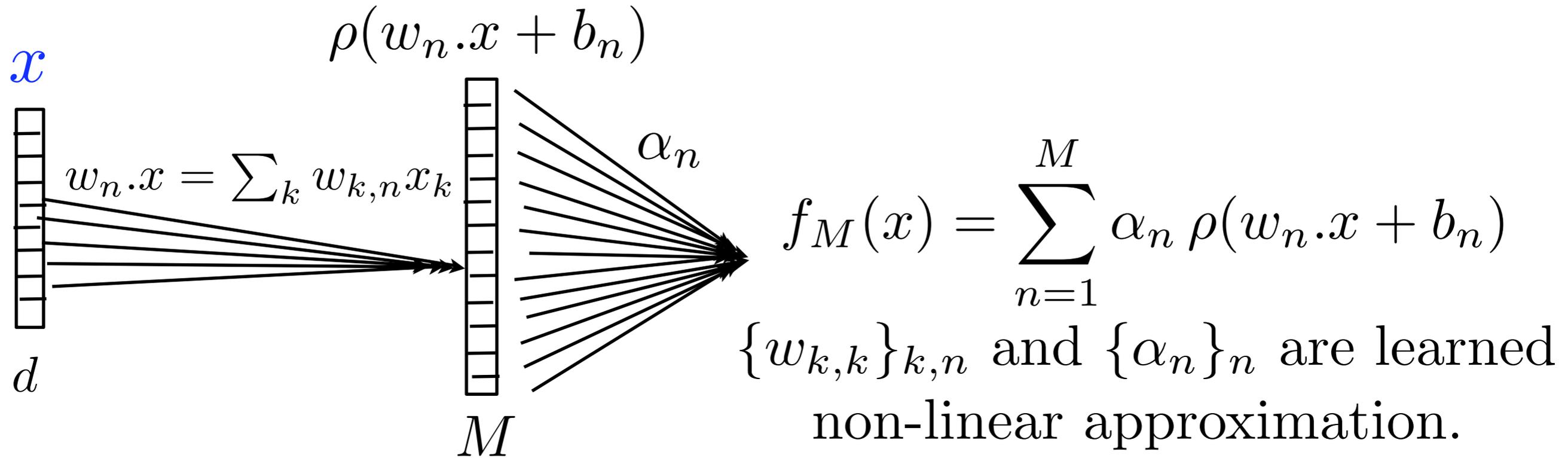
and appropriate choices of $w_{n,k}$ and α_n :

$$\forall f \in \mathbb{L}^2[0, 1]^d \quad \lim_{M \rightarrow \infty} \|f - f_M\| = 0 .$$

No big deal: curse of dimensionality still there.

1 Hidden Layer Neural Networks

One-hidden layer neural network:



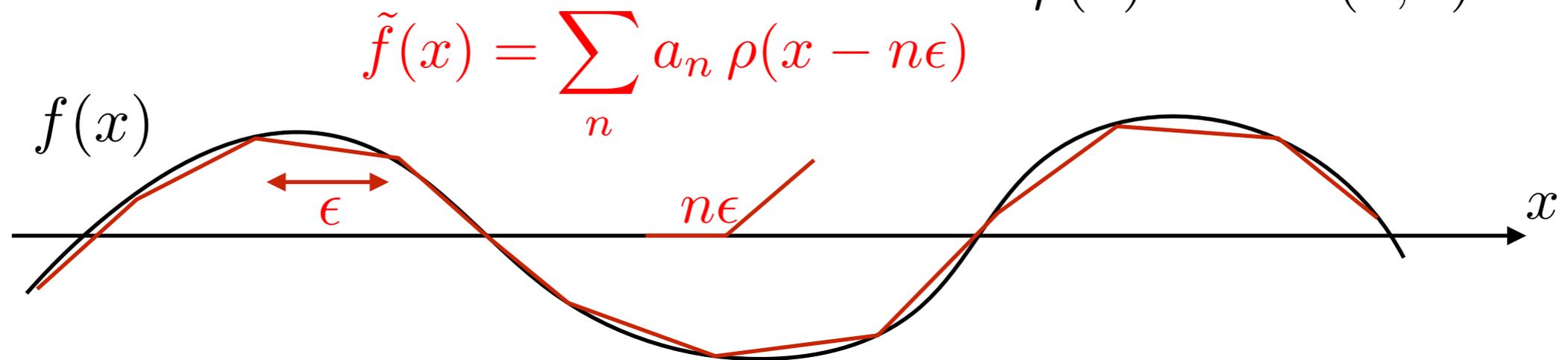
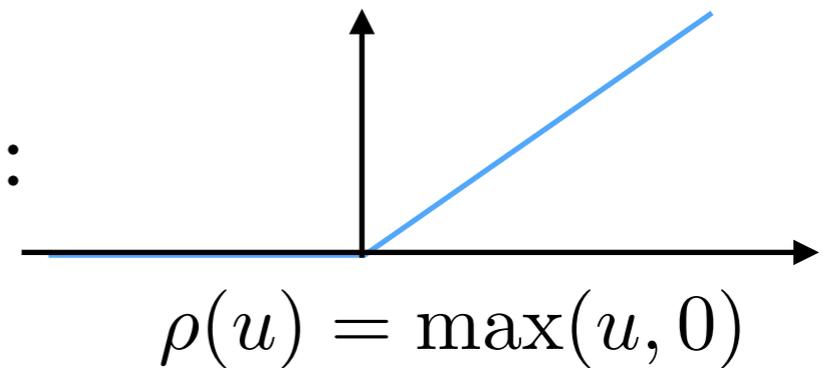
Fourier series: $\rho(u) = e^{iu}$

$$f_M(x) = \sum_{n=1}^M \alpha_n e^{i w_n \cdot x}$$

For nearly all ρ : essentially same approximation results.

Piecewise Linear Approximation

- Piecewise linear approximation:



If f is Lipschitz: $|f(x) - f(x')| \leq C |x - x'|$

$$\Rightarrow |f(x) - \tilde{f}(x)| \leq C \epsilon.$$

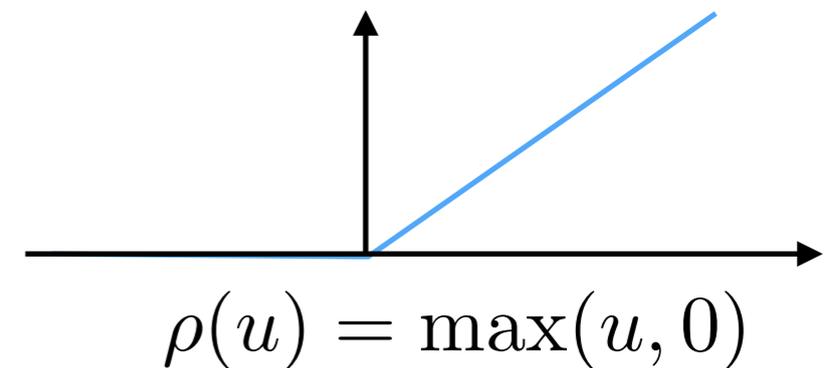
Need $M = \epsilon^{-1}$ points to cover $[0, 1]$ at a distance ϵ

$$\Rightarrow \|f - f_M\| \leq C M^{-1}$$

Linear Ridge Approximation

- Piecewise linear ridge approximation: $x \in [0, 1]^d$

$$\tilde{f}(x) = \sum_n a_n \rho(w_n \cdot x - n\epsilon)$$



If f is Lipschitz: $|f(x) - f(x')| \leq C \|x - x'\|$

Sampling at a distance ϵ :

$$\Rightarrow |f(x) - \tilde{f}(x)| \leq C \epsilon.$$

need $M = \epsilon^{-d}$ points to cover $[0, 1]^d$ at a distance ϵ

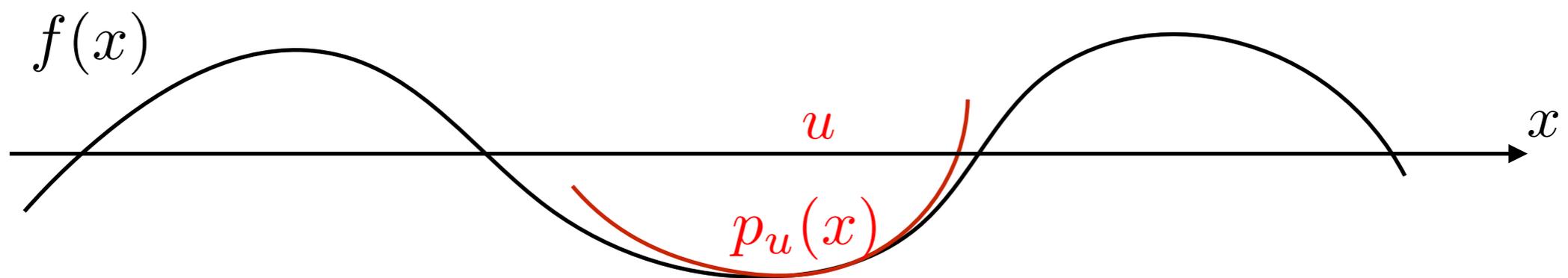
$$\Rightarrow \|f - f_M\| \leq C M^{-1/d}$$

Curse of dimensionality!

Approximation with Regularity

- What prior condition makes learning possible ?
- Approximation of regular functions in $\mathbf{C}^s[0, 1]^d$:

$\forall x, u \quad |f(x) - p_u(x)| \leq C |x - u|^s$ with $p_u(x)$ polynomial



$$|x - u| \leq \epsilon^{1/s} \quad \Rightarrow \quad |f(x) - p_u(x)| \leq C \epsilon$$

Need $M^{-d/s}$ point to cover $[0, 1]^d$ at a distance $\epsilon^{1/s}$

$$\Rightarrow \|f - f_M\| \leq C M^{-s/d}$$

- Can not do better in $\mathbf{C}^s[0, 1]^d$, not good because $s \ll d$.

Failure of classical approximation theory.

Kernel Learning

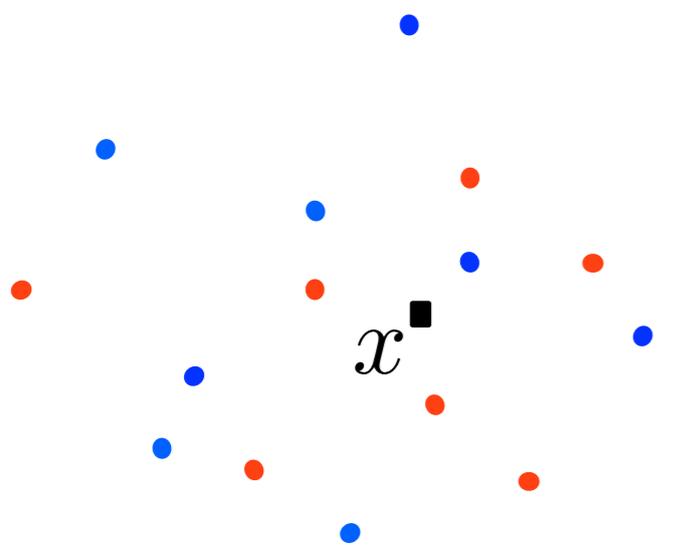
Change of variable $\Phi(x) = \{\phi_k(x)\}_{k \leq d'}$

to nearly linearize $f(x)$, which is approximated by:

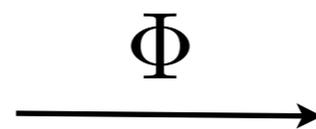
$$\tilde{f}(x) = \langle \Phi(x), w \rangle = \sum_k w_k \phi_k(x) .$$

1D projection

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

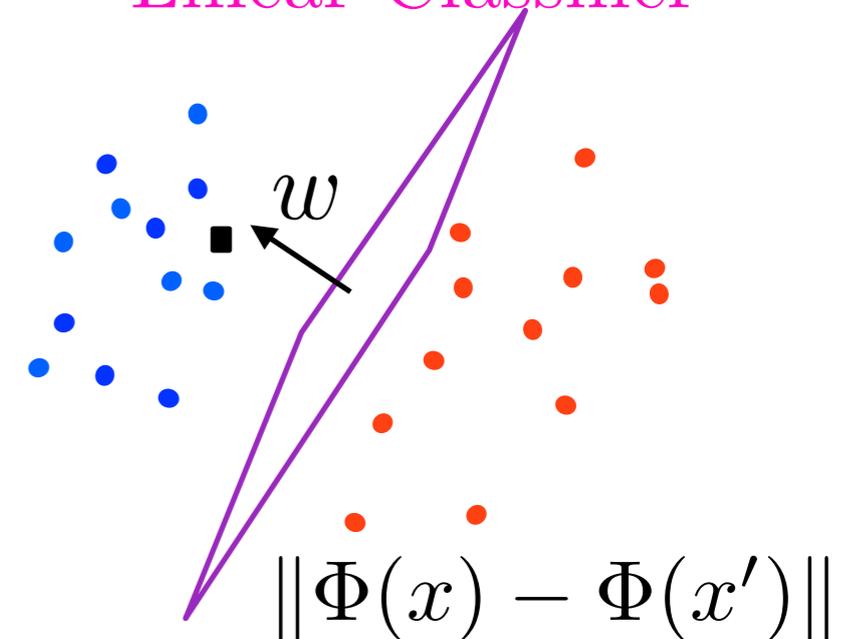


Metric: $\|x - x'\|$



$\Phi(x) \in \mathbb{R}^{d'}$

Linear Classifier



- How and when is possible to find such a Φ ?
- What "regularity" of f is needed ?

Increase Dimensionality

Proposition: There exists a hyperplane separating any two subsets of N points $\{\Phi x_i\}_i$ in dimension $d' > N + 1$ if $\{\Phi x_i\}_i$ are not in an affine subspace of dimension $< N$.

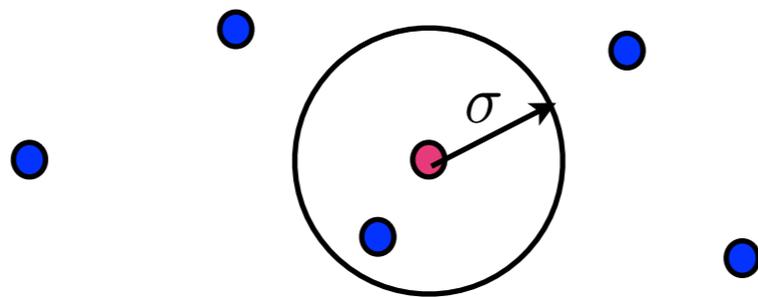
\Rightarrow Choose Φ increasing dimensionality !

Problem: generalisation, overfitting.

Example: Gaussian kernel $\langle \Phi(x), \Phi(x') \rangle = \exp\left(\frac{-\|x - x'\|^2}{2\sigma^2}\right)$

$\Phi(x)$ is of dimension $d' = \infty$

If σ is small, nearest neighbor classifier type:



Reduction of Dimensionality

- Discriminative change of variable $\Phi(x)$:

$$\Phi(x) \neq \Phi(x') \quad \text{if} \quad f(x) \neq f(x')$$

$$\Rightarrow \exists \tilde{f} \quad \text{with} \quad f(x) = \tilde{f}(\Phi(x))$$

- If \tilde{f} is Lipschitz: $|\tilde{f}(z) - \tilde{f}(z')| \leq C \|z - z'\|$

$$z = \Phi(x) \quad \Leftrightarrow \quad |f(x) - f(x')| \leq C \|\Phi(x) - \Phi(x')\|$$

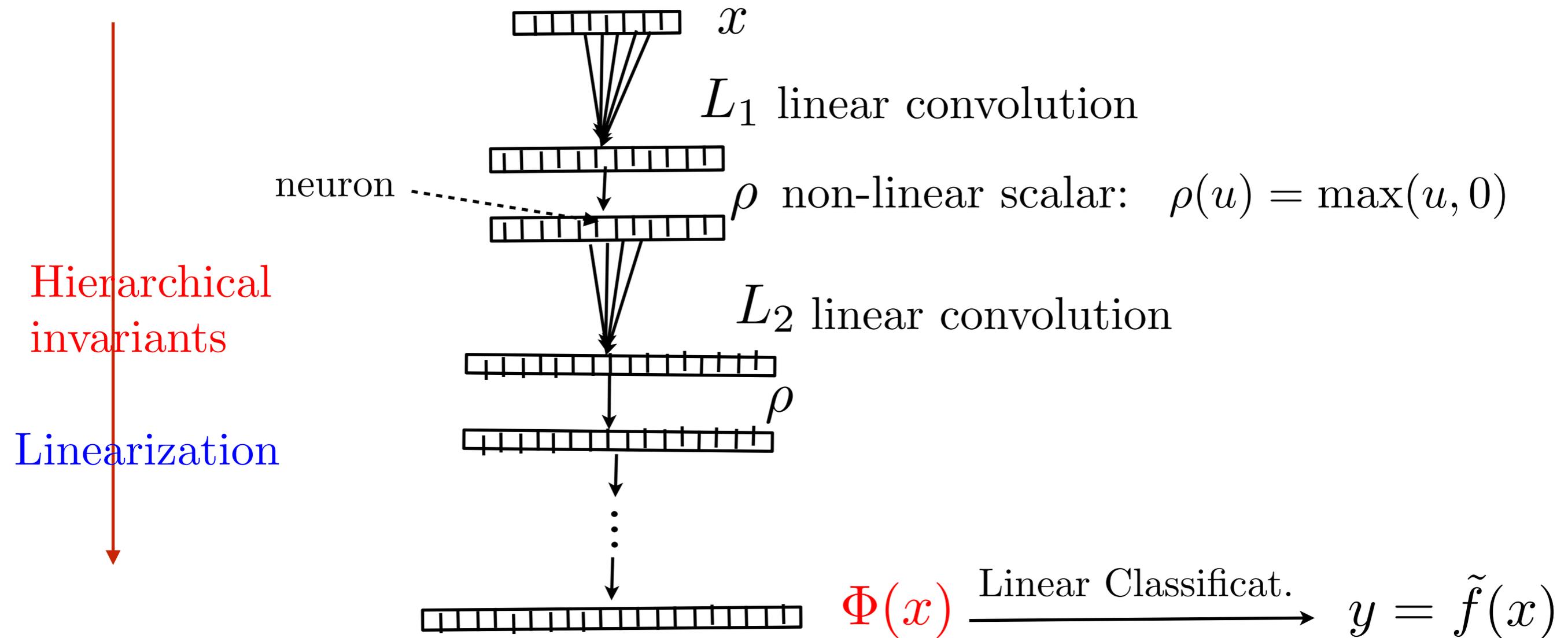
$$\text{Discriminative:} \quad \|\Phi(x) - \Phi(x')\| \geq C^{-1} |f(x) - f(x')|$$

- For $x \in \Omega$, if $\Phi(\Omega)$ is bounded and a low dimension d'

$$\Rightarrow \|f - f_M\| \leq C M^{-1/d'}$$

Deep Convolution Networks

- The revival of neural networks: *Y. LeCun*



Optimize L_j with **architecture constraints**: over 10^9 parameters

Exceptional results for *images, speech, language, bio-data...*

Why does it work so well ? A difficult problem

ImageNet Data Basis

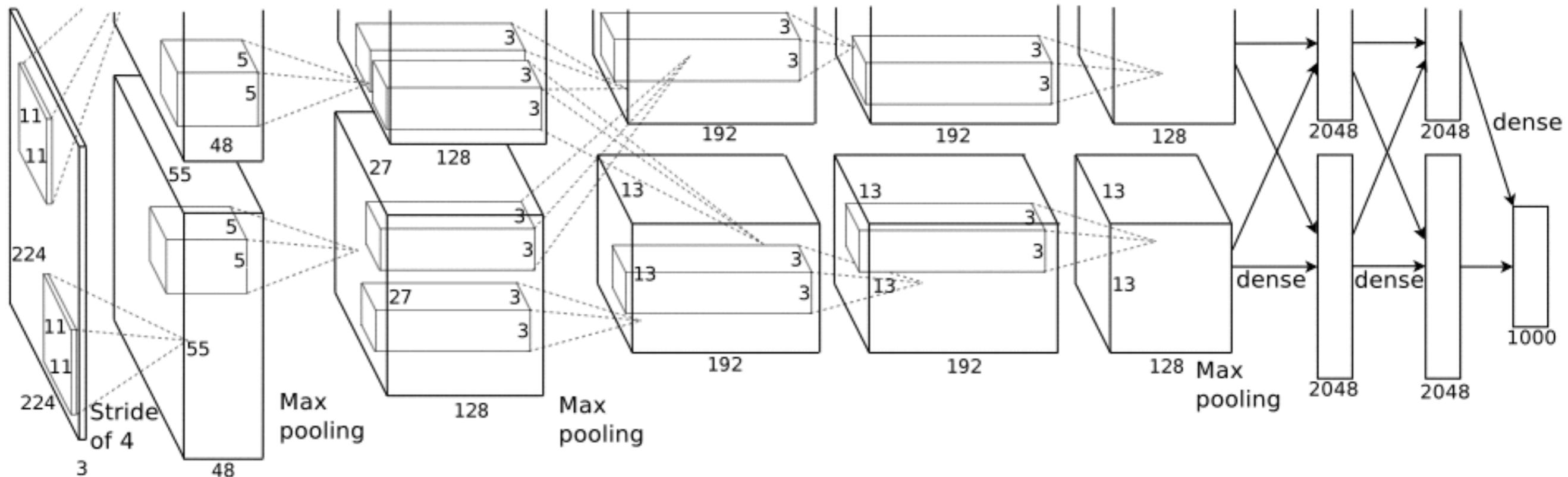
- Data basis with 1 million images and 2000 classes



Alex Deep Convolution Network

A. Krizhevsky, Sutsver, Hinton

- Imagenet supervised training: $1.2 \cdot 10^6$ examples, 10^3 classes
15.3% testing error in 2012



Wavelets

New networks with 5% errors.
Up to 150 layers!

Image Classification



mite



container ship



motor scooter



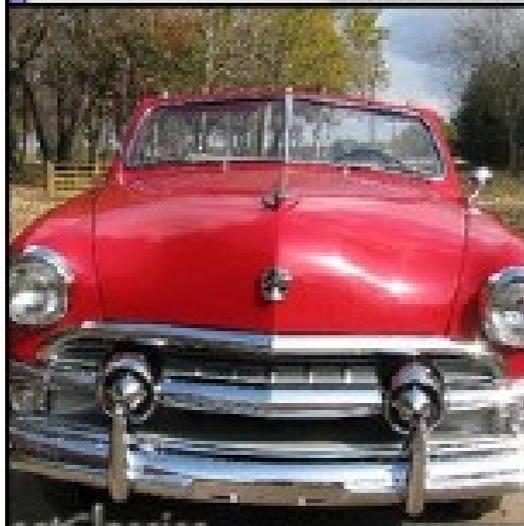
leopard

█	mite
█	black widow
█	cockroach
█	tick
█	starfish

█	container ship
█	lifeboat
█	amphibian
█	fireboat
█	drilling platform

█	motor scooter
█	go-kart
█	moped
█	bumper car
█	golfcart

█	leopard
█	jaguar
█	cheetah
█	snow leopard
█	Egyptian cat



grille



mushroom



cherry



Madagascar cat

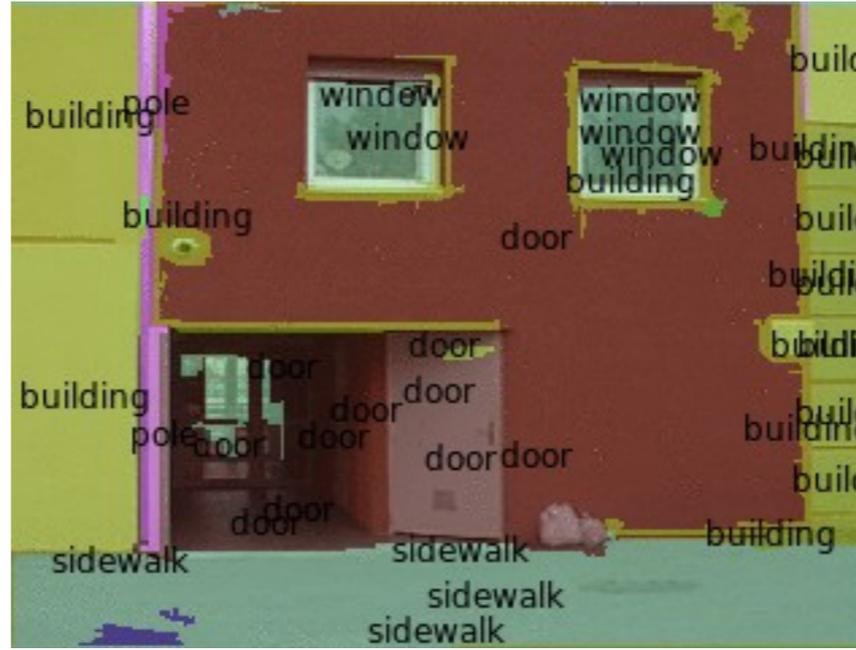
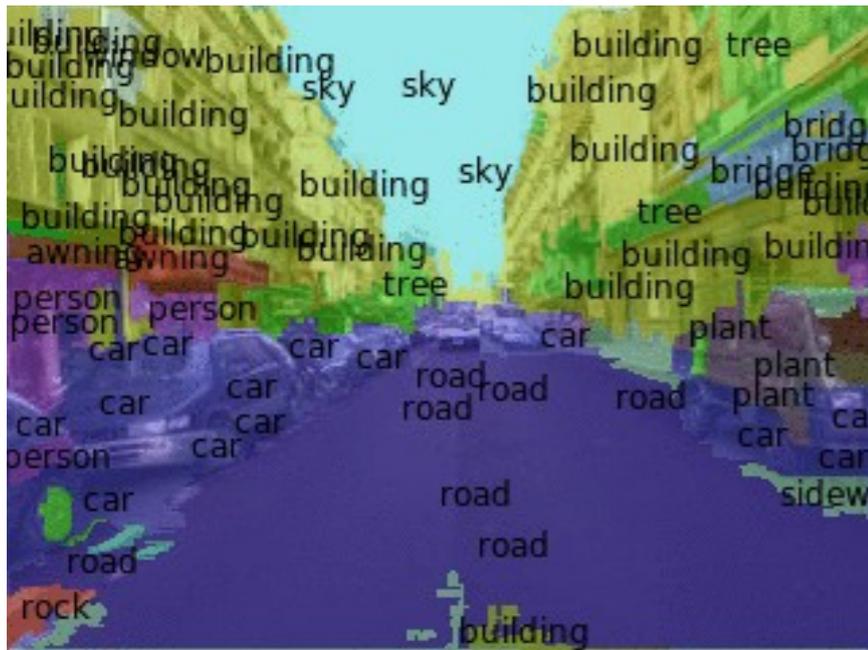
█	convertible
█	grille
█	pickup
█	beach wagon
█	fire engine

