

1

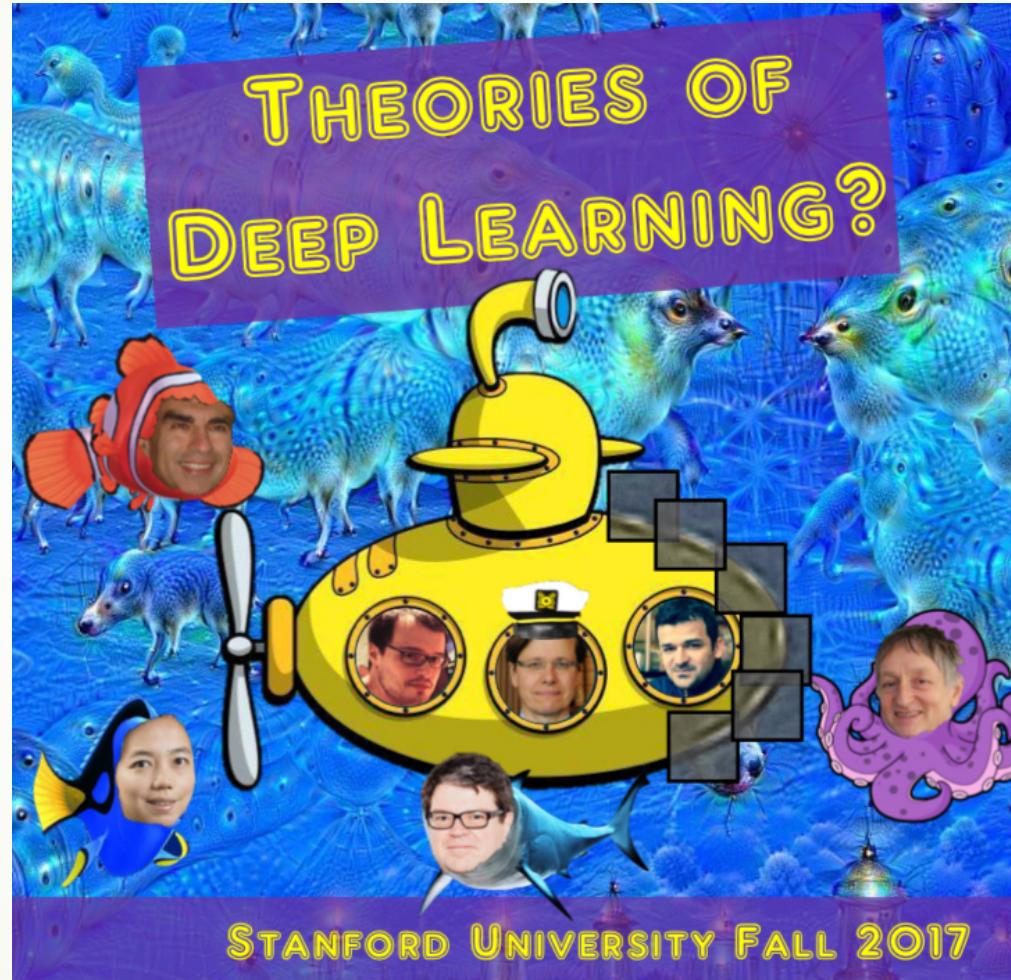
# Harmonic Analysis of Deep Convolutional Networks

Yuan YAO  
HKUST

Based on Mallat and Bolcskei talks etc.