█	agaric
█	mushroom
█	jelly fungus
█	gill fungus
█	dead-man's-fingers

█	dalmatian
█	grape
█	elderberry
█	ffordshire bullterrier
█	currant

█	squirrel monkey
█	spider monkey
█	titi
█	indri
█	howler monkey

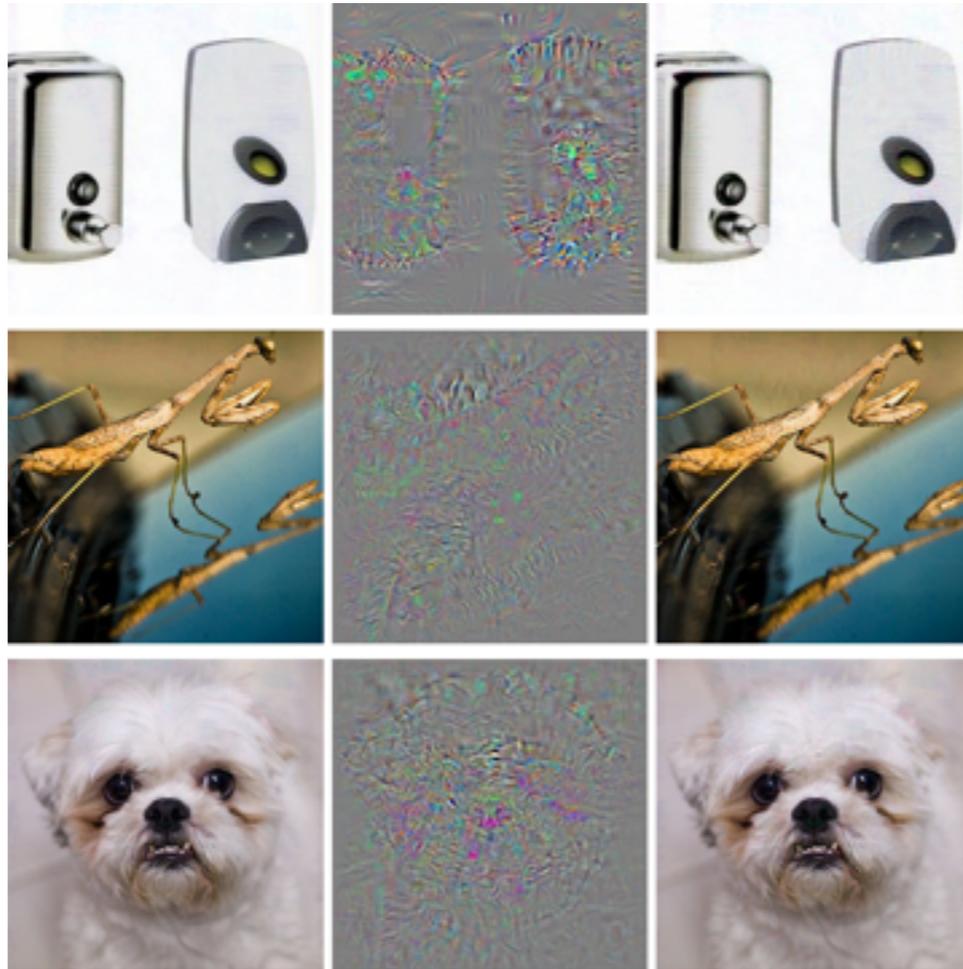
Scene Labeling / Car Driving



Why Understading ?

Szegedy, Zaremba, Sutskever, Bruna, Erhan, Goodfellow, Fergus

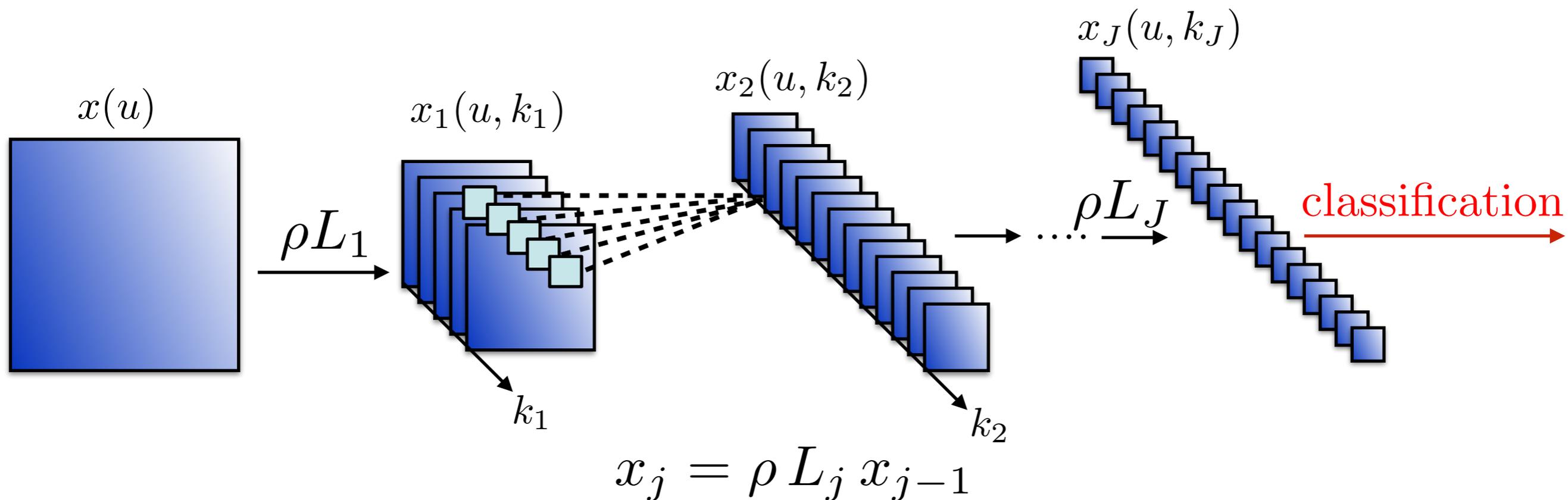
$$x + \epsilon = \tilde{x} \quad \text{with} \quad \|\epsilon\| < 10^{-2} \|x\|$$



correctly
classified

classified as
ostrich

- Trial and error testing can not guarantee reliability.



- L_j is a linear combination of convolutions and subsampling:

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

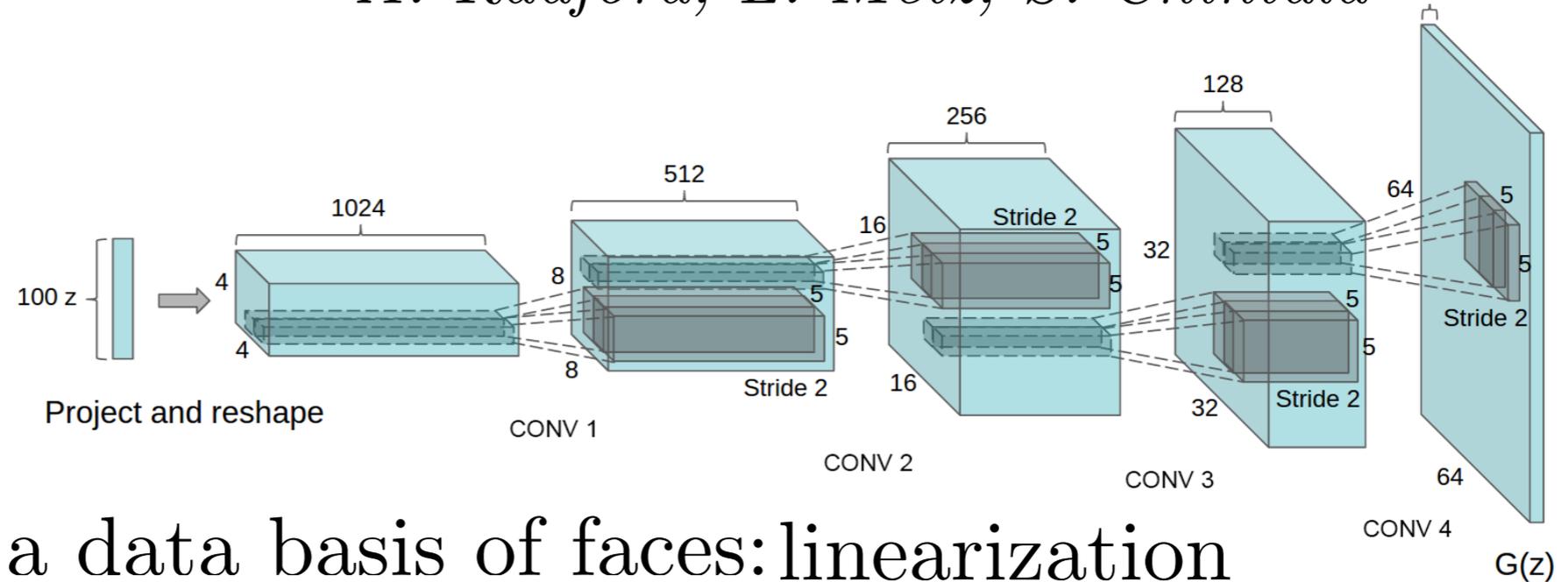
sum across channels

- ρ is contractive: $|\rho(u) - \rho(u')| \leq |u - u'|$

$$\rho(u) = \max(u, 0) \text{ or } \rho(u) = |u|$$

Linearisation in Deep Networks

A. Radford, L. Metz, S. Chintala

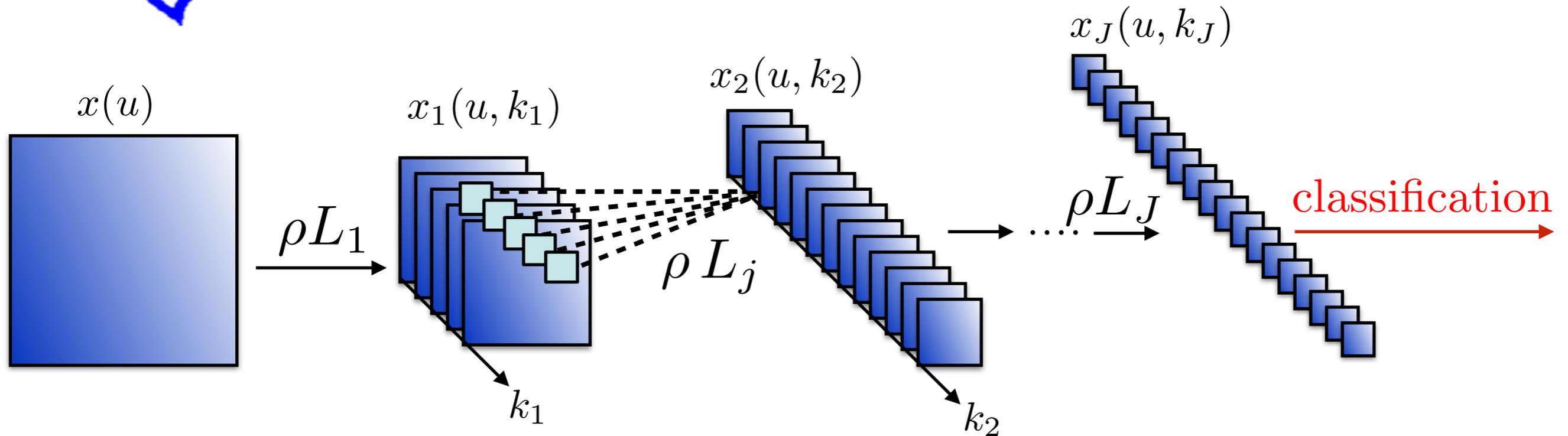


- Trained on a data basis of faces: linearization



- On a data basis including bedrooms: interpolations



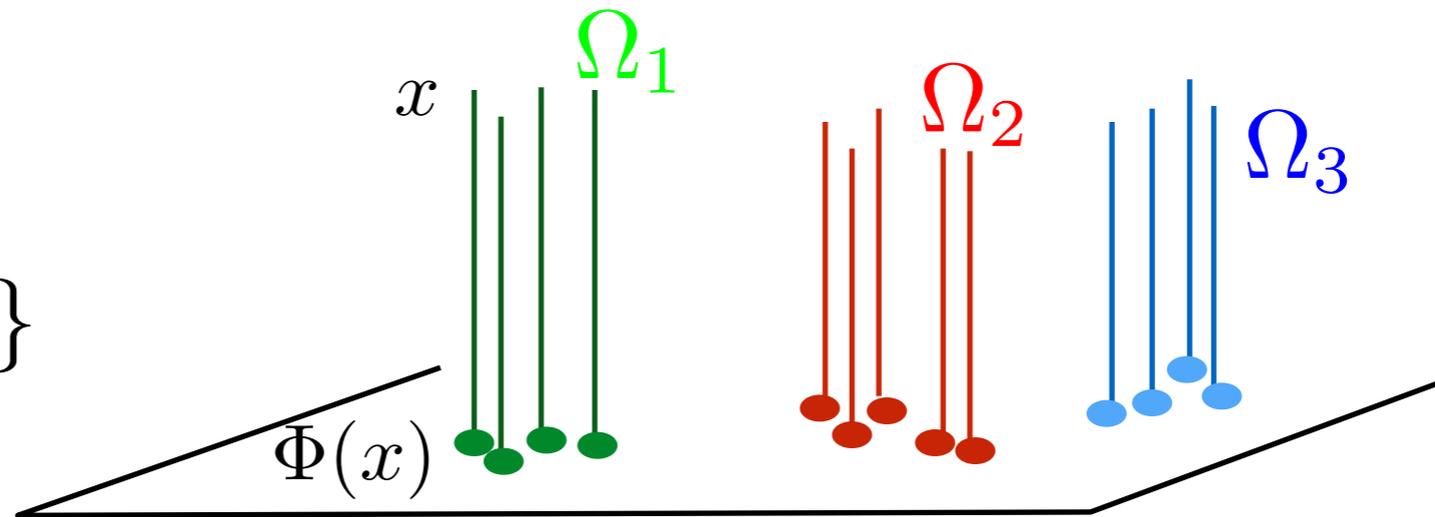


- Why convolutions ? Translation covariance.
- Why no overfitting ? Contractions, dimension reduction
- Why hierarchical cascade ?
- Why introducing non-linearities ?
- How and what to linearise ?
- What are the roles of the multiple channels in each layer ?

Classes

Level sets of $f(x)$

$$\Omega_t = \{x : f(x) = t\}$$

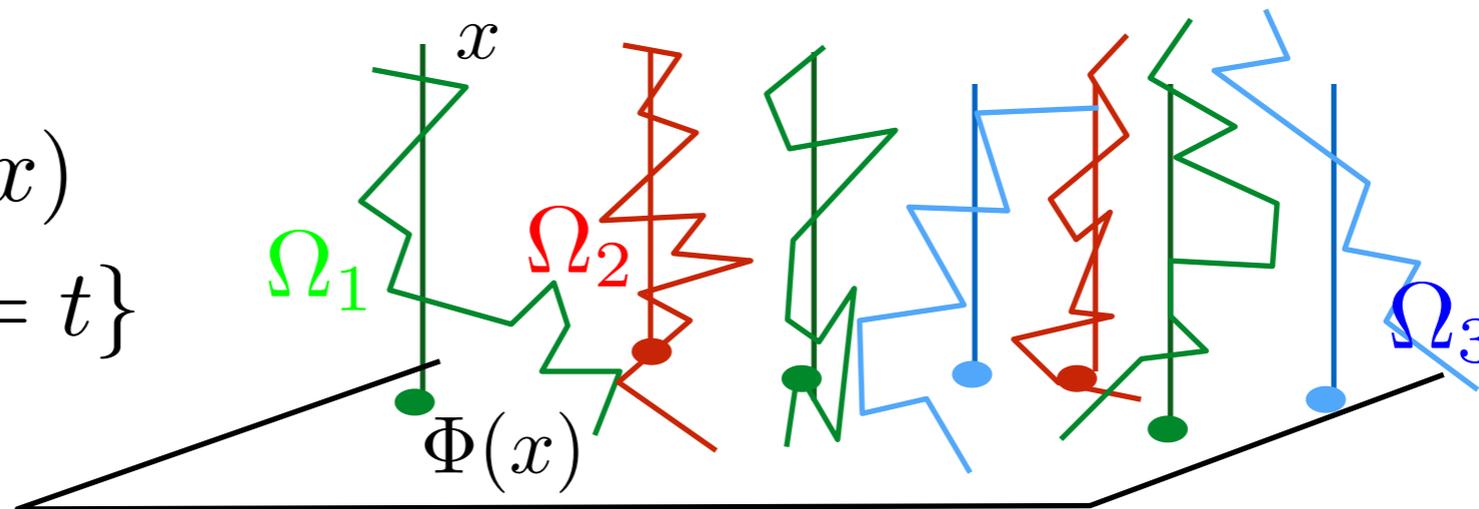


If level sets (classes) are parallel to a linear space
then variables are eliminated by linear projections: *invariants*.

Classes

Level sets of $f(x)$

$$\Omega_t = \{x : f(x) = t\}$$

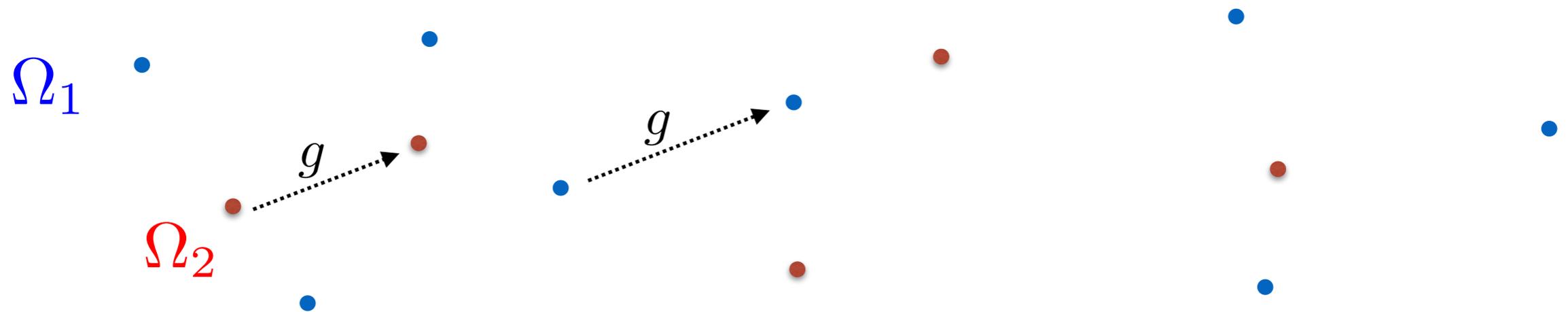


- If level sets Ω_t are not parallel to a linear space
 - Linearise them with a change of variable $\Phi(x)$
 - Then reduce dimension with linear projections
- Difficult because Ω_t are high-dimensional, irregular, known on few samples.

Level Set Geometry: Symmetries

- Curse of dimensionality \Rightarrow not local but global geometry

Level sets: classes, characterised by their global symmetries.



- A symmetry is an operator g which preserves level sets:

$$\forall x \quad , \quad f(g.x) = f(x) : \text{global}$$

If g_1 and g_2 are symmetries then $g_1.g_2$ is also a symmetry

$$f(g_1.g_2.x) = f(g_2.x) = f(x)$$

Groups of symmetries

- $G = \{ \text{all symmetries} \}$ is a group: unknown

$$\forall (g, g') \in G^2 \Rightarrow g.g' \in G$$

Inverse: $\forall g \in G, g^{-1} \in G$

Associative: $(g.g').g'' = g.(g'.g'')$

If commutative $g.g' = g'.g$: Abelian group.

- Group of dimension n if it has n generators:

$$g = g_1^{p_1} g_2^{p_2} \dots g_n^{p_n}$$

- Lie group: infinitely small generators (Lie Algebra)

Translation and Deformations

- Digit classification:

$$x(u) \quad x'(u) = x(u - \tau(u))$$



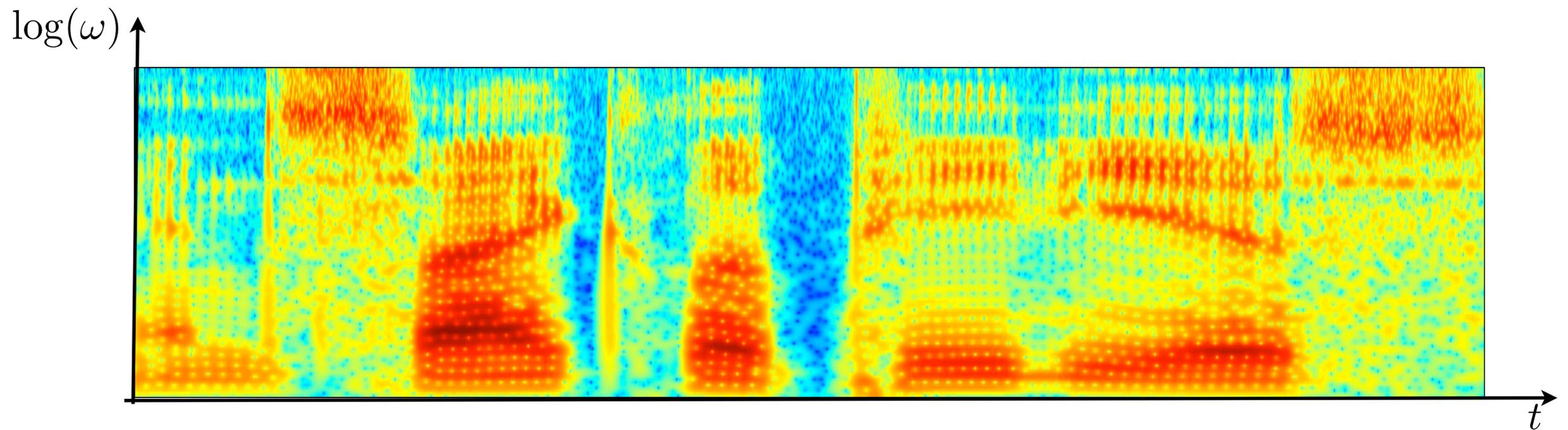
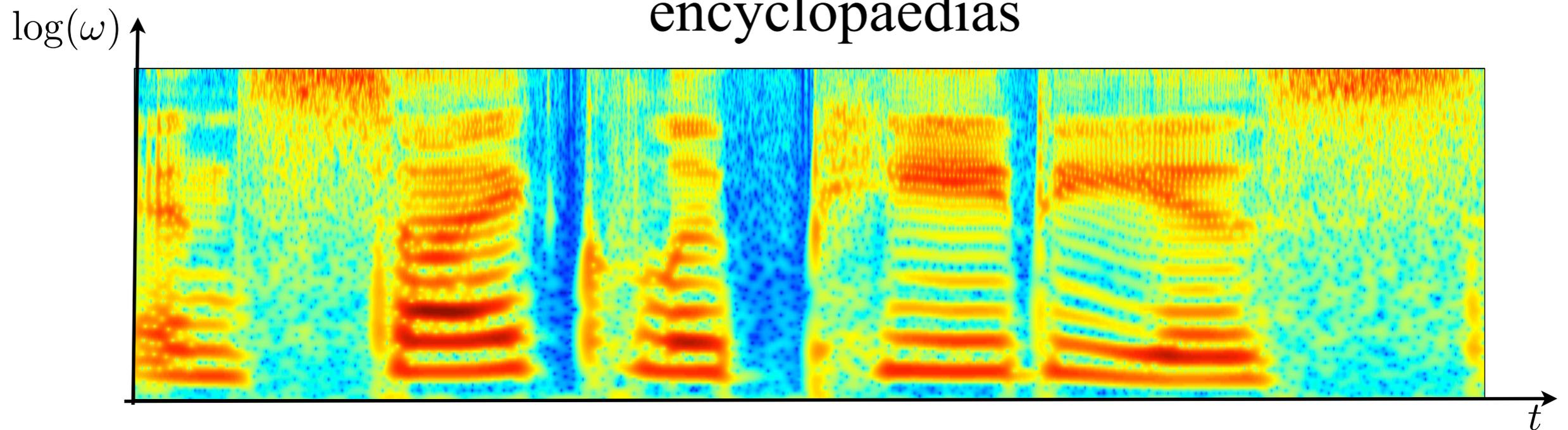
- Globally invariant to the translation group: small
- Locally invariant to small diffeomorphisms: huge group



Video of Philipp Scott Johnson

Frequency Transpositions

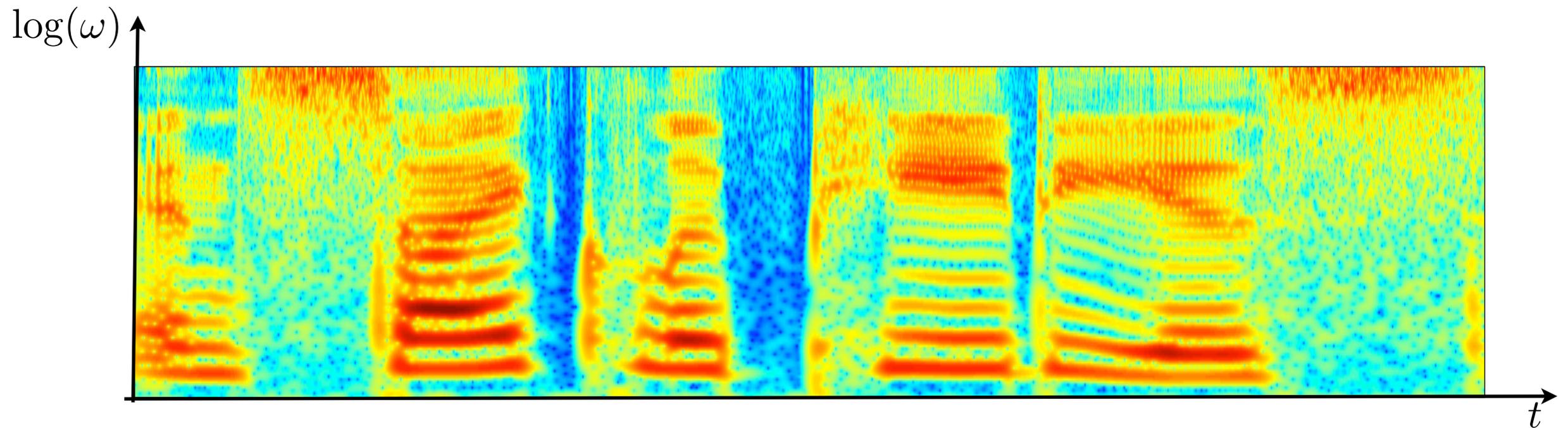
encyclopaedias



H : Heisenberg group of "time-frequency" translations

Frequency Transpositions

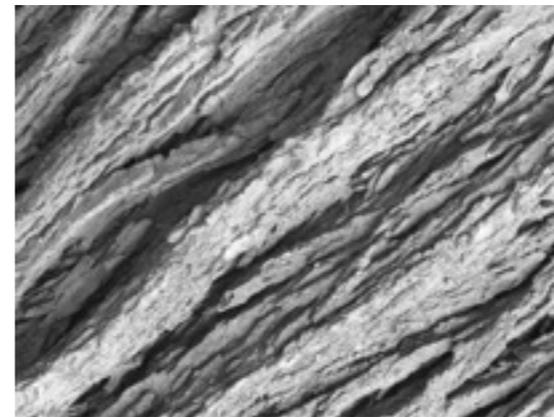
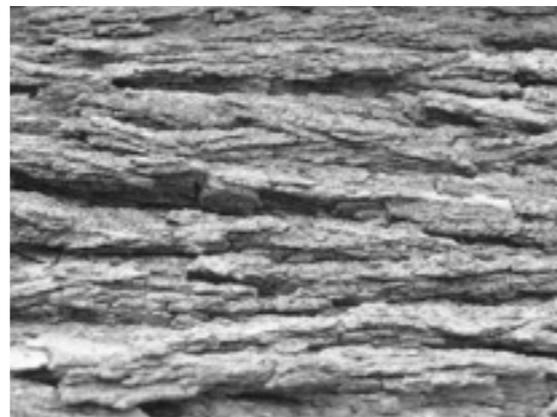
Time and frequency translations and deformations:



- Frequency transposition invariance is needed for speech recognition not for locutor recognition.