# Acknowledgement



A following-up course at HKUST: <https://deeplearning-math.github.io/>

# High Dimensional Natural Image Classification

- 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



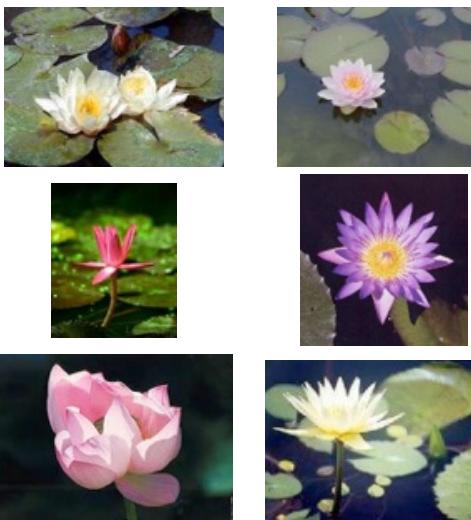
Joshua Tree



Beaver



Lotus



Water Lily

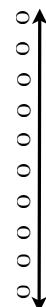
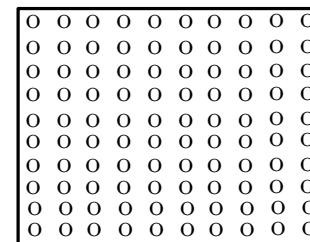


Huge variability  
inside classes

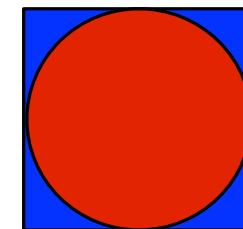
Find invariants

# Curse of Dimensionality

- Analysis in high dimension:  $x \in \mathbb{R}^d$  with  $d \geq 10^6$ .
- Points are far away in high dimensions  $d$ :
  - 10 points cover  $[0, 1]$  at a distance  $10^{-1}$
  - 100 points for  $[0, 1]^2$
  - need  $10^d$  points over  $[0, 1]^d$   
impossible if  $d \geq 20$



$$\lim_{d \rightarrow \infty} \frac{\text{volume sphere of radius } r}{\text{volume } [0, r]^d} = 0$$

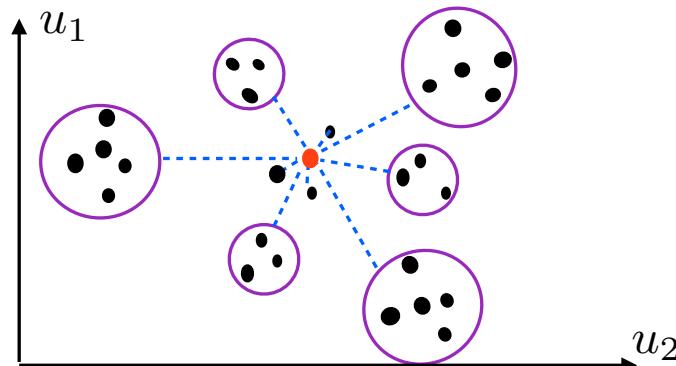


points are  
concentrated  
in  $2^d$  corners!

⇒ Euclidean metrics are not appropriate on **raw data**.

# A Blessing from Physical world? Multiscale “compositional” sparsity

- 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.