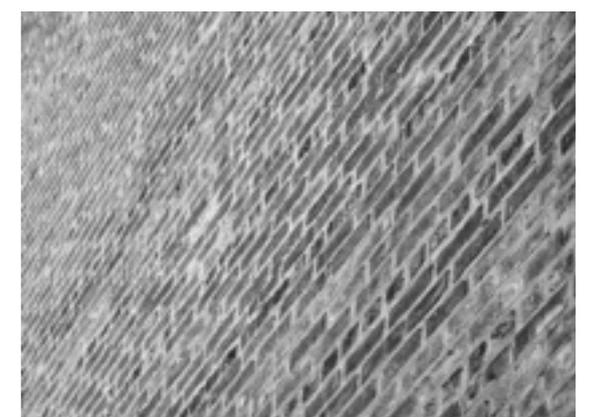
Rotation and Scaling Variability

- Rotation and deformations



Group: $SO(2) \times \text{Diff}(SO(2))$

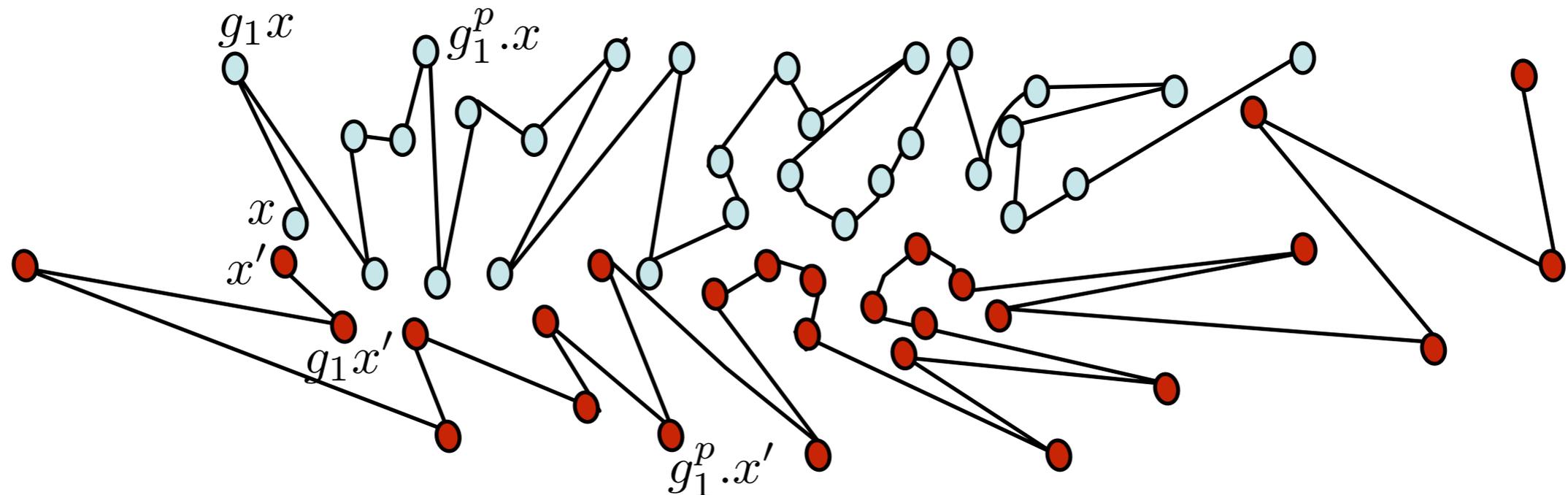
- Scaling and deformations



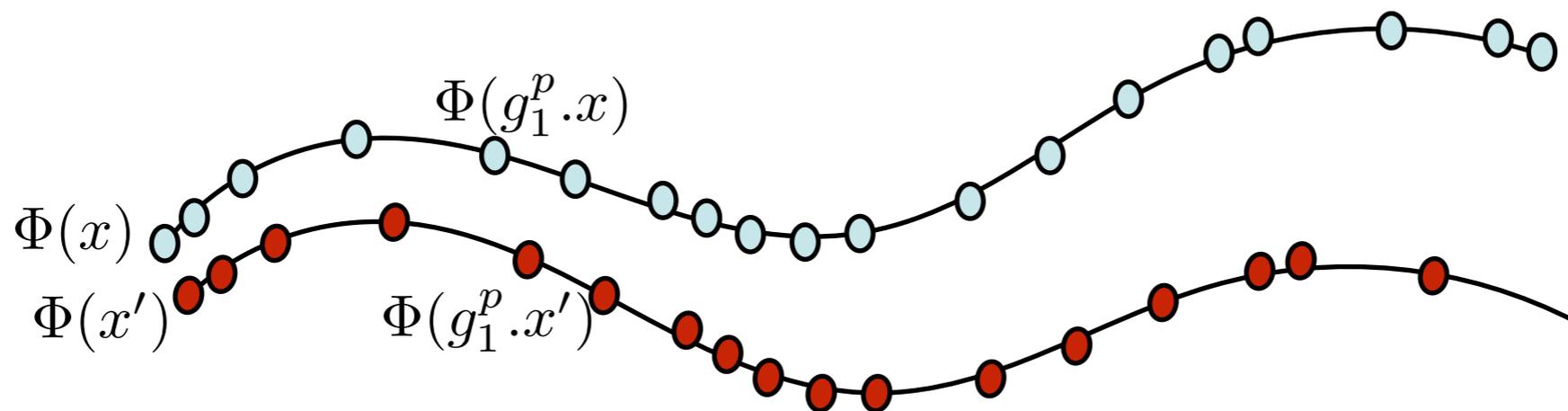
Group: $\mathbb{R} \times \text{Diff}(\mathbb{R})$

Linearize Symmetries

- A change of variable $\Phi(x)$ must linearize the orbits $\{g.x\}_{g \in G}$



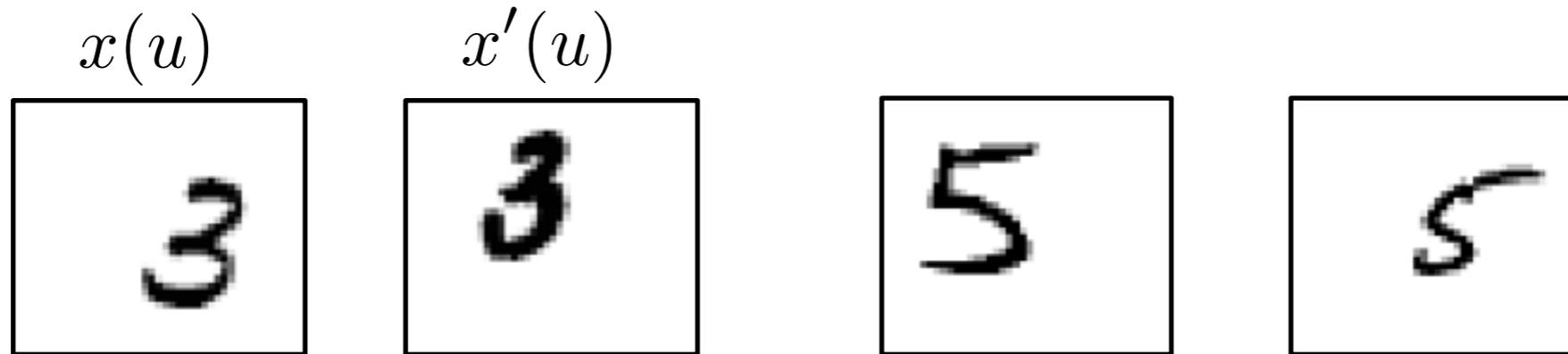
- Linearise symmetries with a change of variable $\Phi(x)$



- Lipschitz: $\forall x, g : \|\Phi(x) - \Phi(g.x)\| \leq C \|g\|$

Translation and Deformations

- Digit classification:



- Globally invariant to the translation group
- Locally invariant to small diffeomorphisms

Linearize small
diffeomorphisms:
 \Rightarrow Lipschitz regular



Video of Philipp Scott Johnson

- Invariance to translations:

$$g.x(u) = x(u - c) \quad \Rightarrow \quad \Phi(g.x) = \Phi(x) .$$

- Small diffeomorphisms: $g.x(u) = x(u - \tau(u))$

Metric: $\|g\| = \|\nabla\tau\|_\infty$ maximum scaling

Linearisation by Lipschitz continuity

$$\|\Phi(x) - \Phi(g.x)\| \leq C \|\nabla\tau\|_\infty .$$

- Discriminative change of variable:

$$\|\Phi(x) - \Phi(x')\| \geq C^{-1} |f(x) - f(x')|$$

- Fourier transform $\hat{x}(\omega) = \int x(t) e^{-i\omega t} dt$

$$x_c(t) = x(t - c) \Rightarrow \hat{x}_c(\omega) = e^{-ic\omega} \hat{x}(\omega)$$

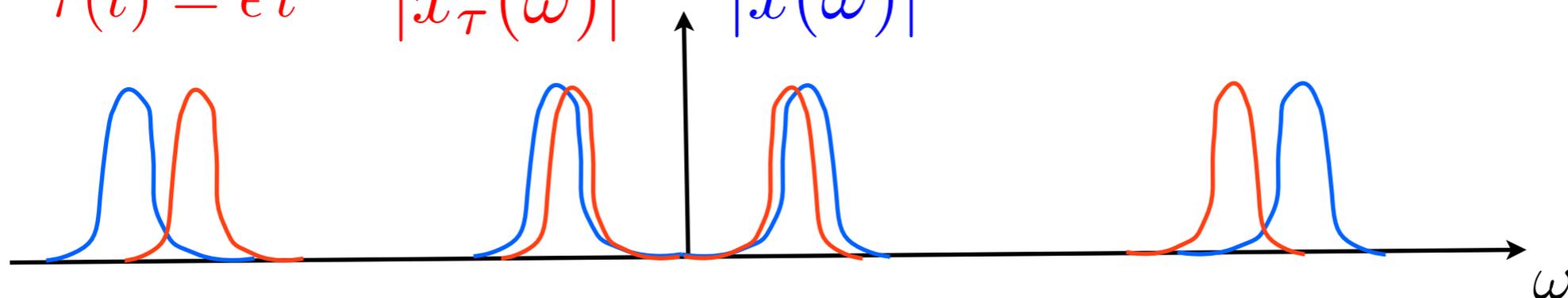
The modulus is invariant to translations:

$$\Phi(x) = |\hat{x}| = |\hat{x}_c|$$

- Instabilites to small deformations $x_\tau(t) = x(t - \tau(t))$:

$||\hat{x}_\tau(\omega)| - |\hat{x}(\omega)||$ is big at high frequencies

$$\tau(t) = \epsilon t \quad |\hat{x}_\tau(\omega)| \quad |\hat{x}(\omega)|$$



$$\Rightarrow |||\hat{x}| - |\hat{x}_\tau||| \gg \|\nabla \tau\|_\infty \|x\|$$

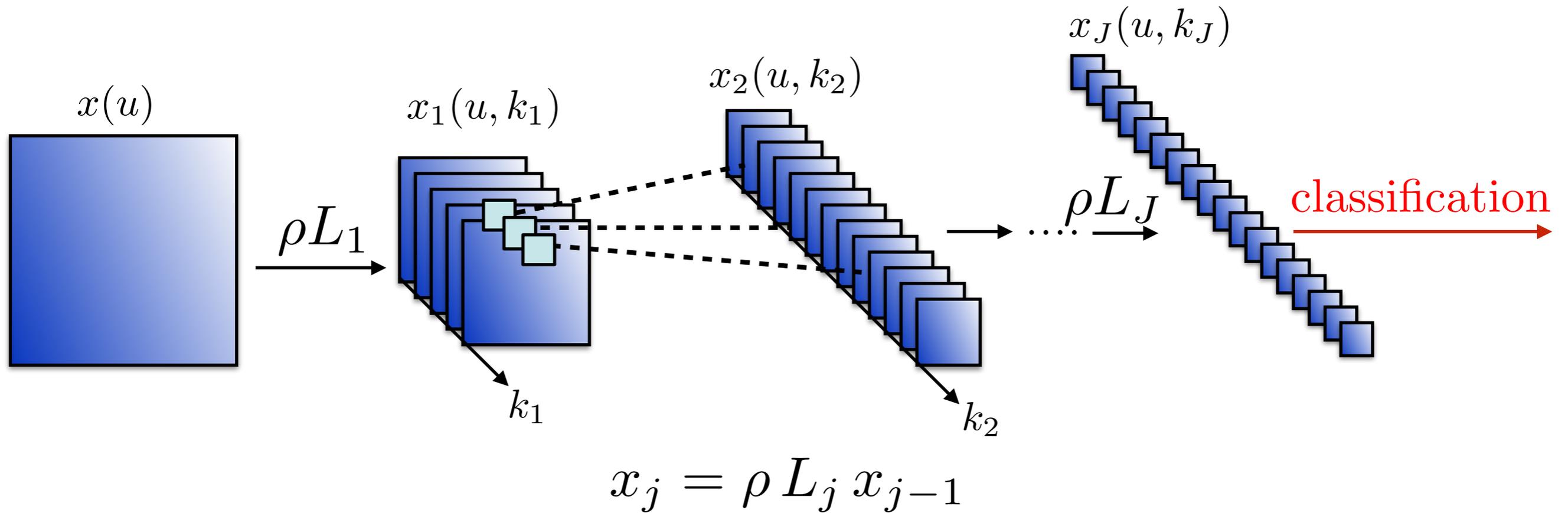
Deep Neural Network Mathematical Mysteries for High Dimensional Learning



Stéphane Mallat

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

Deep Convolutional Trees



L_j is composed of convolutions and subs samplings:

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

No channel communication: how far can we go ?

Why hierachical cascade ?

- Invariance to translations:

$$g.x(u) = x(u - c) \quad \Rightarrow \quad \Phi(g.x) = \Phi(x) .$$

- Small diffeomorphisms: $g.x(u) = x(u - \tau(u))$

Metric: $\|g\| = \|\nabla\tau\|_\infty$ maximum scaling

Linearisation by Lipschitz continuity

$$\|\Phi(x) - \Phi(g.x)\| \leq C \|\nabla\tau\|_\infty .$$

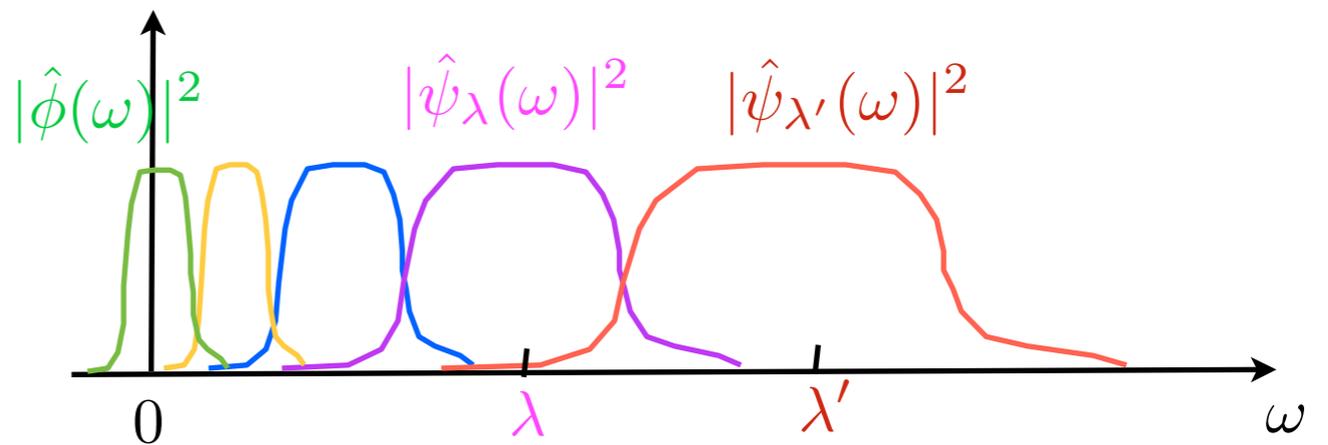
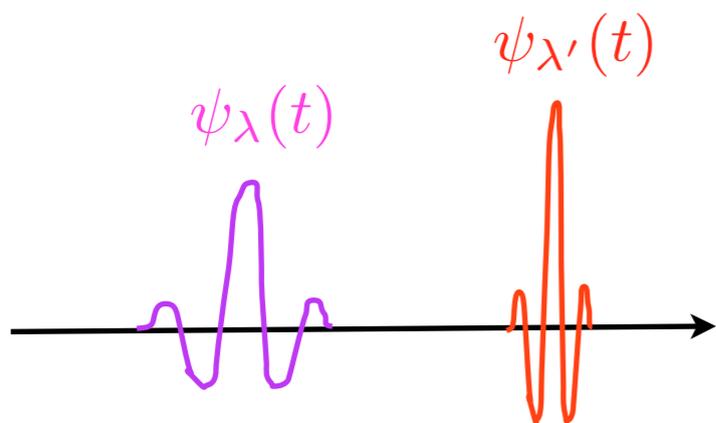
- Discriminative change of variable:

$$\|\Phi(x) - \Phi(x')\| \geq C^{-1} |f(x) - f(x')|$$

- Wavelet Scattering transform along translations
- Generation of textures and random processes
- Channel connections for more general groups
- Image and audio classification with small training sets
- Quantum chemistry
- Open problems

Multiscale 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 $\hat{\psi}_\lambda$

$$x \star \psi_\lambda(t) = \int x(u) \psi_\lambda(t - u) du \Rightarrow \widehat{x \star \psi_\lambda}(\omega) = \hat{x}(\omega) \hat{\psi}_\lambda(\omega)$$

- 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: $\|Wx\|^2 = \|x\|^2$.

Why Wavelets ?

- Wavelets are uniformly stable to deformations:

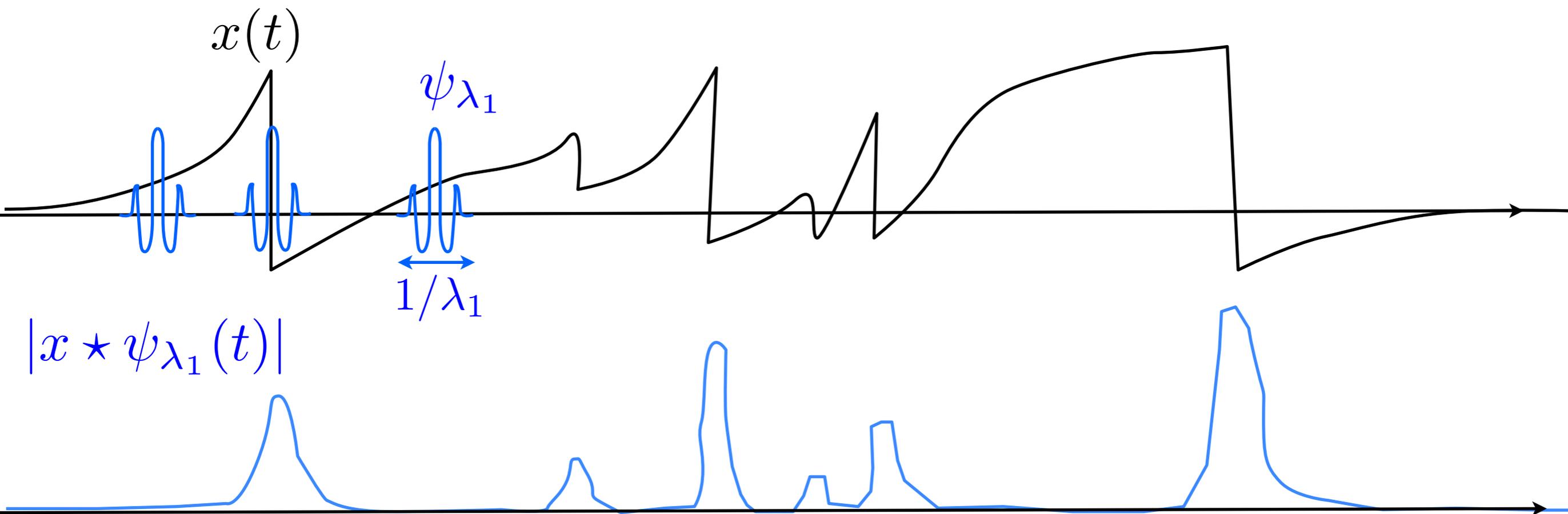
if $\psi_{\lambda,\tau}(t) = \psi_{\lambda}(t - \tau(t))$ then

$$\|\psi_{\lambda} - \psi_{\lambda,\tau}\| \leq C \sup_t |\nabla \tau(t)| .$$

- Wavelets separate multiscale information.
- Wavelets provide sparse representations.

Singular Functions

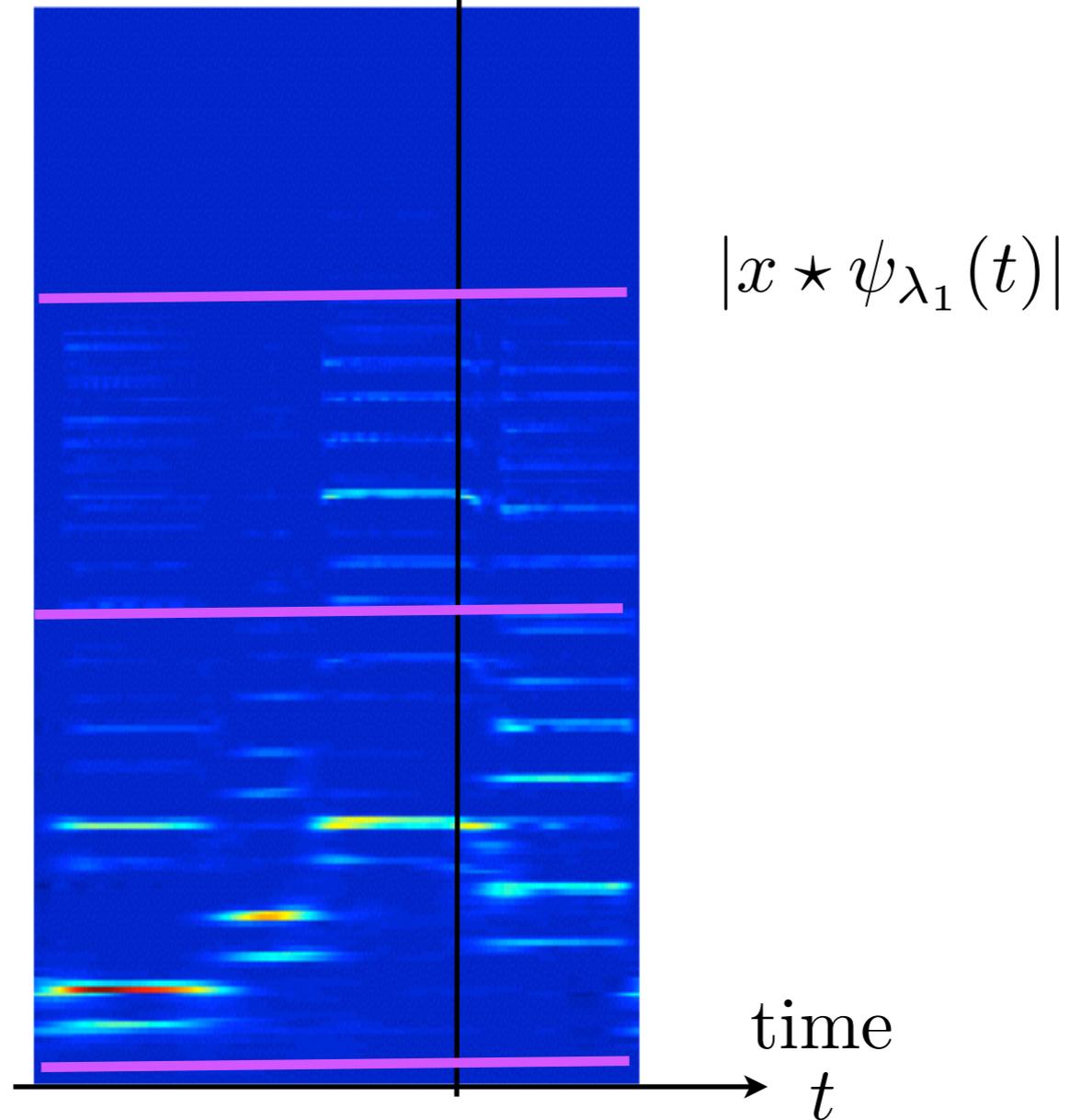
$$|x \star \psi_{\lambda_1}(t)| = \left| \int x(u) \psi_{\lambda_1}(t-u) du \right|$$



Time-Frequency Fibers

Wavelet transform modulus: $|W|$

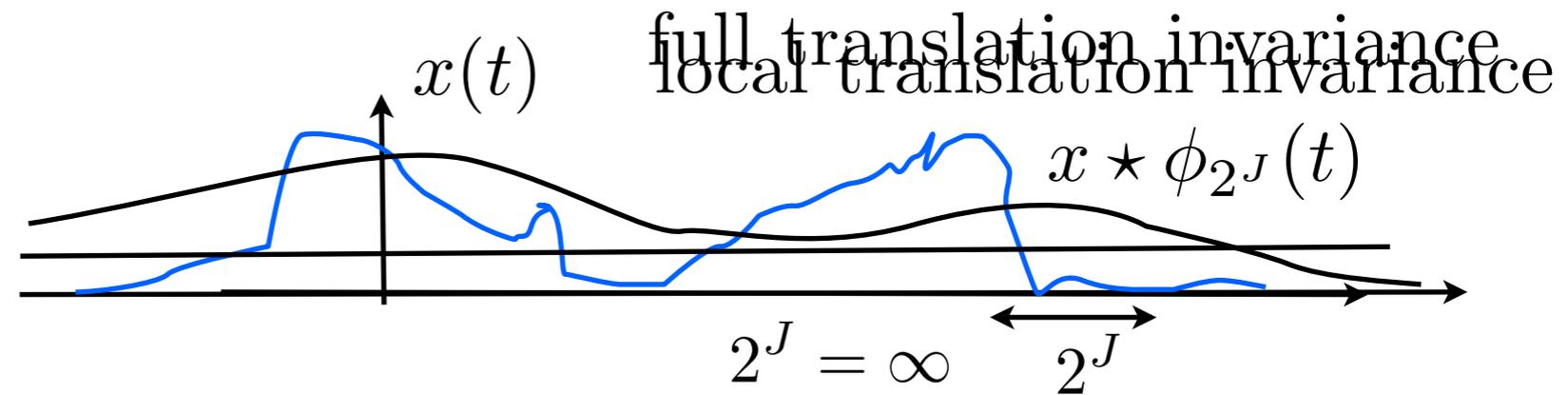
frequency
 $\log \omega = \lambda_1$



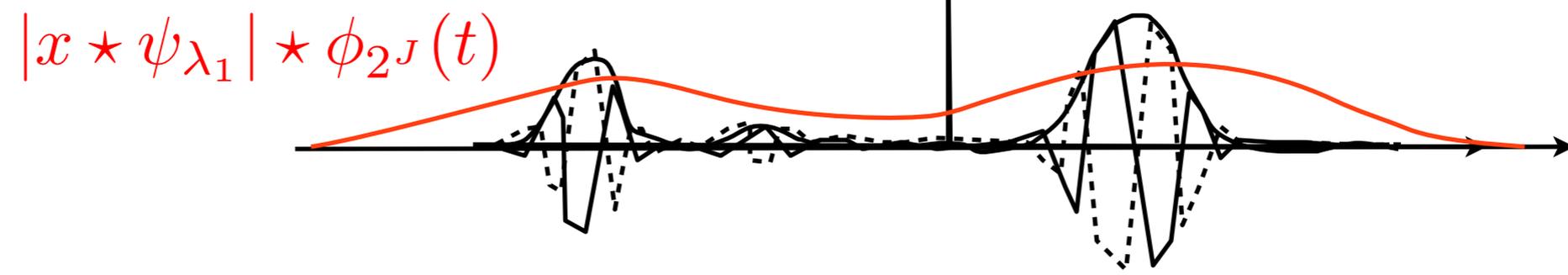
Wavelet Translation Invariance

First wavelet transform

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



Modulus improves invariance: $|x \star \psi_{\lambda_1}| \star \phi_{2^J}(t) = \sqrt{|x \star \psi_{\lambda_1}|^2 \star \phi_{2^J}(t)}$

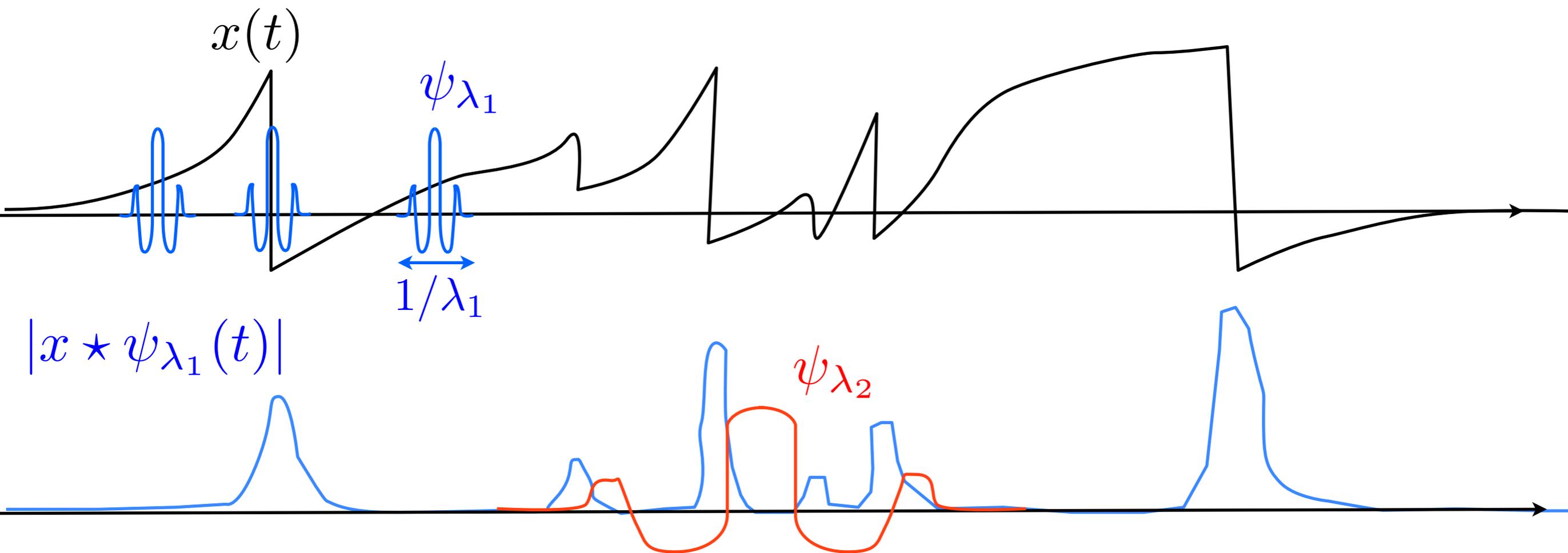


Second wavelet transform modulus

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

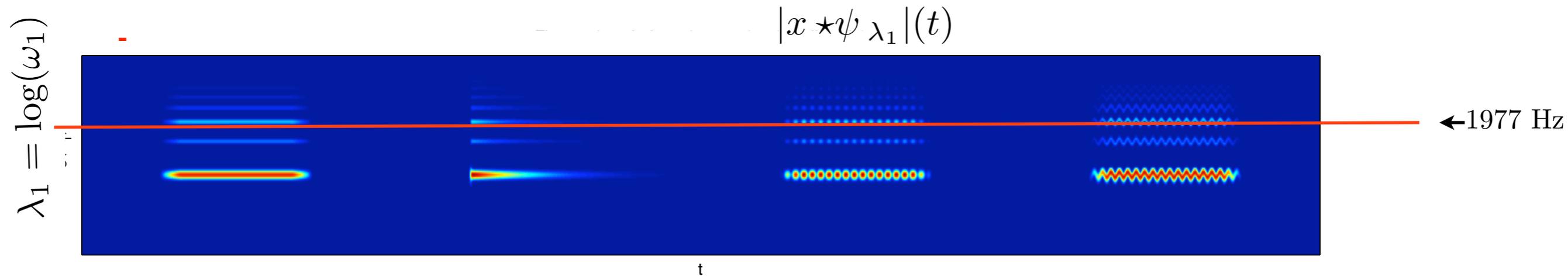
Singular Functions

$$|x \star \psi_{\lambda_1}(t)| = \left| \int x(u) \psi_{\lambda_1}(t - u) du \right|$$

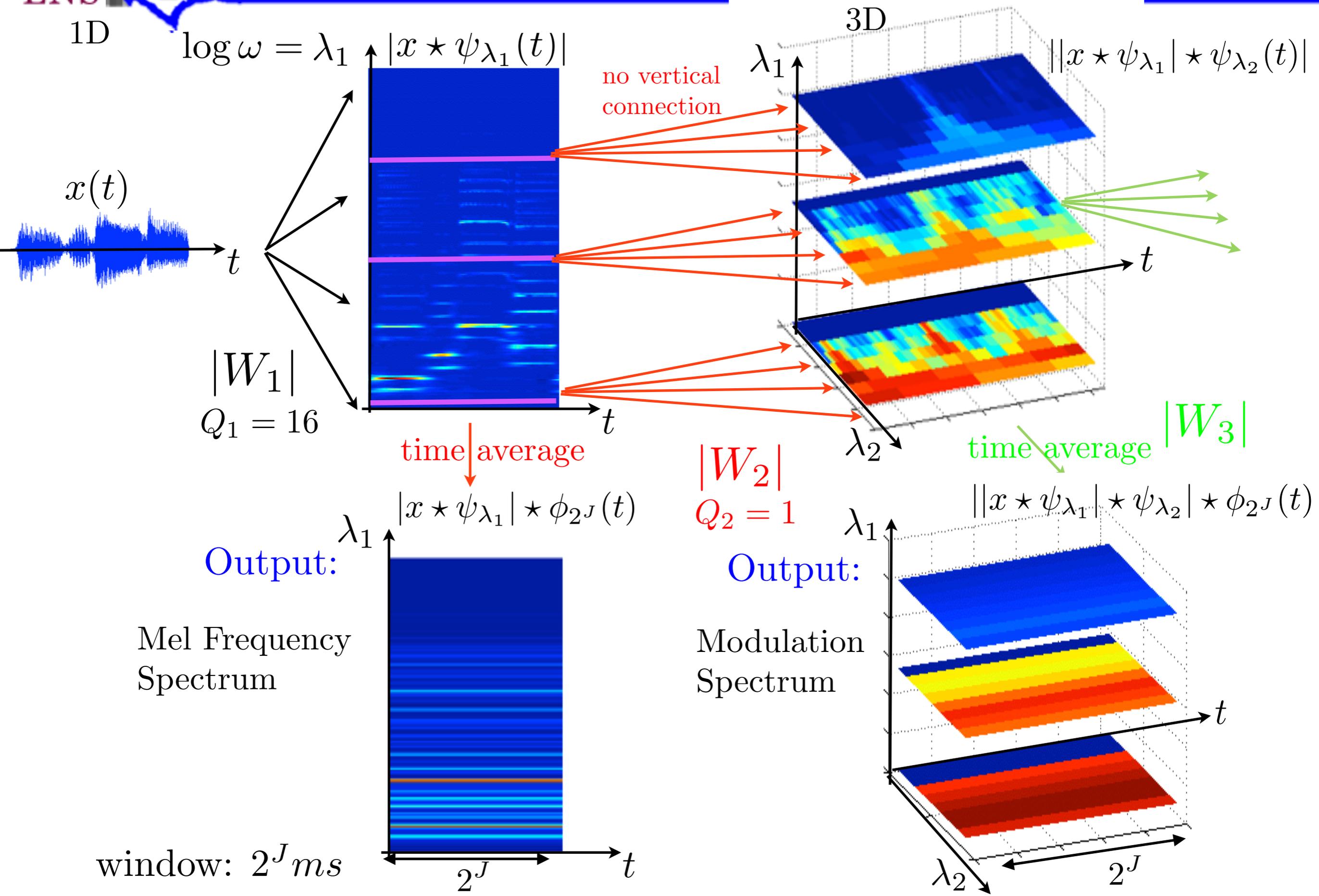


Amplitude Modulation

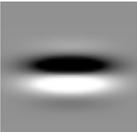
Harmonic sound: $x(t) = a(t) e \star h(t)$ with varying $a(t)$



ScatteringConvolution Network



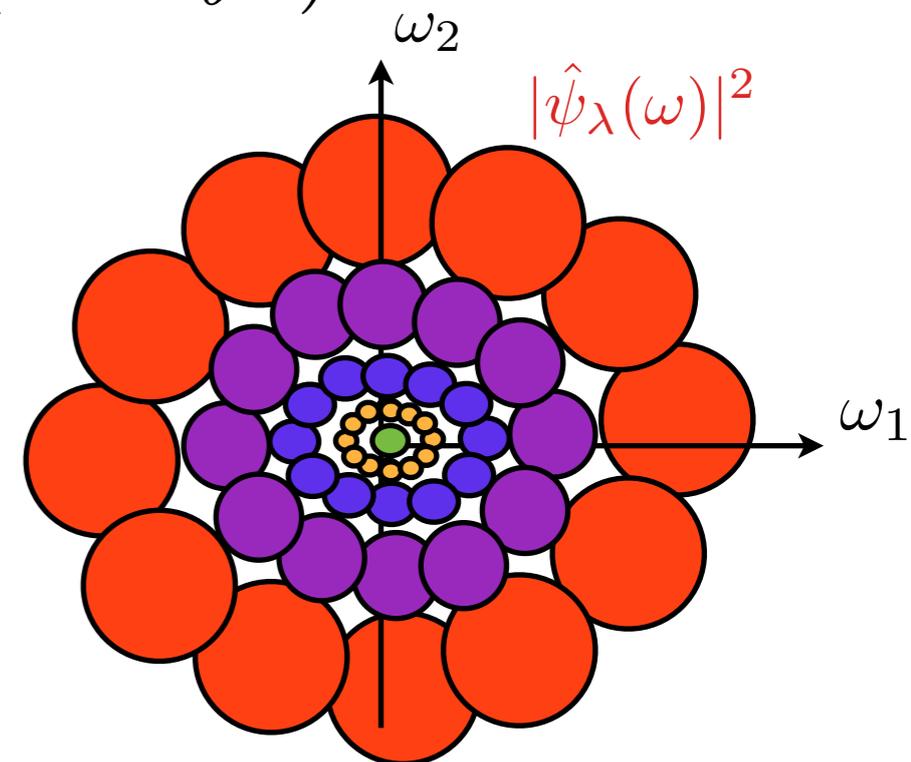
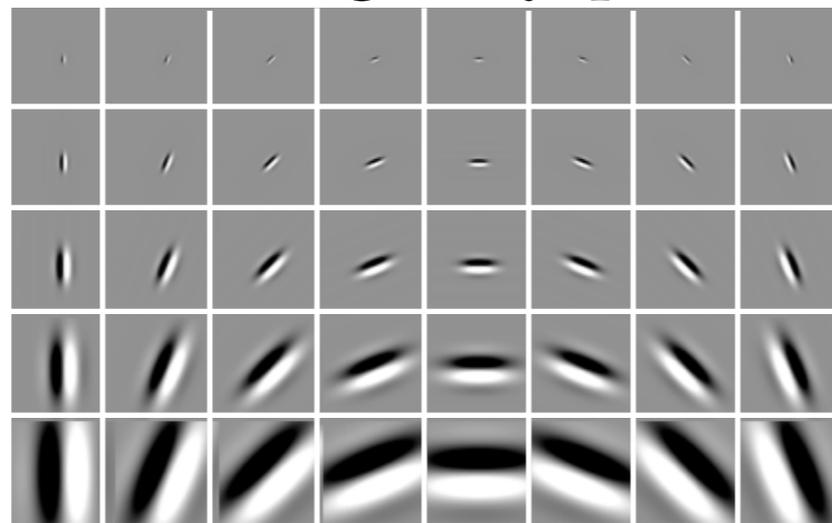
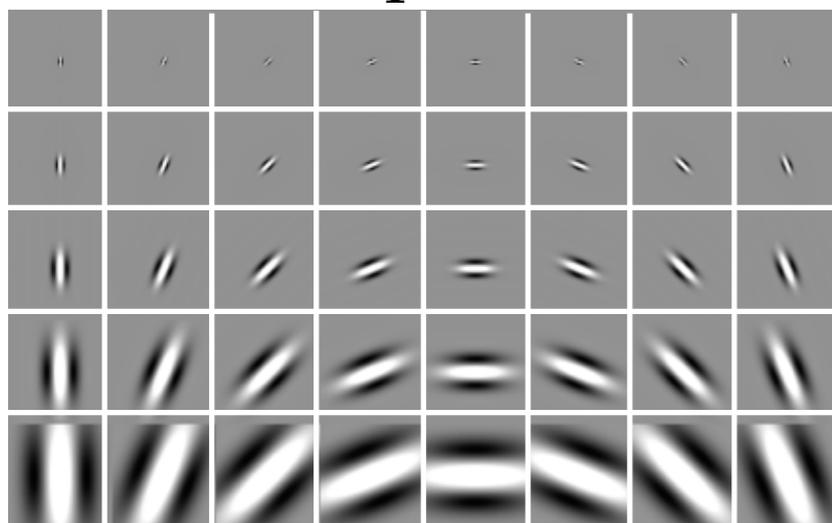
Scale separation with Wavelets

- Wavelet filter $\psi(u)$:  + i 

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

real parts

imaginary parts



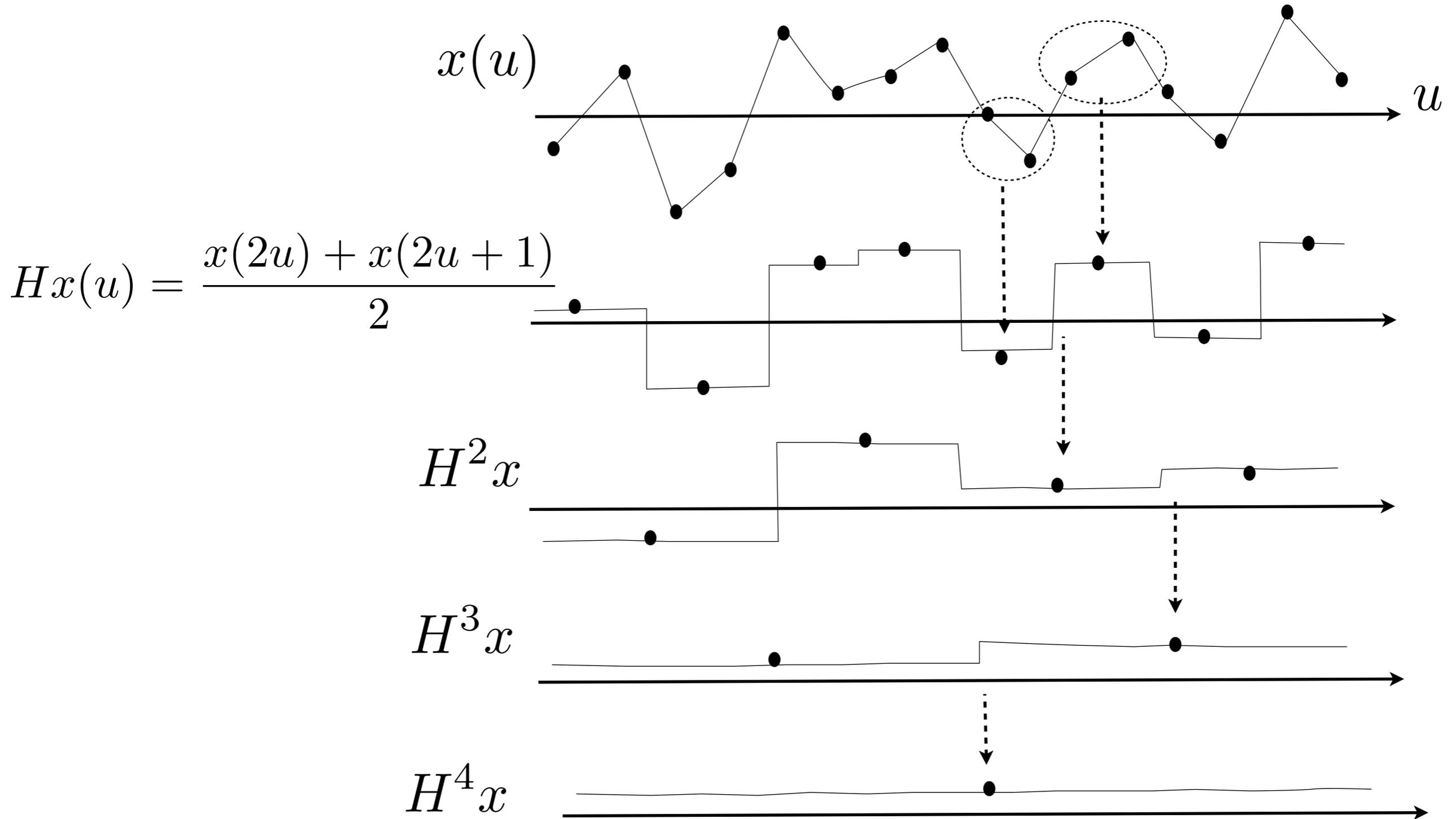
$$x \star \psi_{2^j, \theta}(u) = \int x(v) \psi_{2^j, \theta}(u - v) dv$$

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

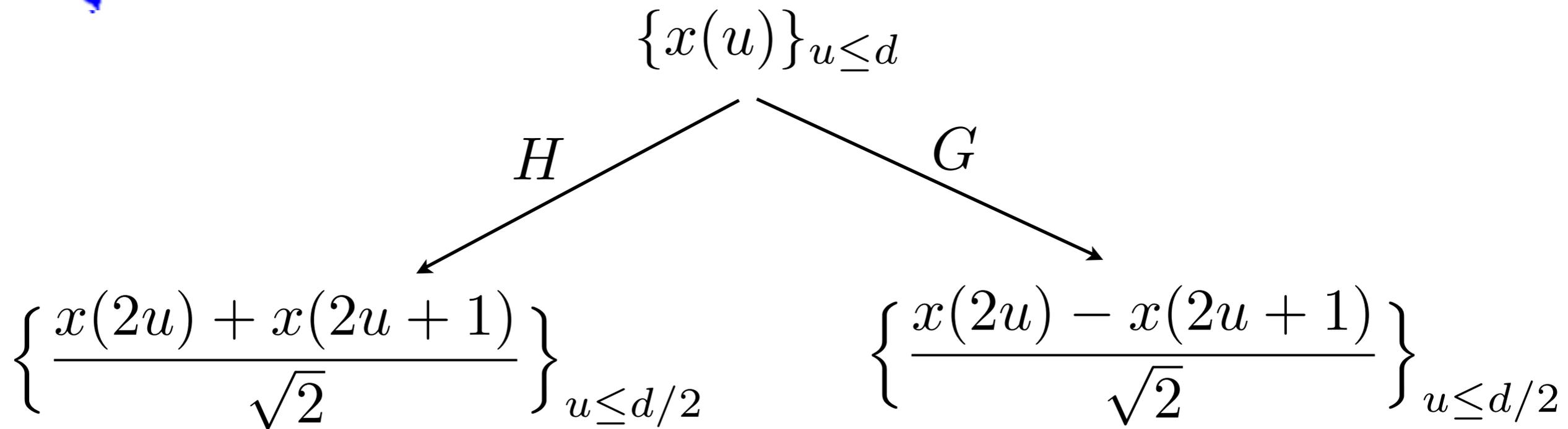
Preserves norm: $\|Wx\|^2 = \|x\|^2$.

Averaging Pyramid

- Multiscale averaging by cascade of pair averaging:



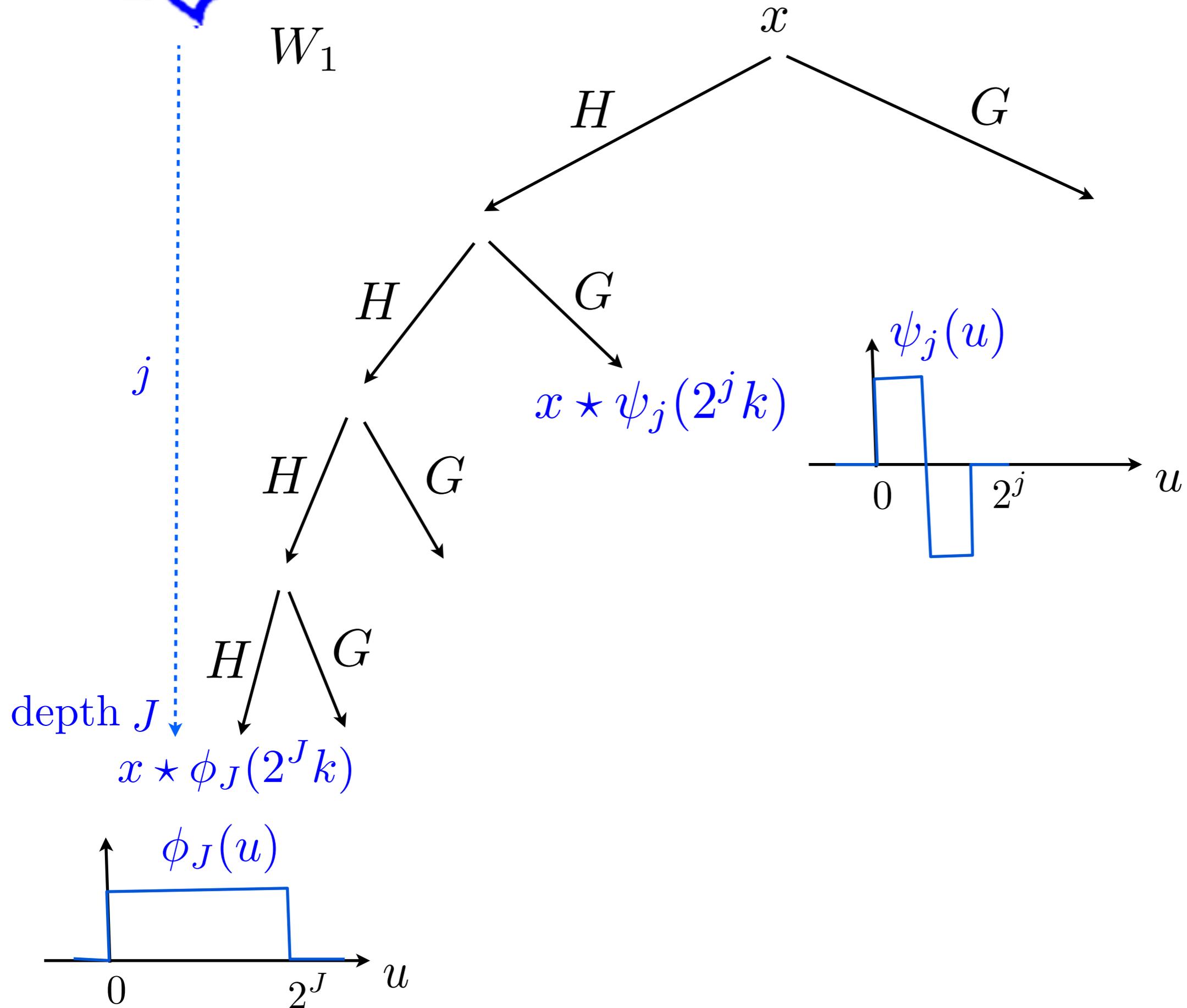
Haar Filtering



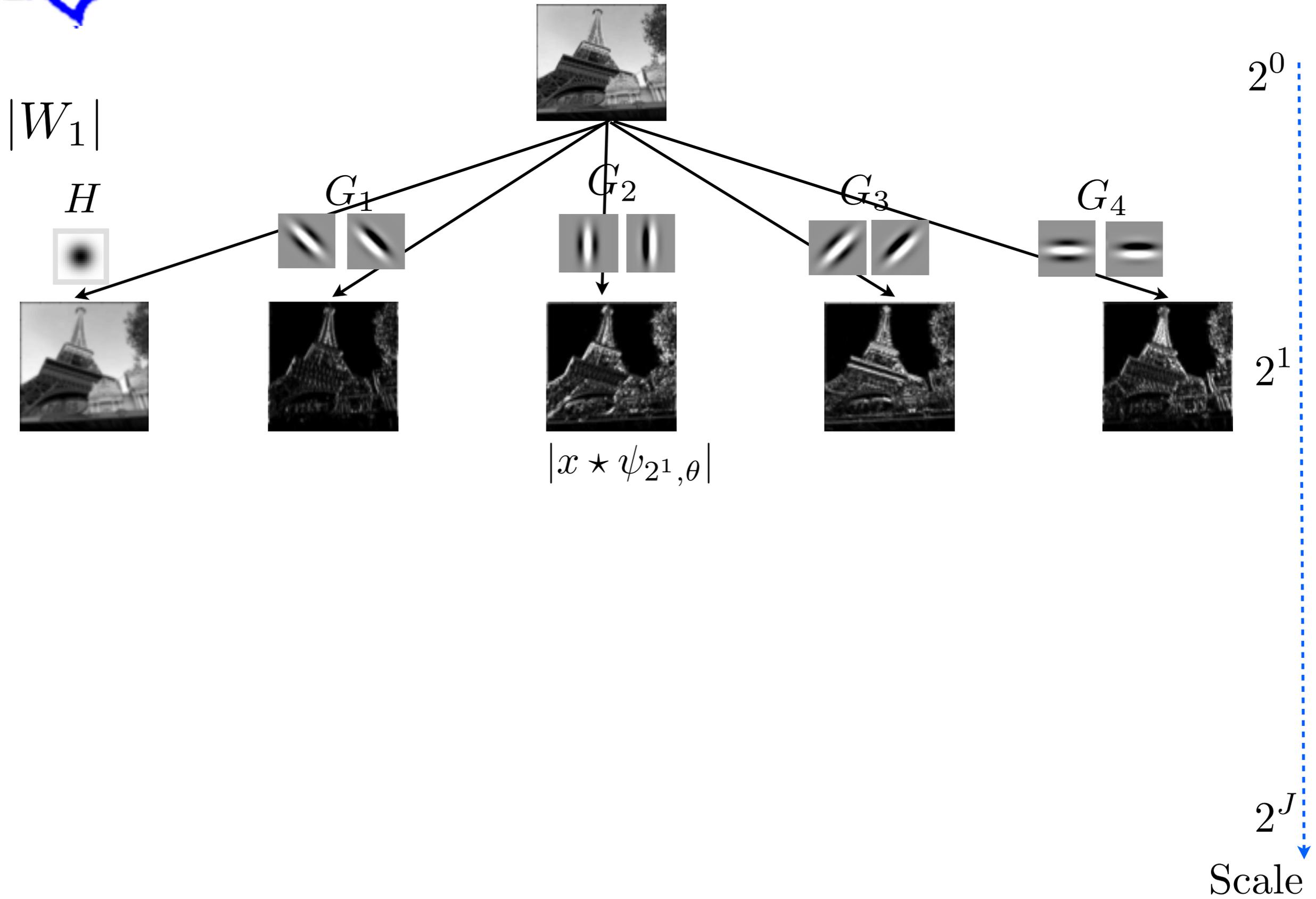
$$Hx(u) = x \star h(2u) \quad \text{and} \quad Gx(u) = x \star g(2u)$$

where h is a low frequency and g is a high frequency filter.

Haar Wavelet Transform



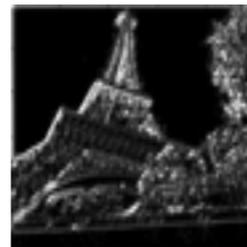
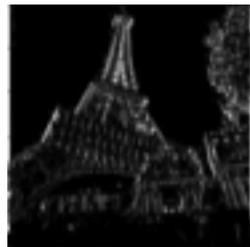
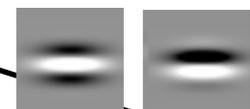
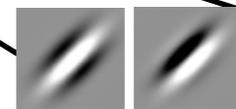
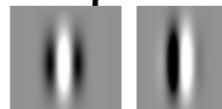
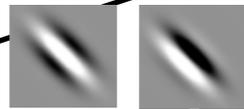
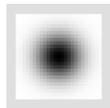
Fast Wavelet Filter Bank



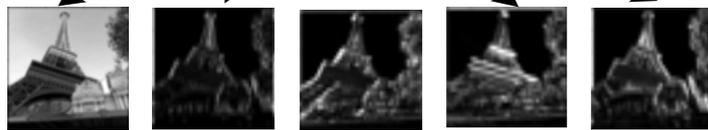
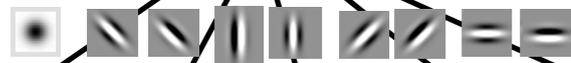
Wavelet Filter Bank

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

$$|W_1|$$


 $x(u)$
 2^0

 2^1

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



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

 2^2

If $u \geq 0$ then $\rho(u) = u$
 ρ has no effect after an averaging.



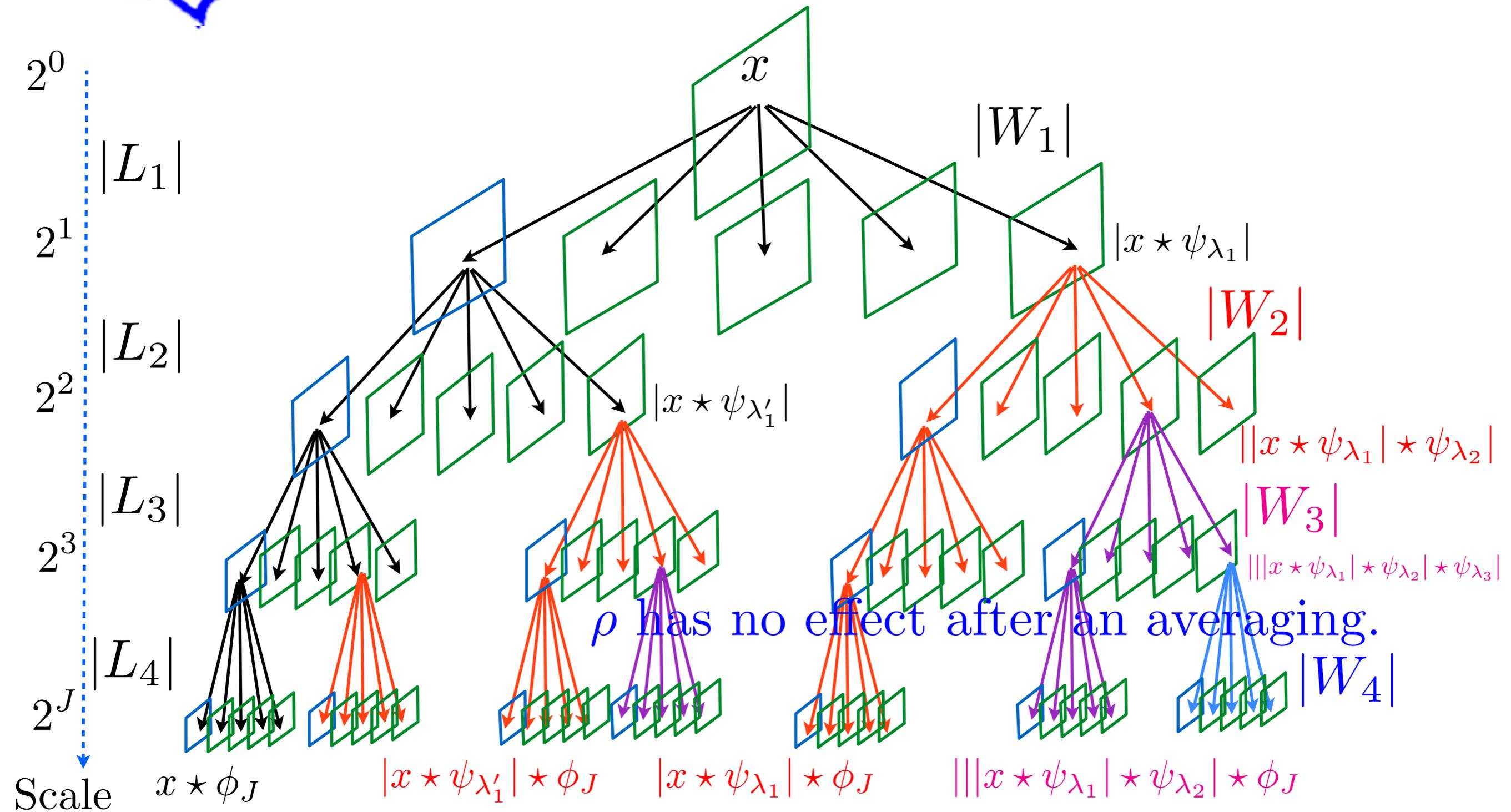
$$|x \star \psi_{2^j, \theta}|$$

 2^j

Scale

- Sparse representation

Wavelet Convolution Network Tree



$$S_4 x = |L_4| |L_3| |L_2| |L_1| x = |W_4| |W_3| |W_2| |W_1| x$$

Contraction

$$Wx = \begin{pmatrix} x \star \phi(t) \\ x \star \psi_\lambda(t) \end{pmatrix}_{t,\lambda} \quad \text{is linear and } \|Wx\| = \|x\|$$

$$\rho(u) = |u|$$

$$|W|x = \begin{pmatrix} x \star \phi(t) \\ |x \star \psi_\lambda(t)| \end{pmatrix}_{t,\lambda} \quad \text{is non-linear}$$

- it is contractive $\| |W|x - |W|y \| \leq \|x - y\|$

because for $(a, b) \in \mathbb{C}^2$ $\| |a| - |b| \| \leq |a - b|$

- it preserves the norm $\| |W|x \| = \|x\|$

Scattering Properties

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

~~Lemma:~~ $\|x\|_{W_k, D_\tau} \leq C' \|\nabla \tau\|_\infty \|x\|_{W_k, D_\tau}$

Theorem: *For appropriate wavelets, a scattering is*

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

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

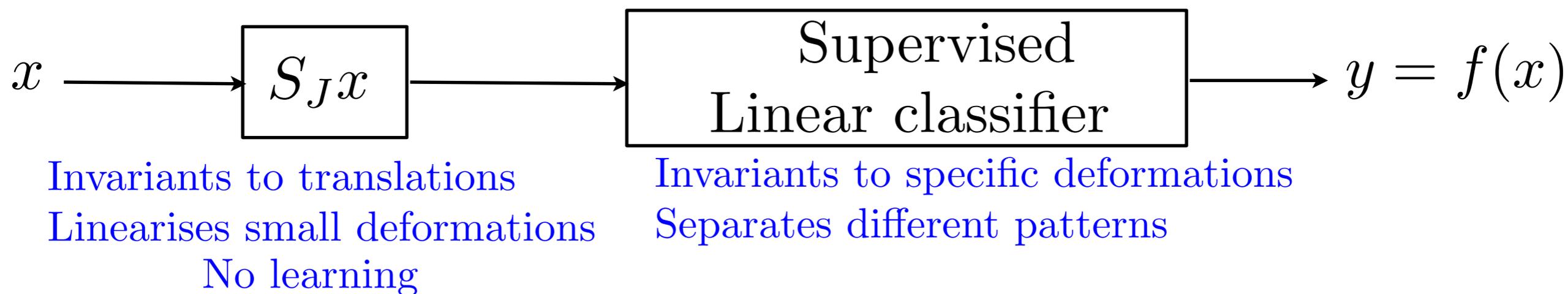
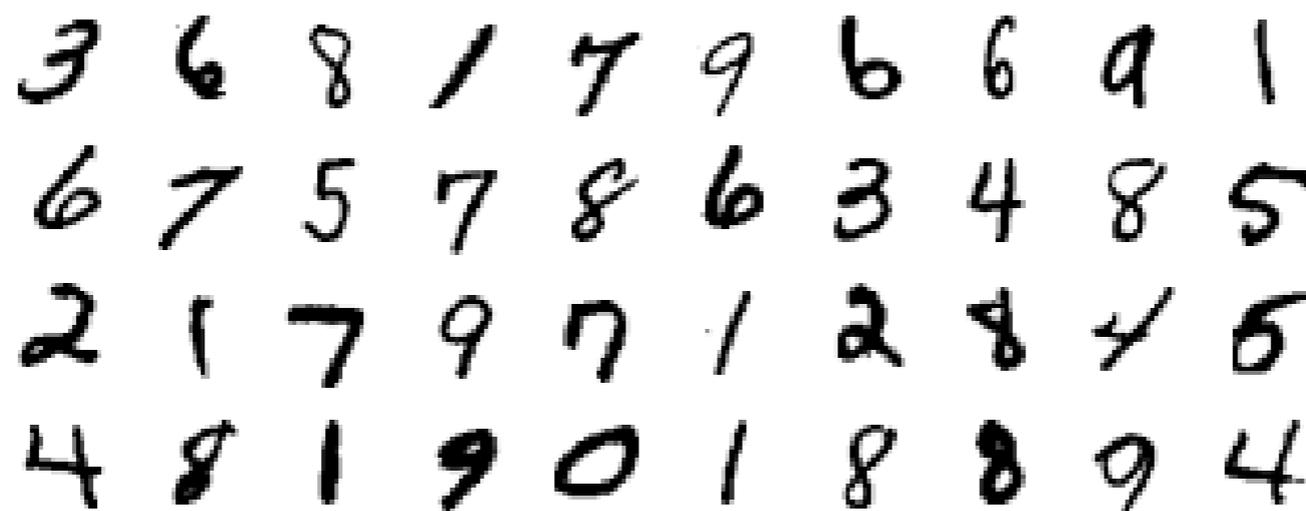
translations invariance and deformation stability:

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

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

Digit Classification: MNIST

Joan Bruna



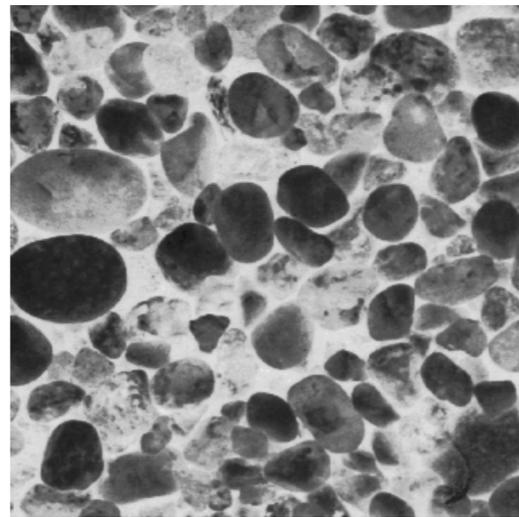
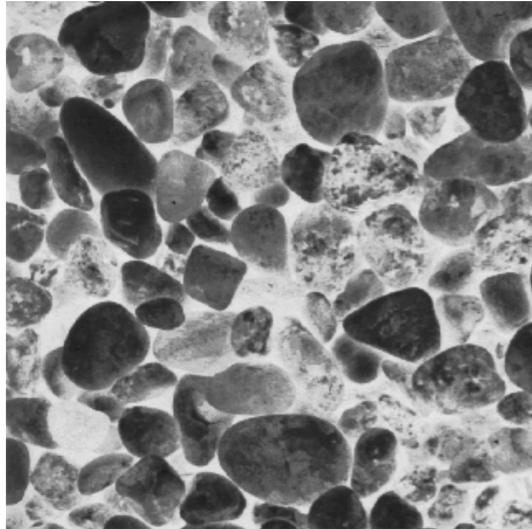
Classification Errors

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

LeCun et. al.

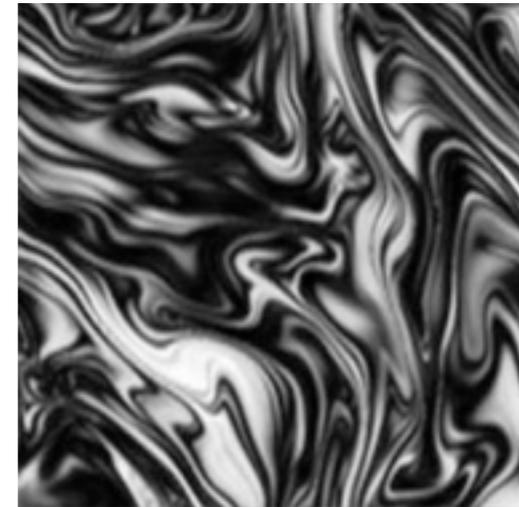
Classification of Stationary Textures

Ω_1



2D Turbulence

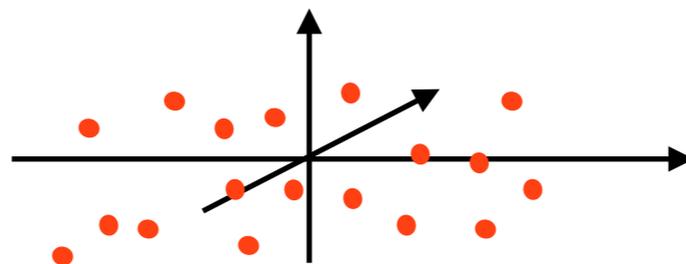
Ω_2



- What stochastic models ?

Non Gaussian with long-range dependance.

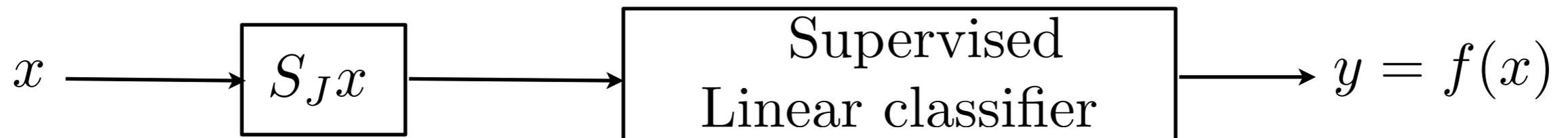
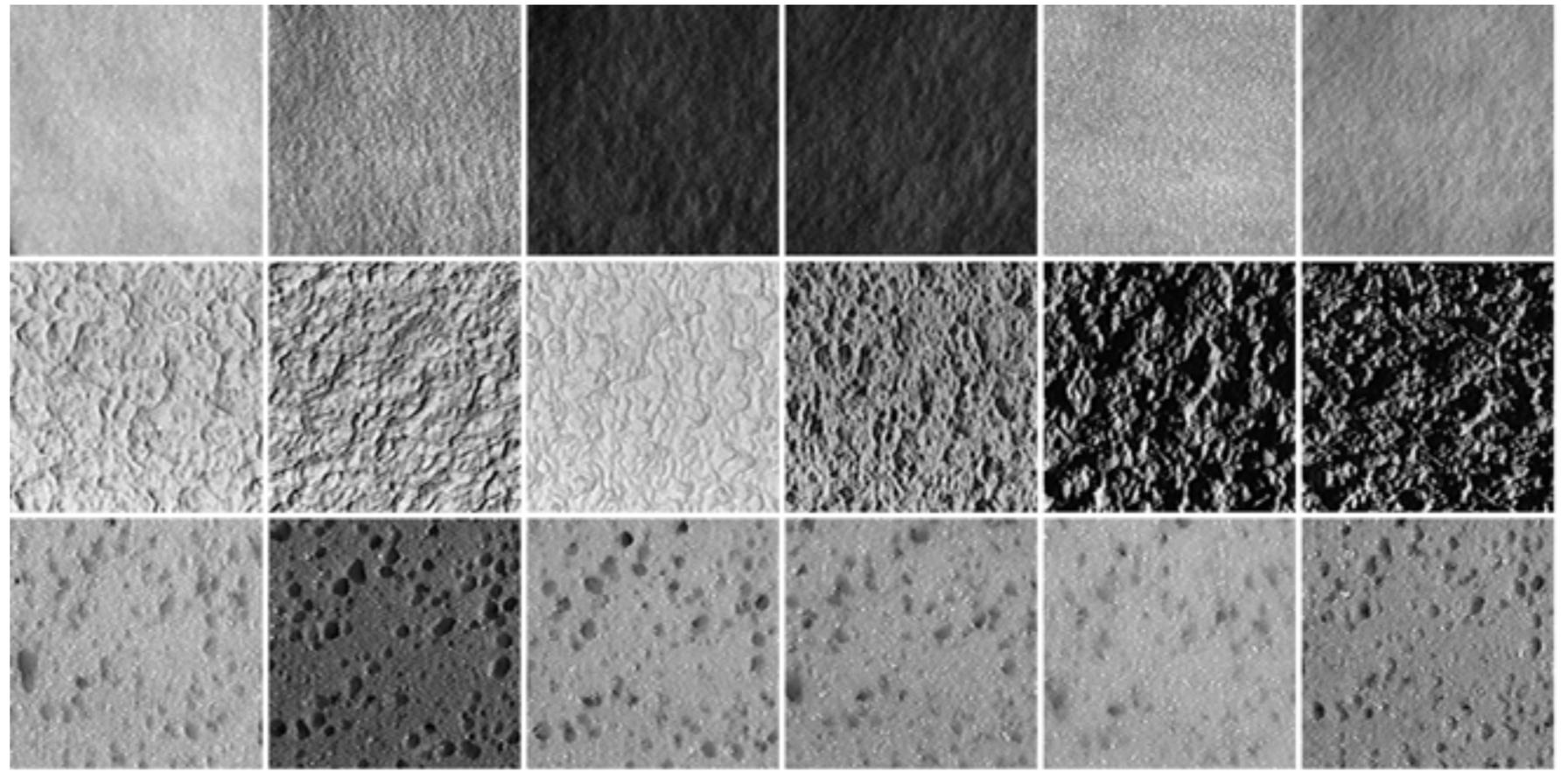
- Can we "Gaussianize" (linearize) such distributions in a reduced dimensional space ?



Classification of Textures

J. Bruna

CUREt database



Classification Errors

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

Scattering Moments of Processes

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

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

$J \rightarrow \infty$

Central limit theorem
with "weak" ergodicity conditions

J. Bruna

Gaussian distribution: $\mathcal{N}(\mathbb{E}(SX), \Sigma_J \rightarrow 0)$

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

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

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

$J \rightarrow \infty$

Central limit theorem
with "weak" ergodicity conditions

Gaussian distribution: $\mathcal{N}(\mathbb{E}(S X), \Sigma_J \rightarrow 0)$

- Reconstruction: compute \tilde{X} which minimises

$$\|S_J \tilde{X} - S_J X\|^2$$

- Gradient descent

Representation of Audio Textures

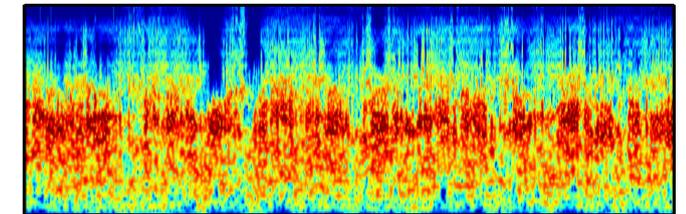
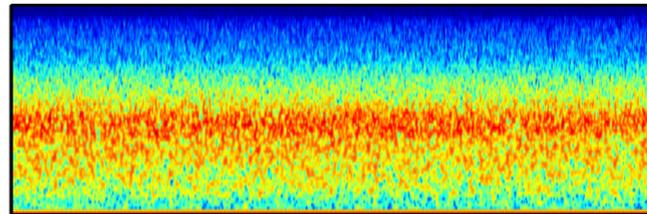
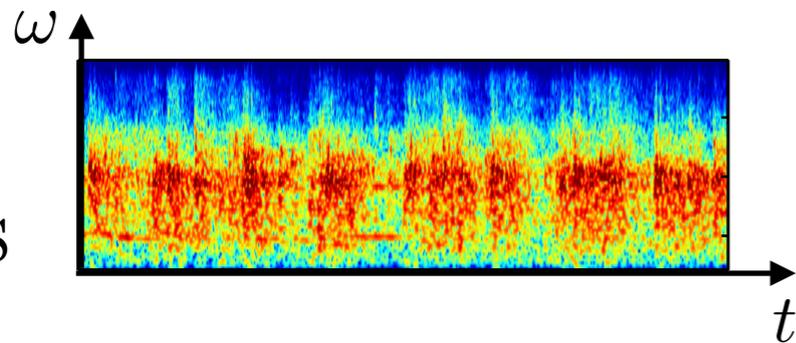
Joan Bruna

Original

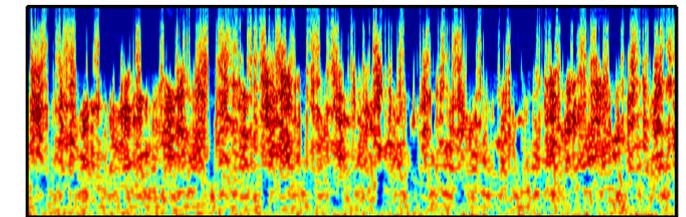
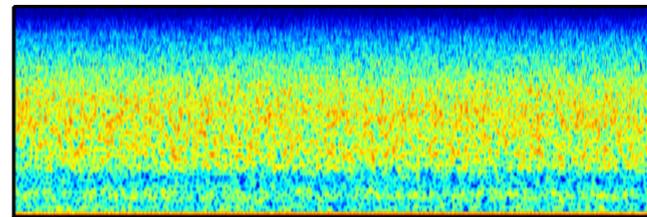
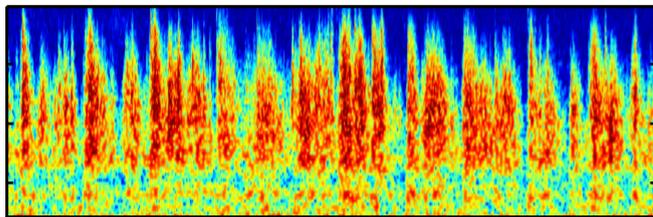
Gaussian
in time

Gaussian
in scattering

Applauds



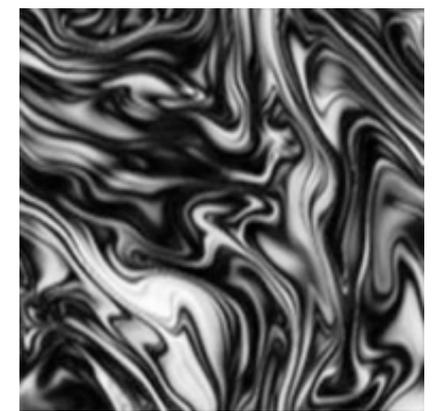
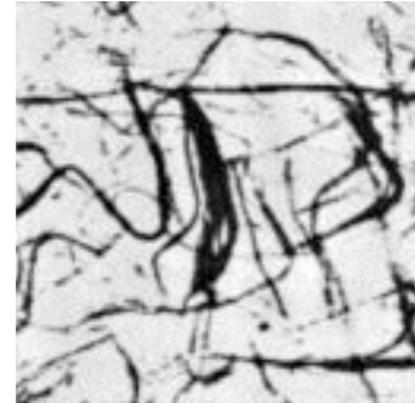
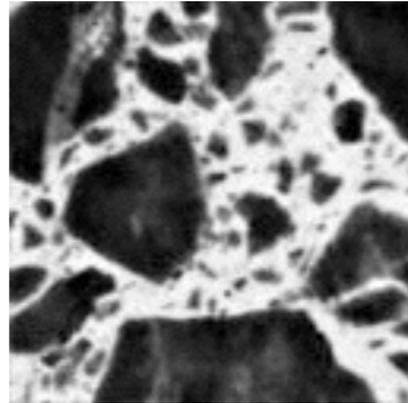
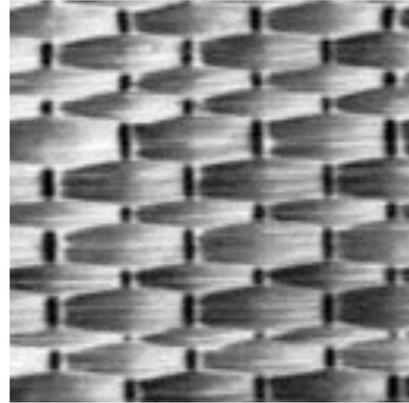
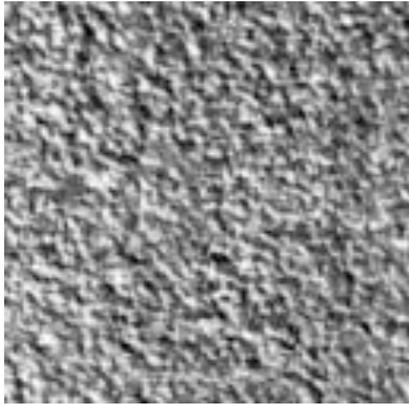
Paper



Cocktail Party

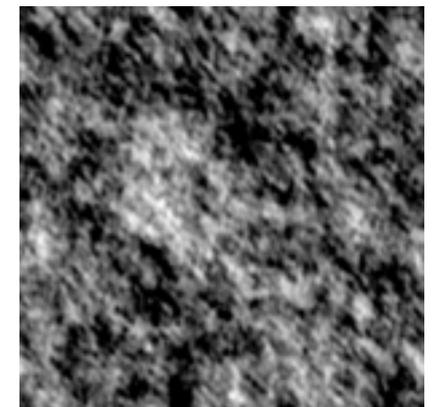
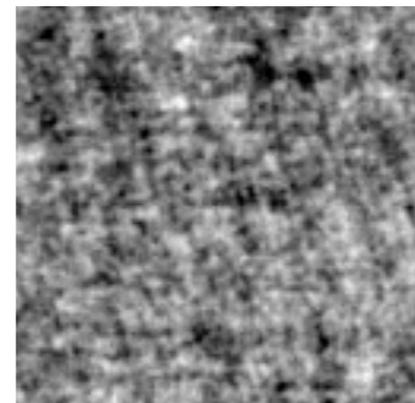
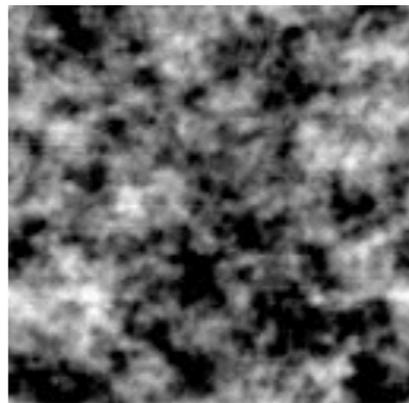
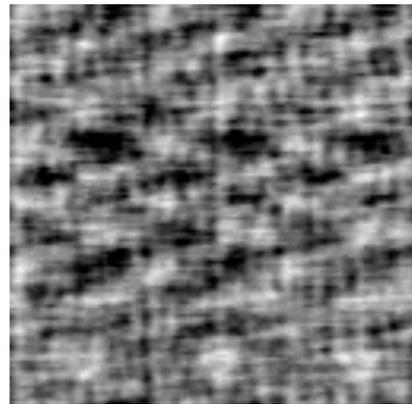
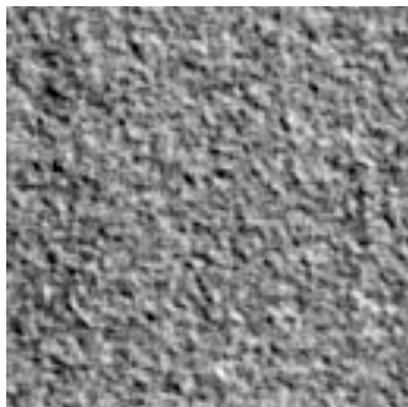
Joan Bruna

Textures of N pixels



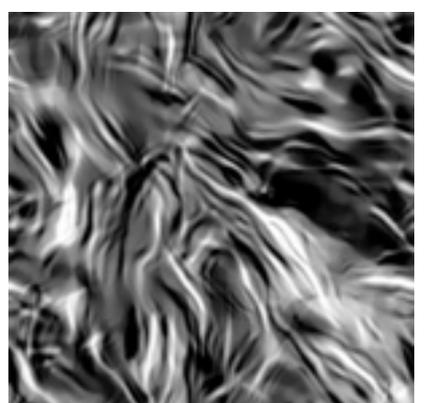
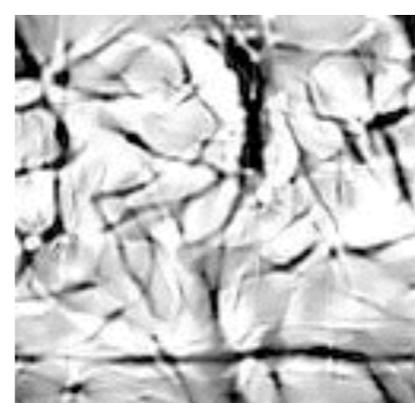
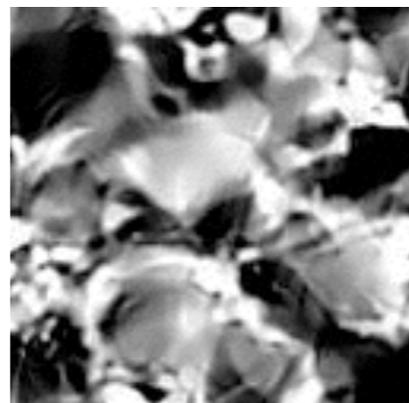
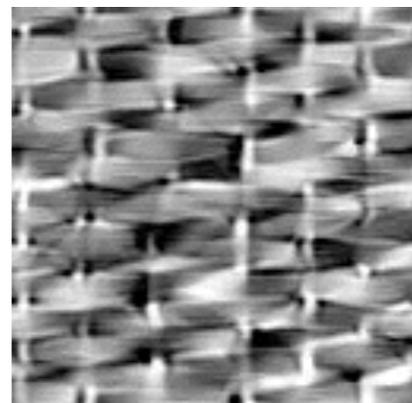
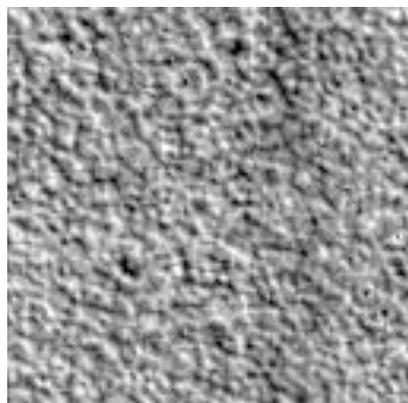
2D Turbulence

Gaussian process model with N second order moments

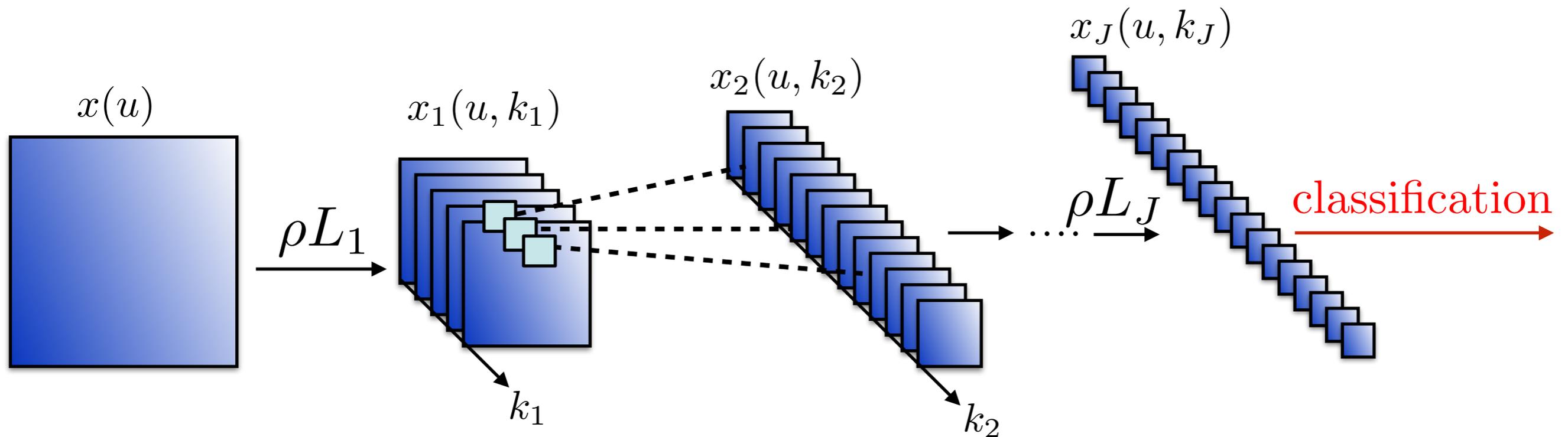


Second order Gaussian Scattering: $O(\log N^2)$ moments

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



Deep Convolutional Trees



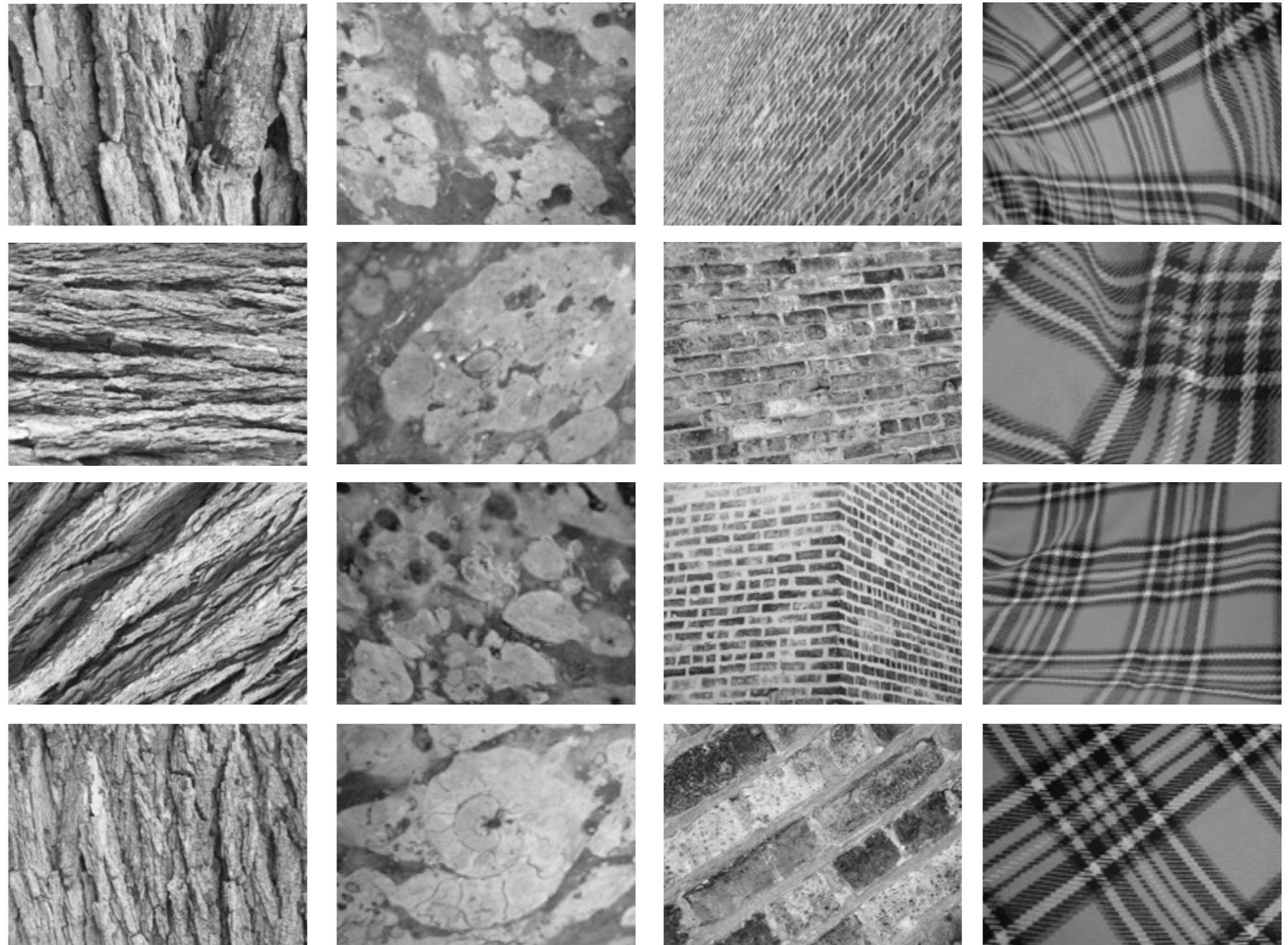
$$x_j = \rho L_j x_{j-1}$$

L_j is composed of convolutions and subs samplings:

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

No channel communication: what limitations ?

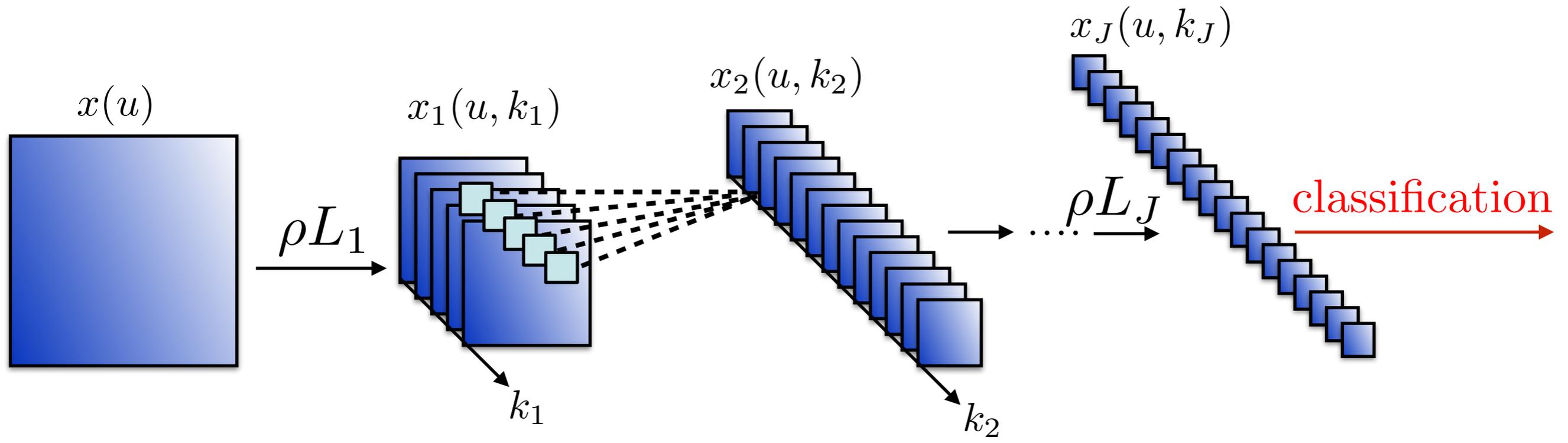
UIUC database:
25 classes



Scattering classification errors

Training	Scat. Translation
20	20 %

Deep Convolutional Networks



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

- L_j is a linear combination of convolutions and subsampling:

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

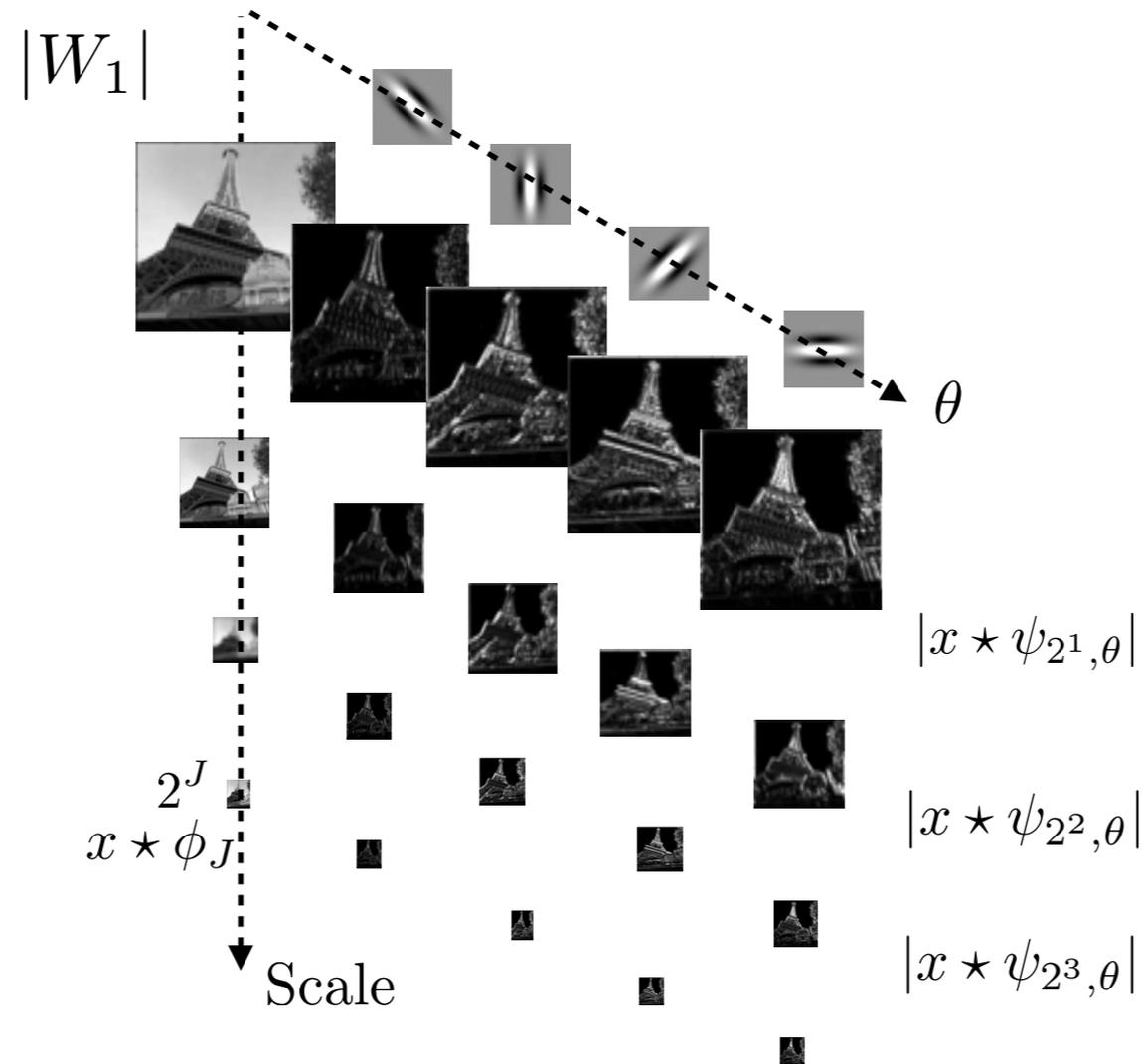
sum across channels

What is the role of channel connections ?

Linearize other symmetries beyond translations.

Rotation Invariance

- Channel connections linearize other symmetries.



- Invariance to rotations are computed by convolutions along the rotation variable θ with wavelet filters.
 \Rightarrow invariance to rigid movements.

Extension to Rigid Movements

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

$$x(r_\alpha(u, \mathbf{x}(u))) \longrightarrow \boxed{|W_1|} \longrightarrow x_j(r_\alpha(\theta, \mathbf{x}, \psi_2, \alpha, \theta)(u))$$

\downarrow
 $\int x(u) du$

Extension to Rigid Movements

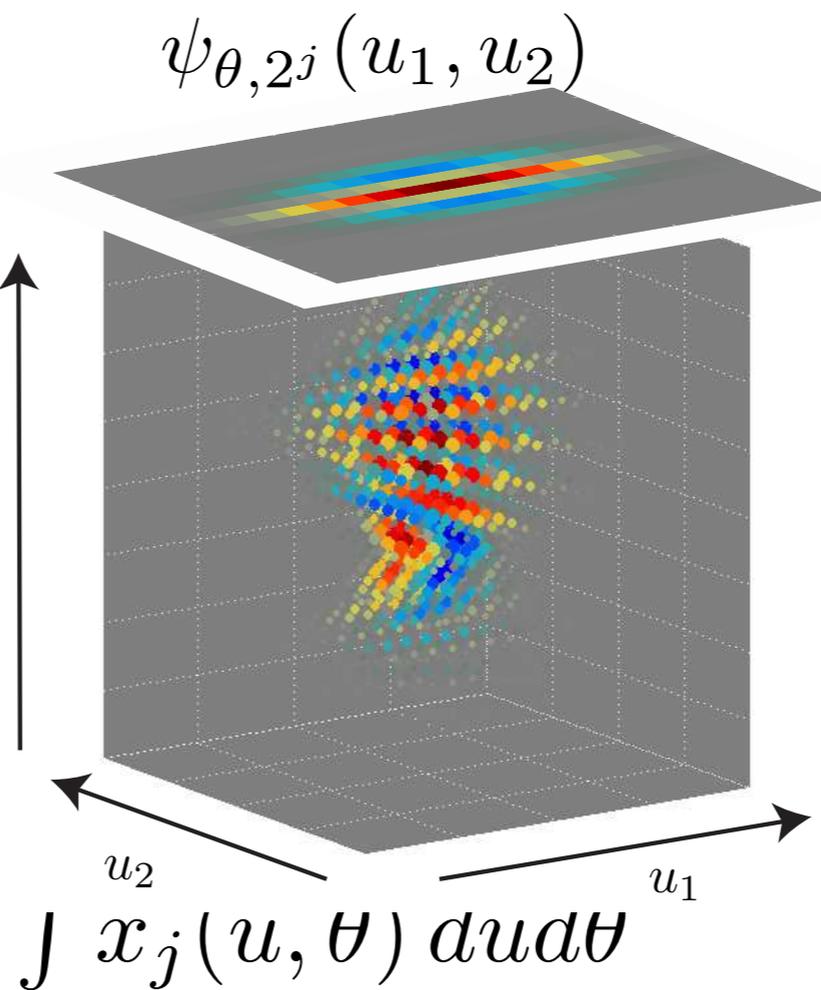
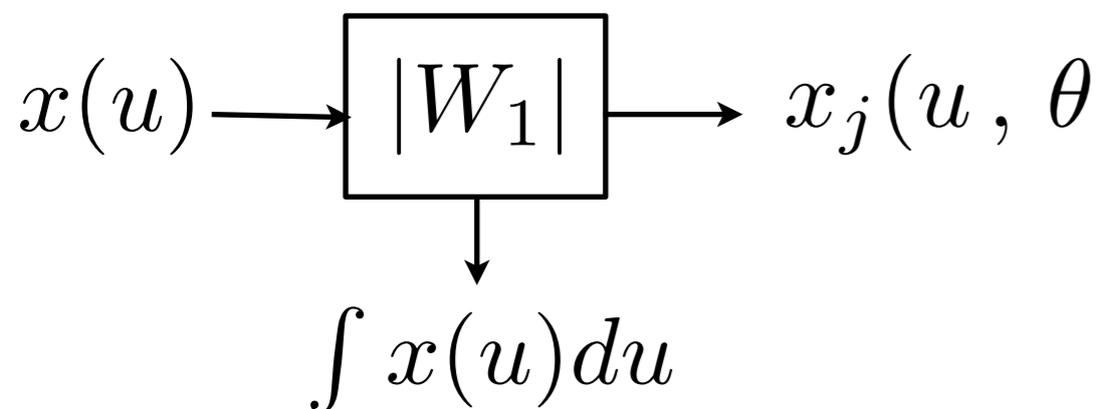
Laurent Sifre

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

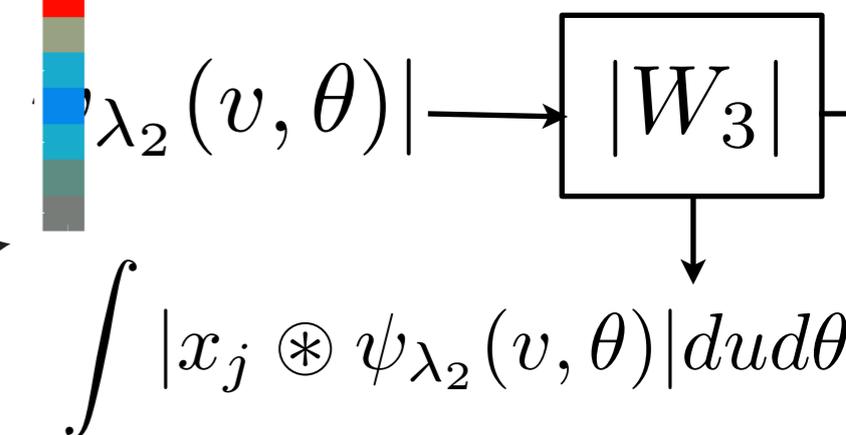
$$x \circledast \psi_{\lambda}(u, \theta) = \int_0^{2\pi} \left(\int_{\mathbb{R}^2} x(u', \theta') \psi_{\theta, 2^j}(r_{-\theta'}(u - u')) \right) \psi_{2^k}(\theta - \theta') d\theta' dt'$$

- Scattering on rigid movements

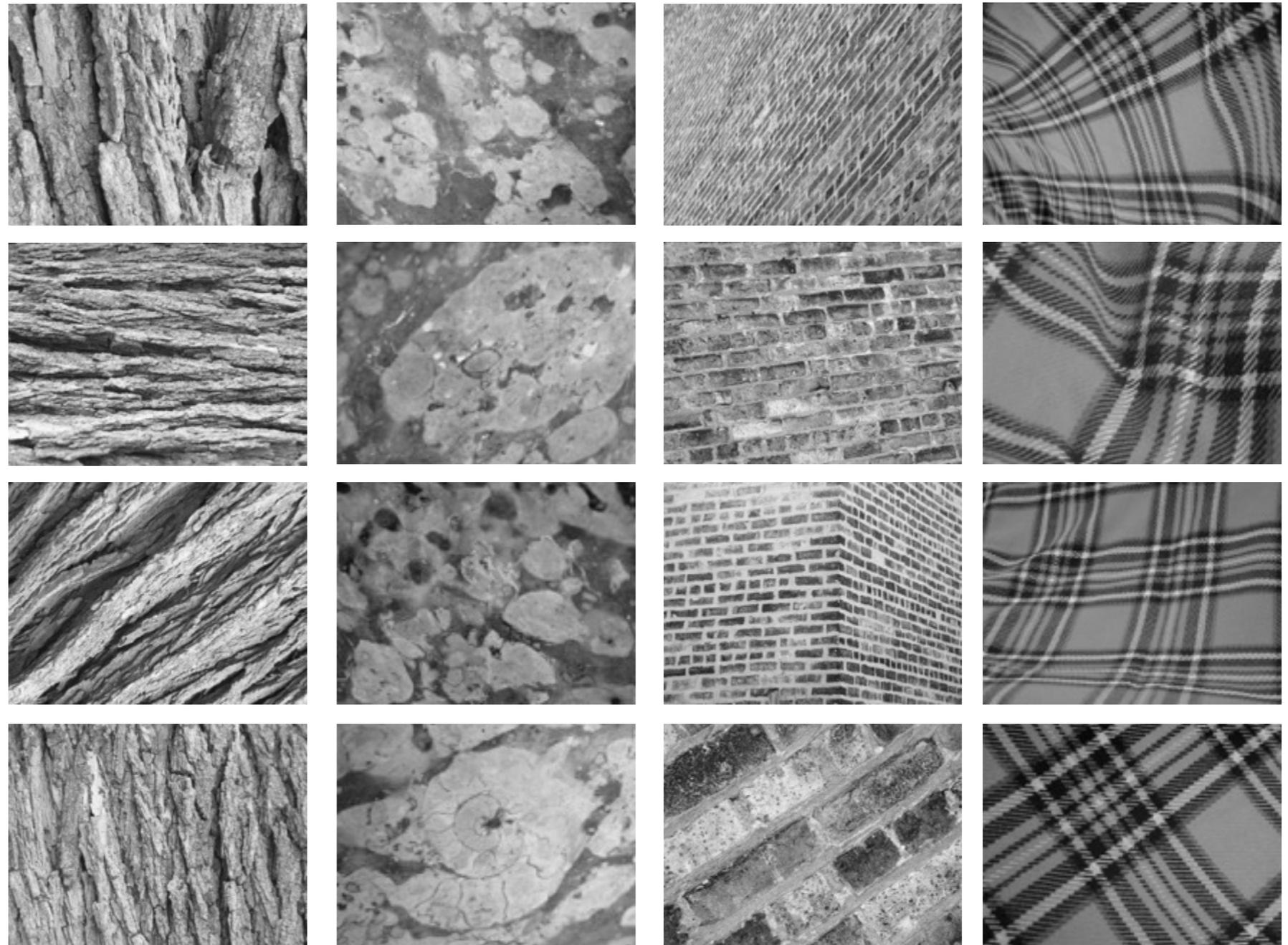
Wavelets on Translations W



Wavelets on Rigid Mvt. W



UIUC database:
25 classes



Scattering classification errors

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

*N. Poilvert
Matthew Hirn*

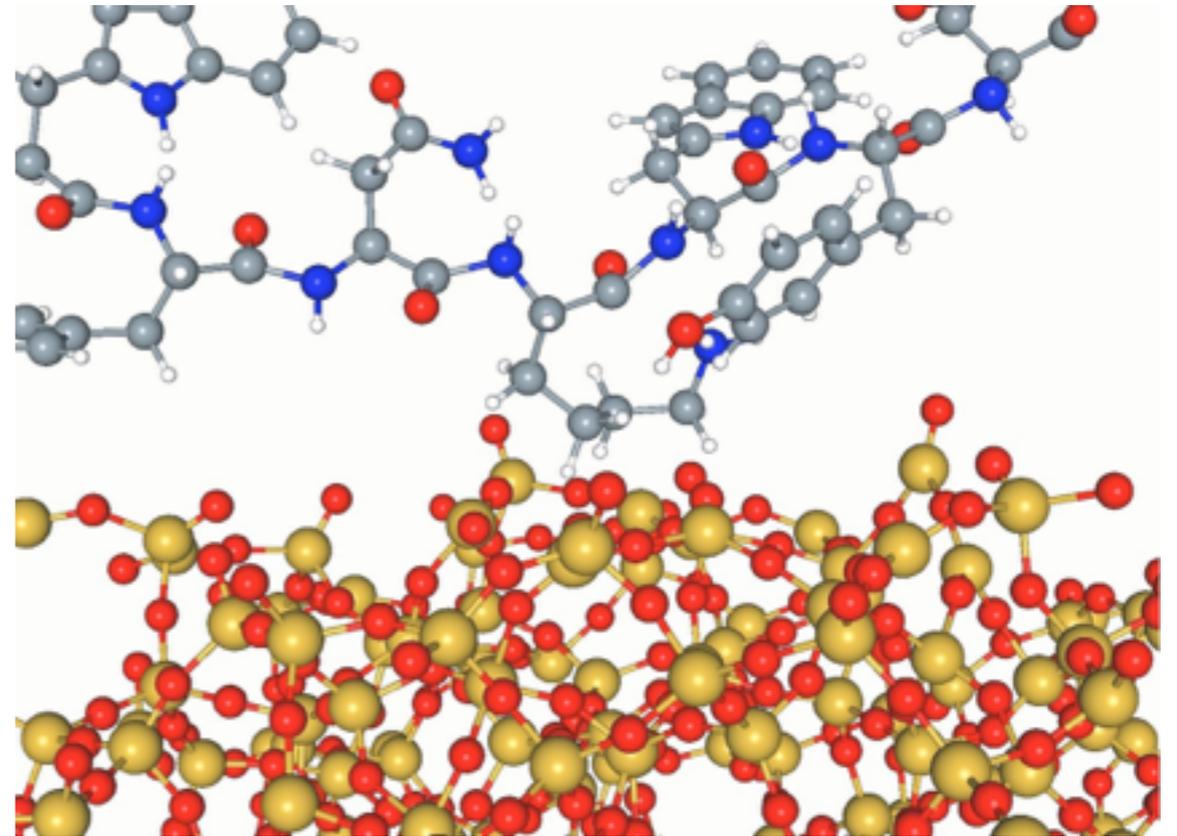
- Energy of d interacting bodies:

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

Astronomy



Quantum Chemistry



Kohn-Sham model:

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

↓

Molecular energy

↓

Kinetic energy

↓

electron-nuclei attraction

↓

electron-electron Coulomb repulsion

↓

Exchange correlat. energy

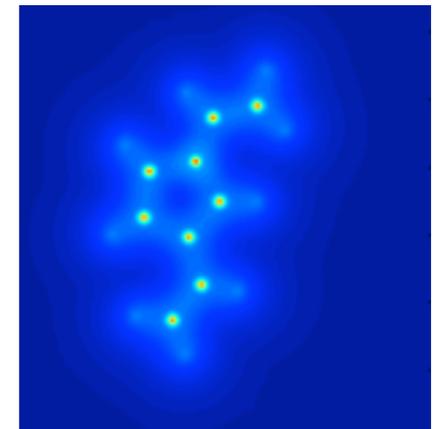
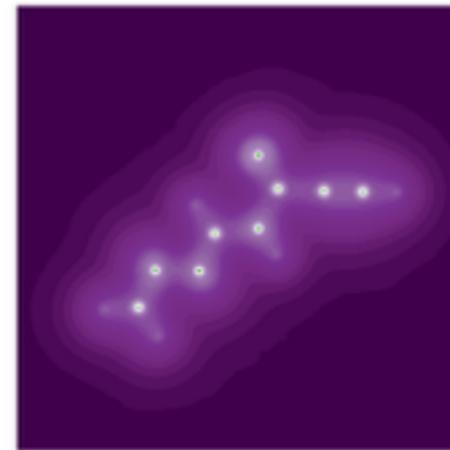
At equilibrium:

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

Quantum chemistry: $f(x)$ is invariant to rigid movements,
stable to deformations.

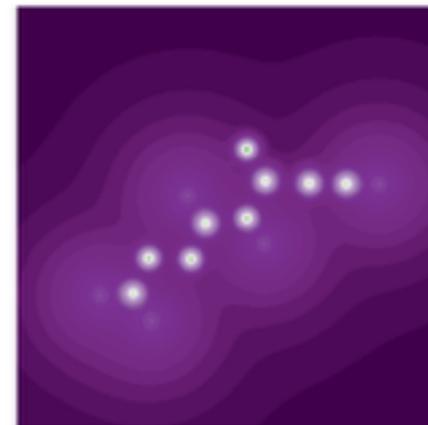
Depends on the true electronic density (Kohn-Sham)

Ground state
electronic density
computed with Schroedinger



- Can we estimate $f(x)$ from a naive electronic density ?

Density $\tilde{\rho}_x$ computed
as a sum of blobs



- Linear regressions computed with invariant change of variables:

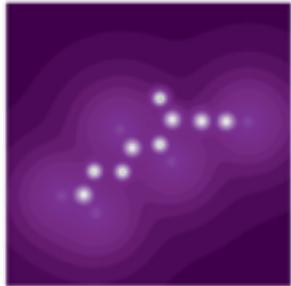
$$\Phi x = \{\phi_n(\tilde{\rho}_x)\}_n : \left| \begin{array}{l} \text{Fourier modulus coefficients and squared} \\ \text{or} \\ \text{scattering coefficients and squared} \end{array} \right.$$

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

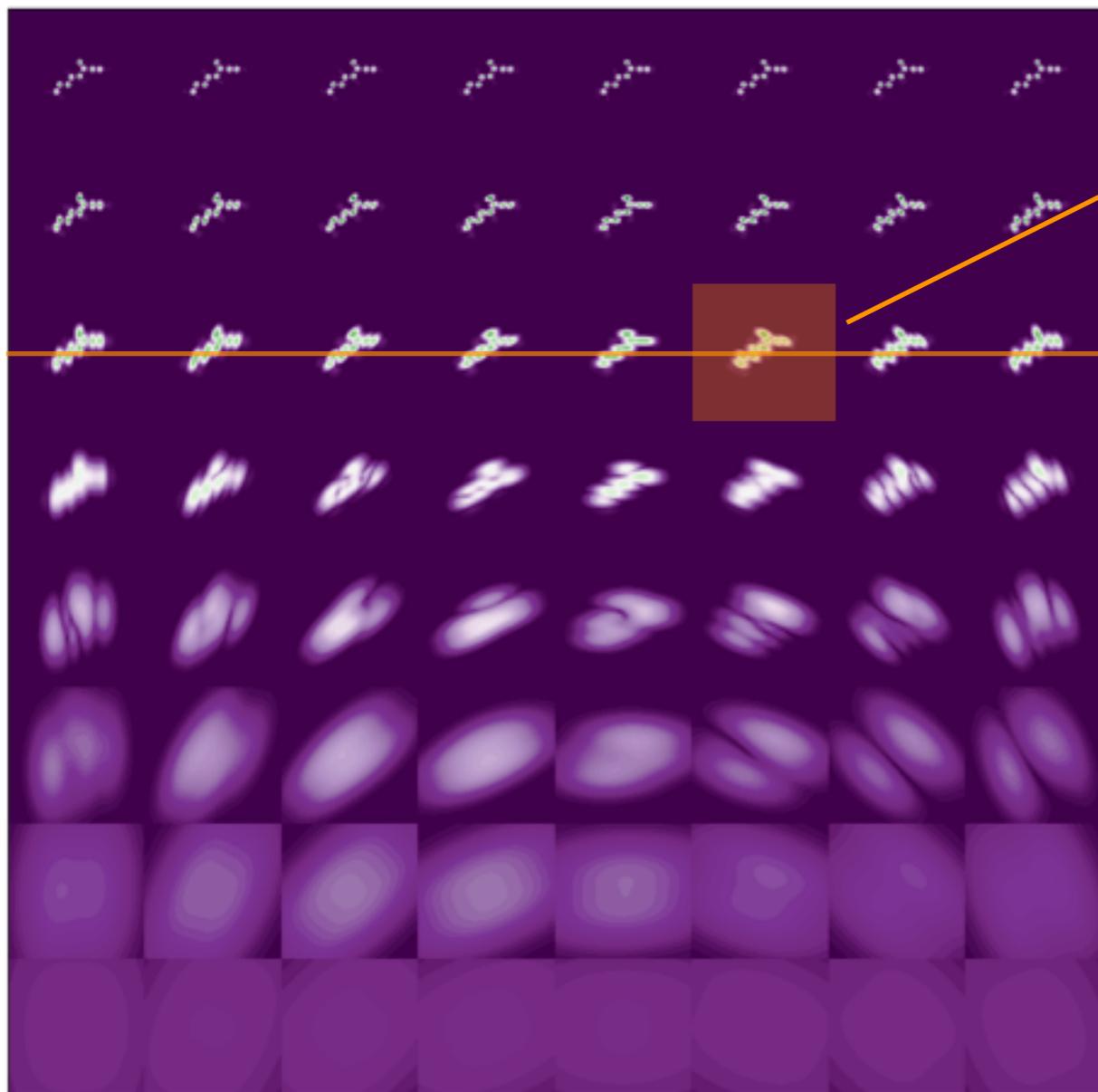
Regression coefficients w_k : equivalent potential.

Scattering Dictionary

$\rho(u)$



$$\downarrow |\rho * \psi_{j_1, \theta_1}(u)|$$



Rotations θ_1

2nd Order Interferences

Recover translation variability:

$$|\rho * \psi_{j_1, \theta_1}| * \psi_{j_2, \theta_2}(u)$$

Recover rotation variability:

$$|\rho * \psi_{j_1, \cdot}(u)| \otimes \bar{\psi}_{l_2}(\theta_1)$$

Combine to recover roto-translation variability:

$$||\rho * \psi_{j_1, \cdot}| * \psi_{j_2, \theta_2}(u) \otimes \bar{\psi}_{l_2}(\theta_1)|$$

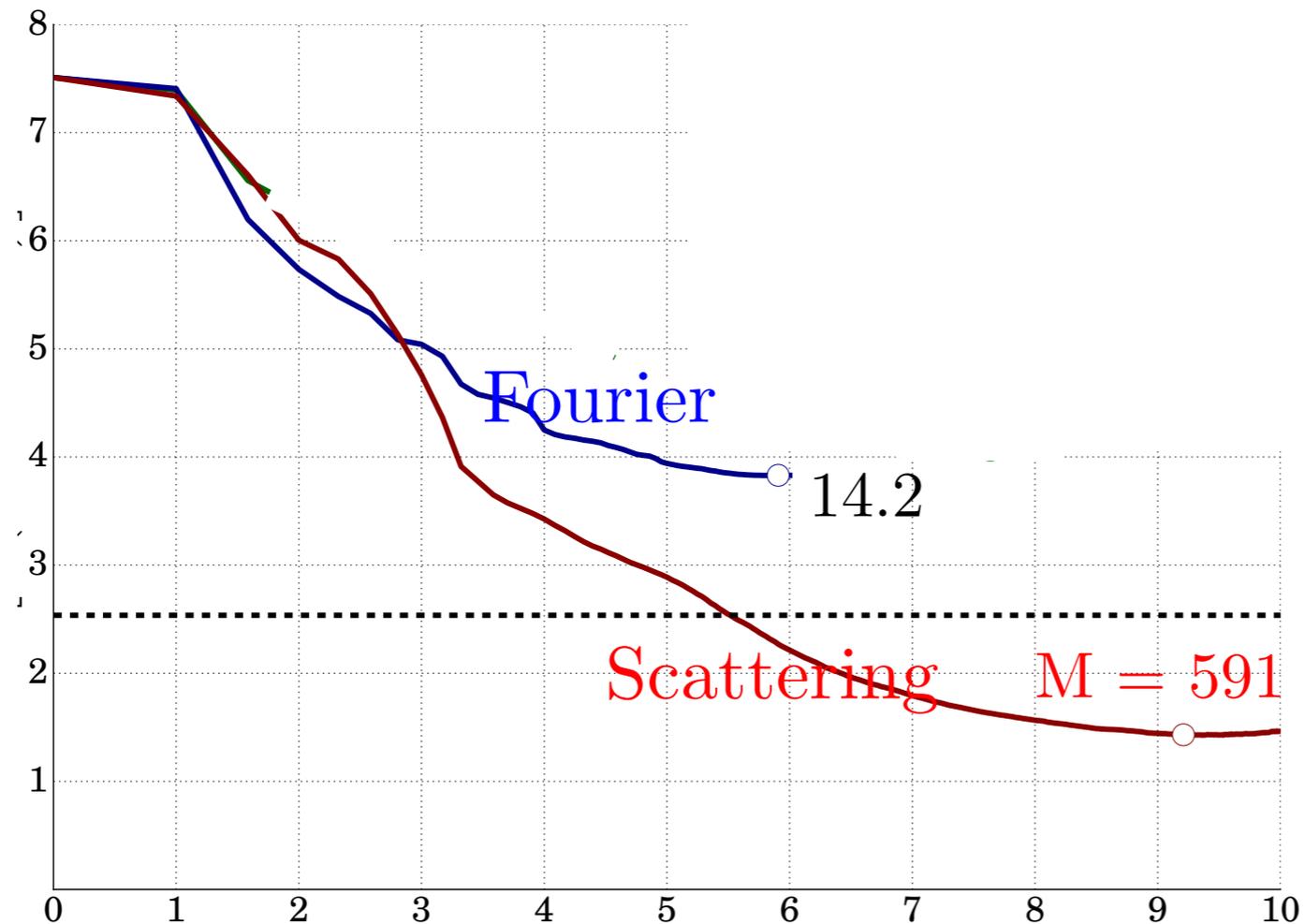
Scattering Regression

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

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

Interaction terms
across scales

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



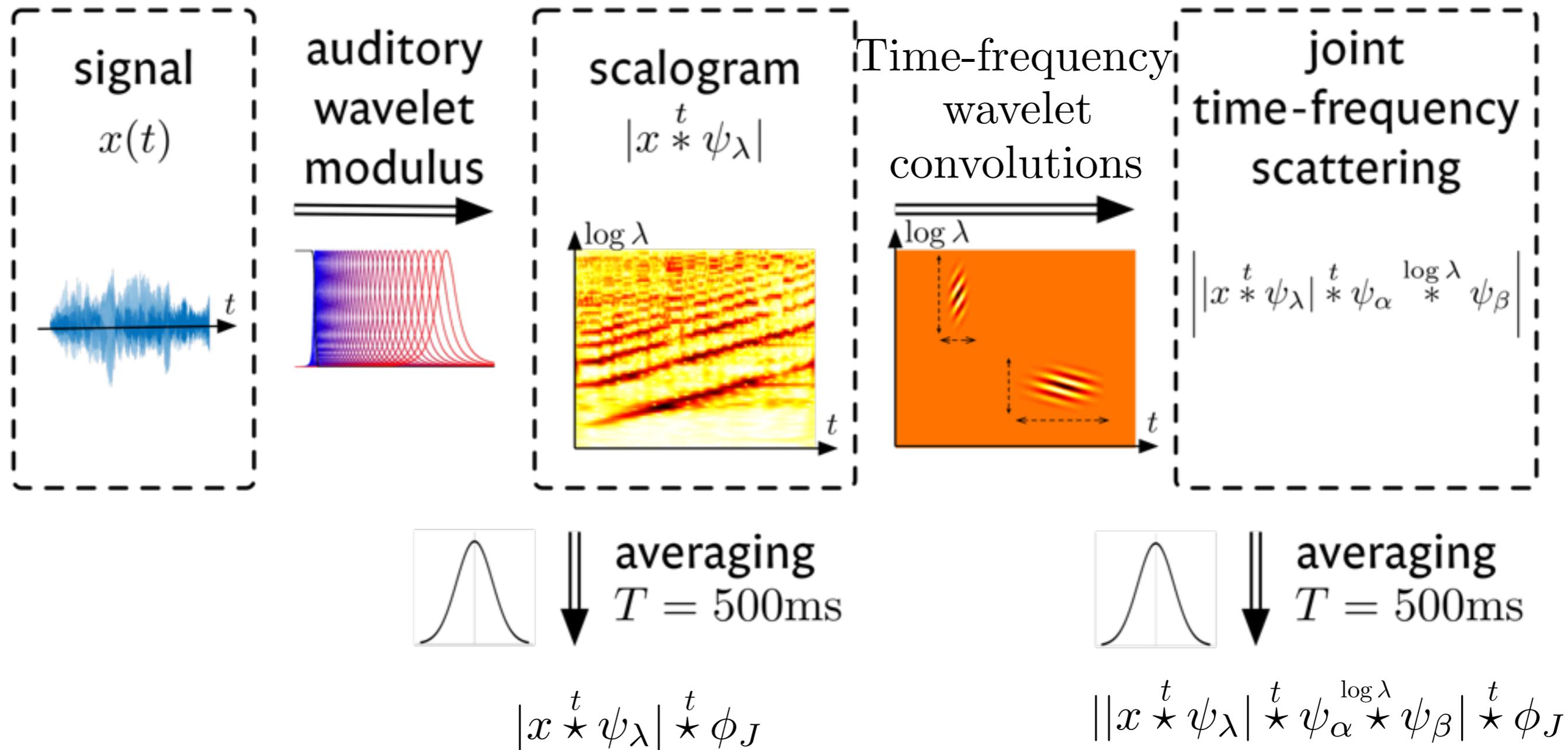
5.8 : State of the art

1.8 kcal/mol

$\log_2 M$

Time-Frequency Translation Group

J. Anden and V. Lostanlen



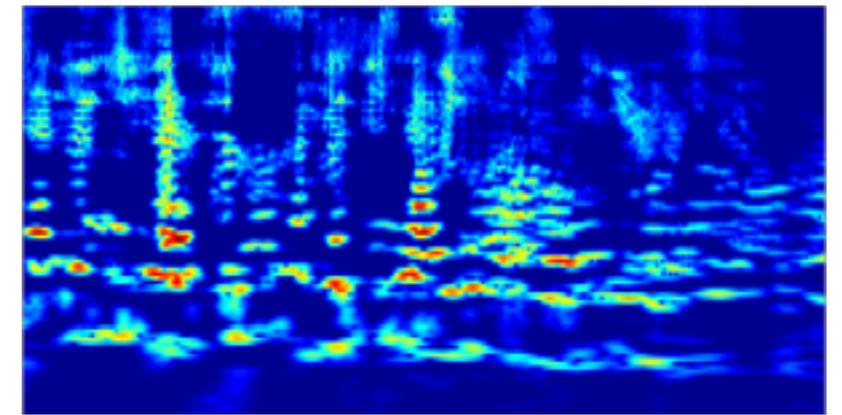
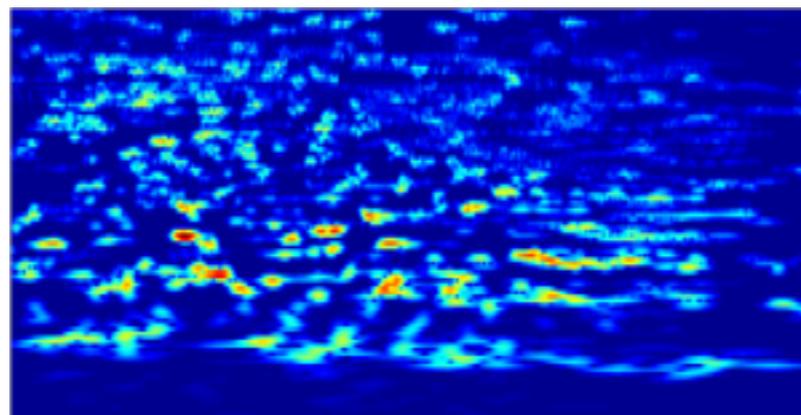
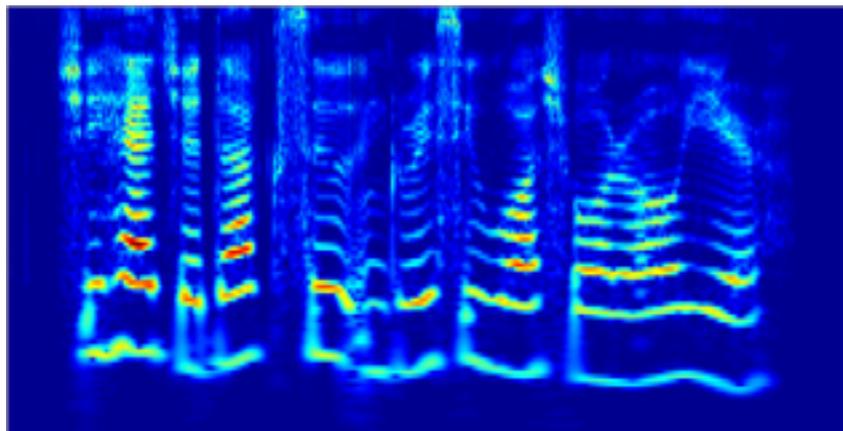
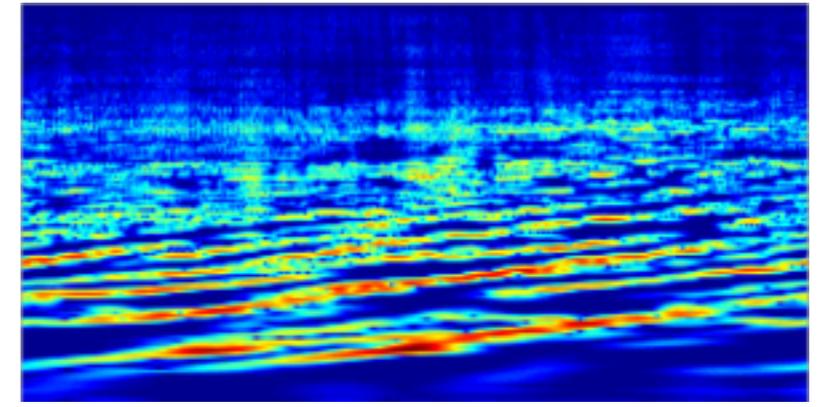
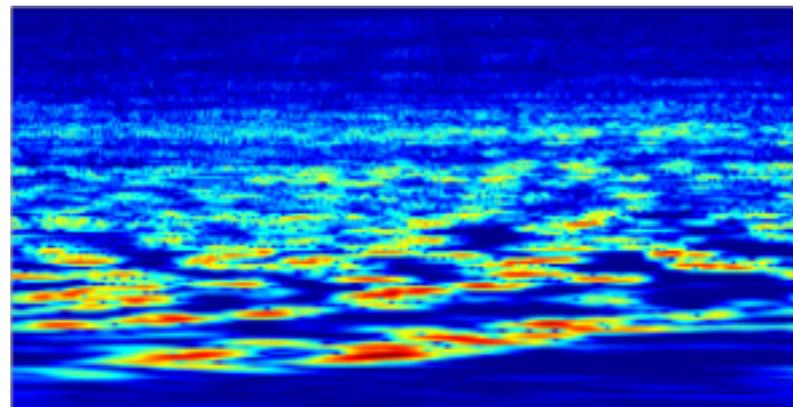
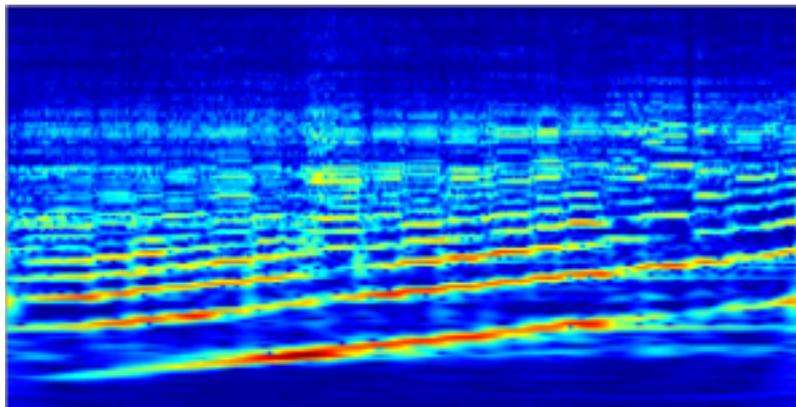
Joint Time-Frequency Scattering

J. Anden and V. Lostanlen

Original

Time Scattering

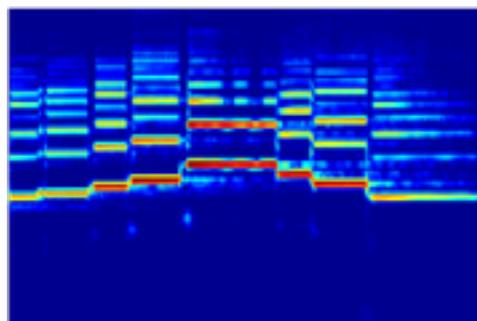
Time/Freq Scattering



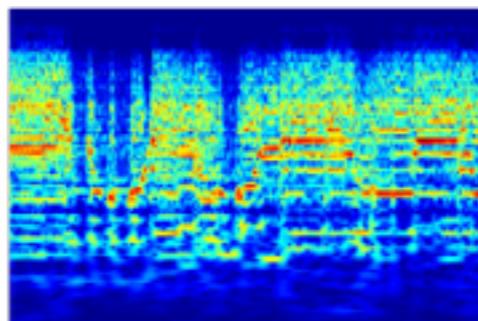
Musical Instrument Classification

J. Anden and V. Lostanlen

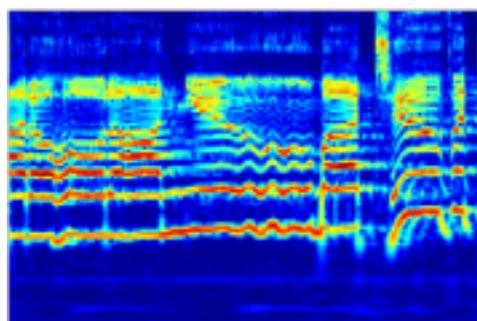
clarinet



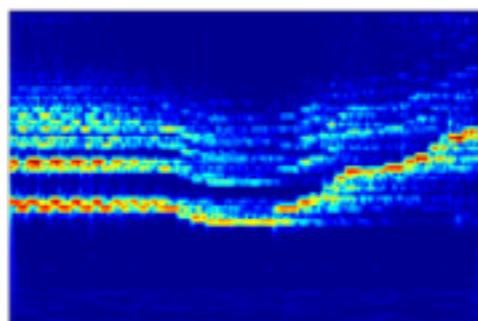
electric guitar



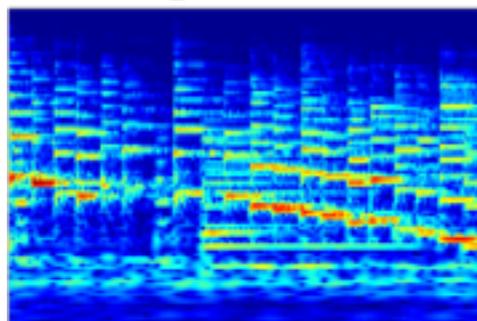
female singer



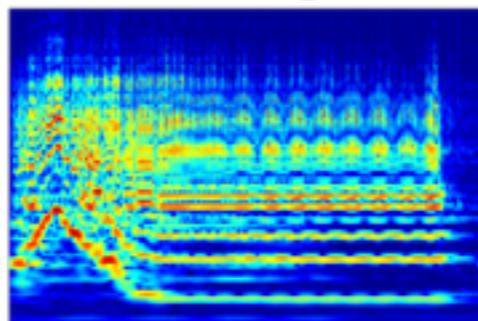
flute



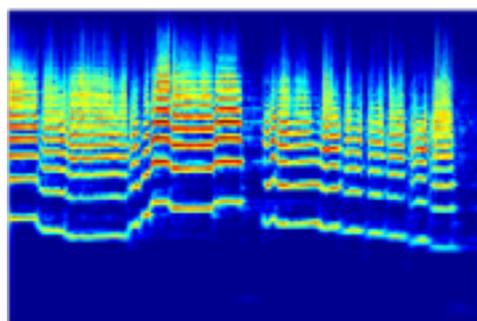
piano



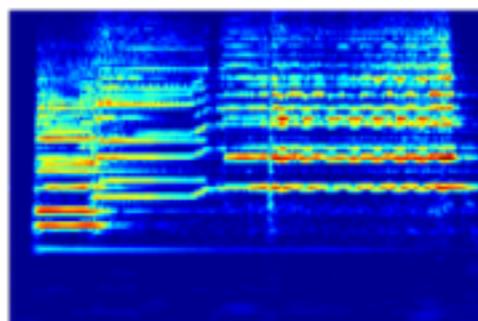
tenor saxophone



trumpet



violin



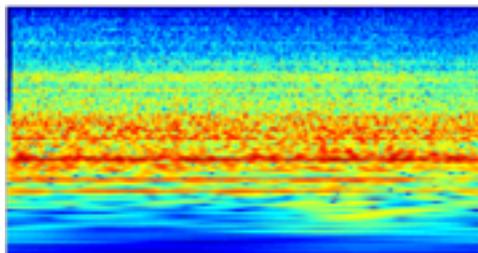
MedleyDB: 8 classes
10k training examples

class-wise average error

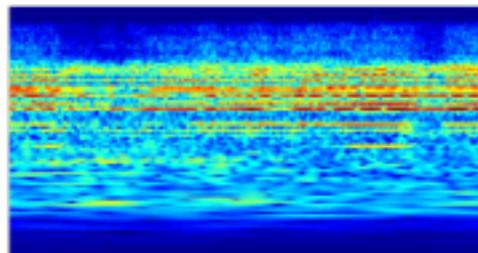
MFCC audio descriptors	0,39
time scattering	0,31
ConvNet	0,31
time-frequency scattering	0,18

Environmental Sound Classification

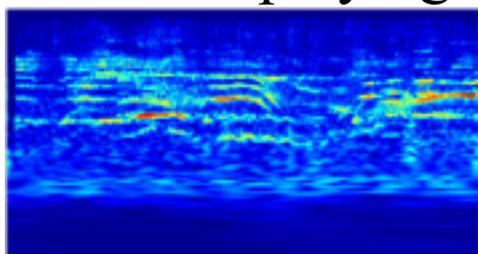
air conditioner



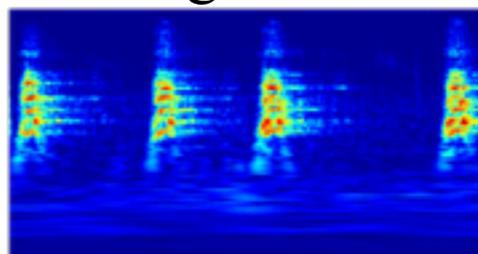
car horns



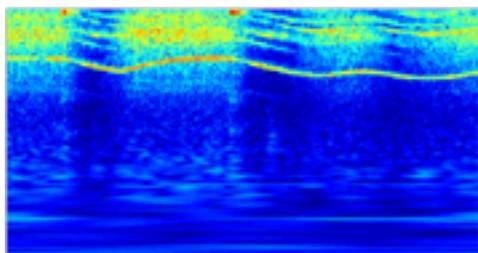
children playing



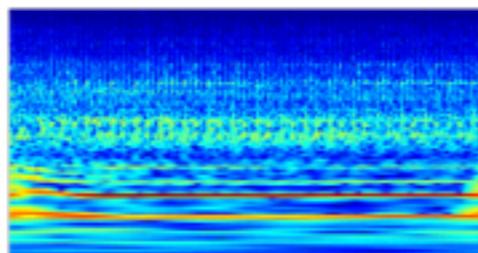
dog barks



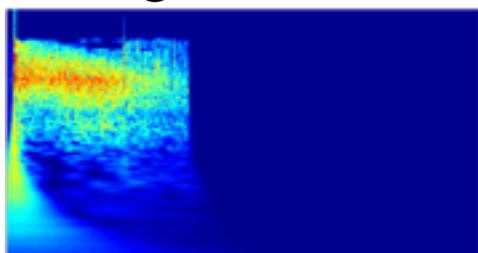
drilling



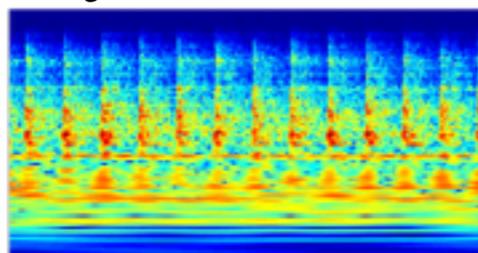
engine at idle



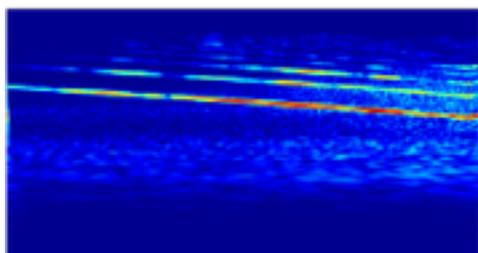
gunshot



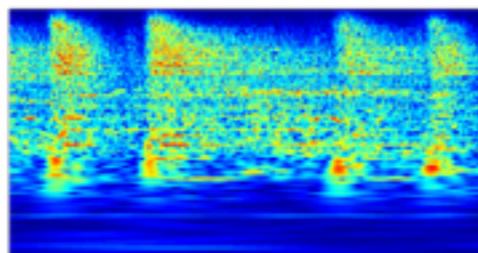
jackhammer



siren



street music



J. Anden and V. Lostanlen

UrbanSound8k: 10 classes
8k training examples

class-wise average error

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

Complex Image Classification

Edouard Oyallon

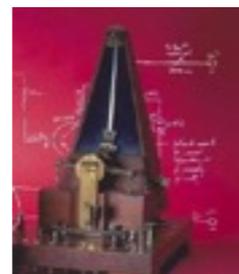
Arbre de Joshua



Ancre



Metronome



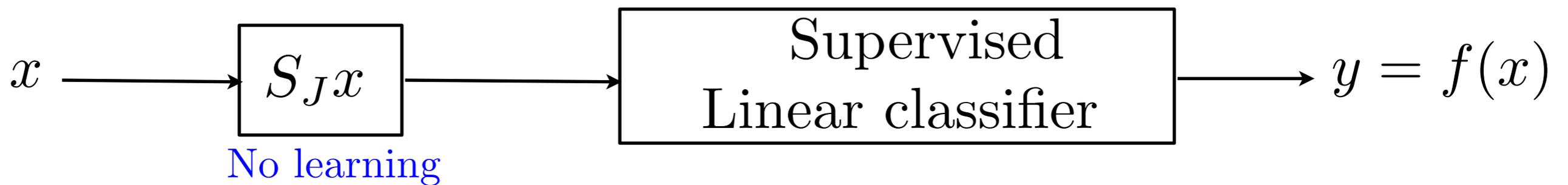
Castore



Nénuphare



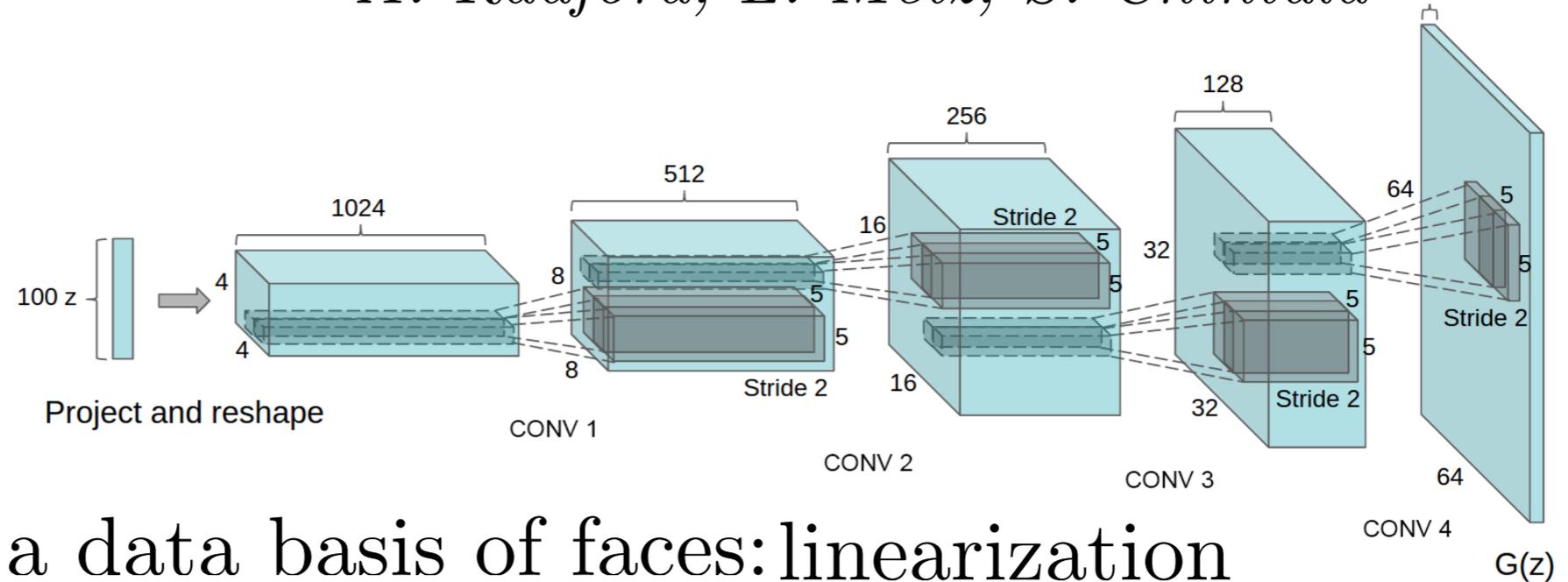
Bateau



Data Basis	Deep-Net	Scat/Unsupervised
CIFAR-10	7%	20%

Linearisation in Deep Networks

A. Radford, L. Metz, S. Chintala

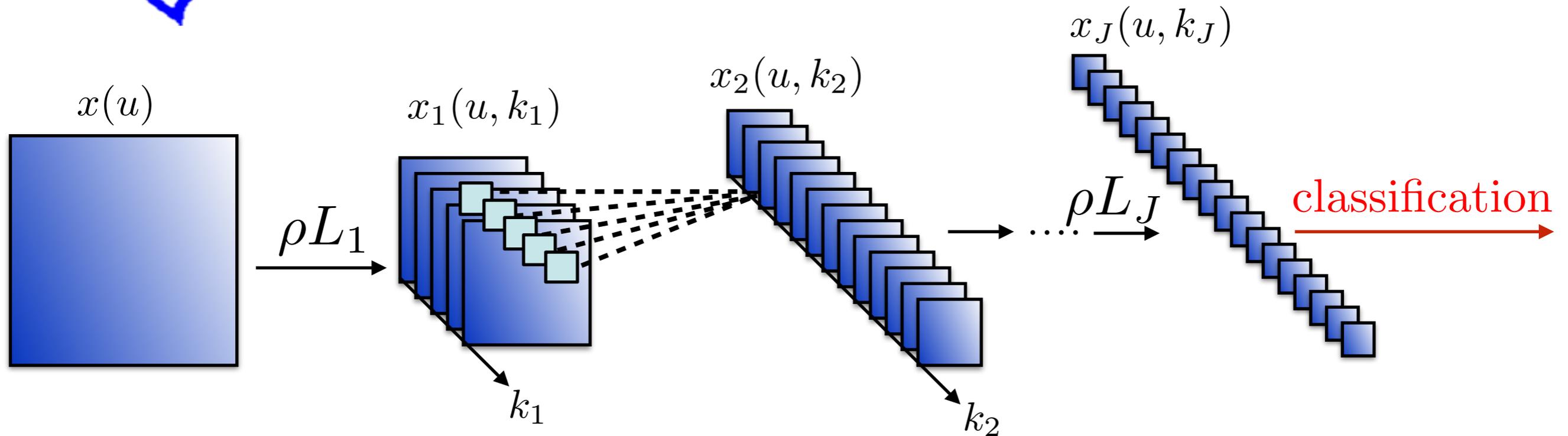


- Trained on a data basis of faces: linearization

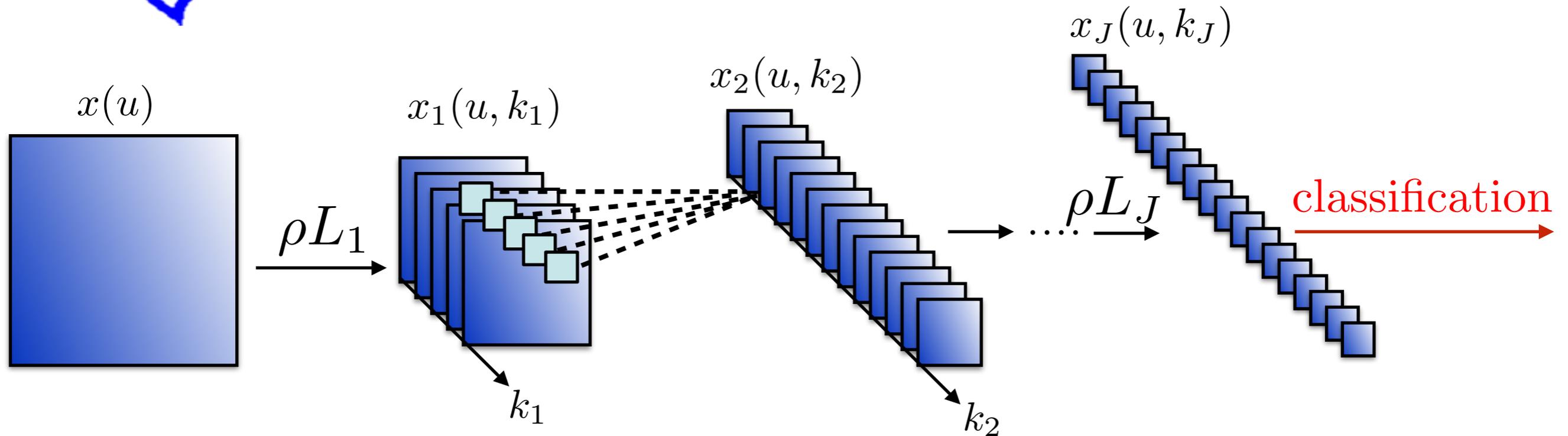


- On a data basis including bedrooms: interpolations





- The convolution network operators L_j have many roles:
 - Linearize non-linear transformations (symmetries)
 - Reduce dimension with projections
 - Memory storage of « characteristic » structures
- Difficult to separate these roles when analyzing learned networks



- Can we recover symmetry groups from the matrices L_j ?
- What kind of groups ?
- Can we characterise the regularity of $f(x)$ from these groups ?
- Can we define classes of high-dimensional « regular » functions that are well approximated by deep neural networks ?
- Can we get approximation theorems giving errors depending on number of training examples, with a fast decay ?

Conclusions

- Deep convolutional networks have spectacular high-dimensional approximation capabilities.
- Seem to compute hierarchical invariants of complex symmetries
- Used as models in physiological vision and audition
- Close link with particle and statistical physics
- Outstanding mathematical problem to understand them:
notions of complexity, regularity, approximation theorems...

Understanding Deep Convolutional Networks, arXiv 2016.