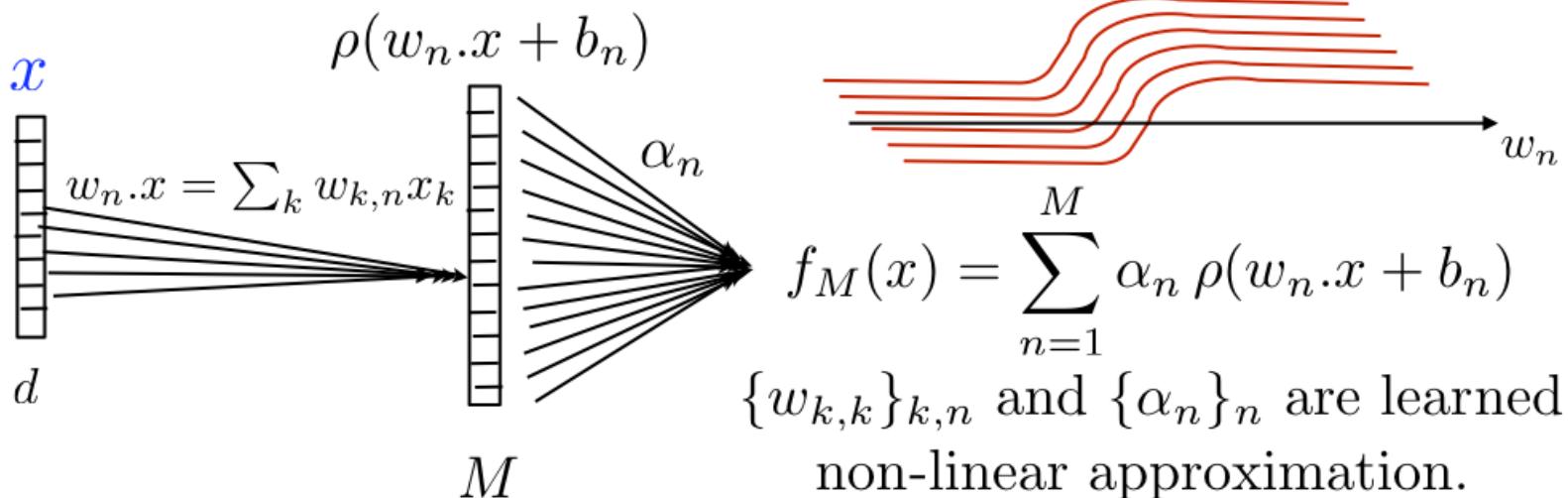
# Learning as an Approximation



- 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  $\rho(x.w_n + b_n)$



Cybenko, Hornik, Stinchcombe, White

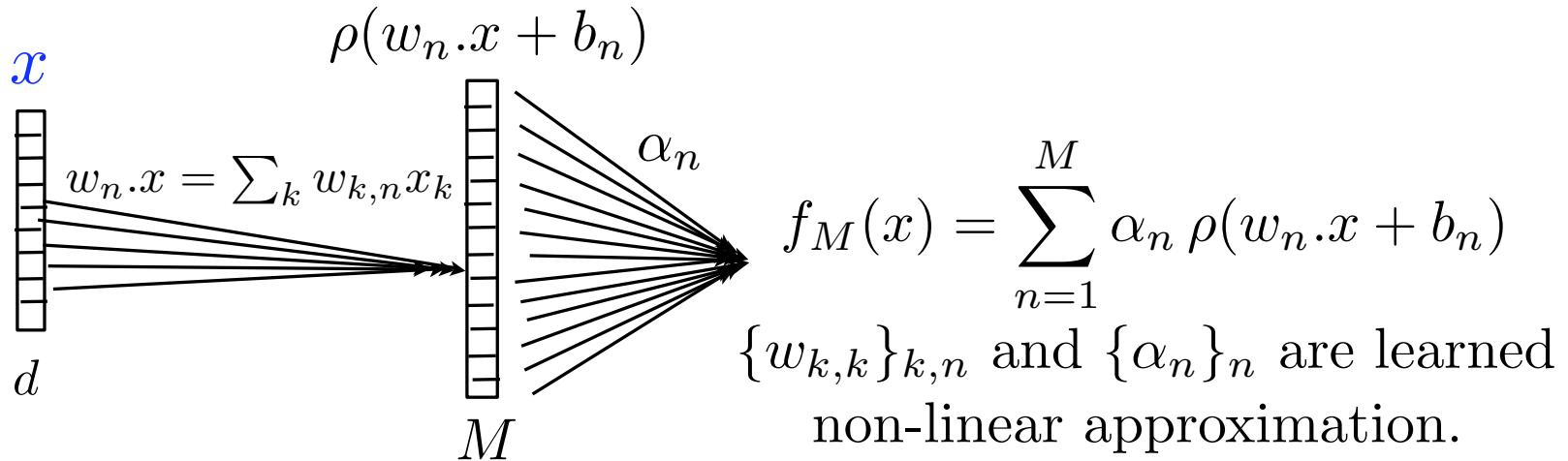
**Theorem:** For "reasonable" bounded  $\rho(u)$   
and appropriate choices of  $w_{n,k}$  and  $\alpha_n$ :

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

No big deal: curse of dimensionality still there.

# 1 Hidden Layer Neural Networks

One-hidden layer neural network:



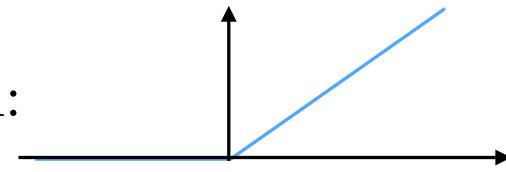
Fourier series:  $\rho(u) = e^{iu}$

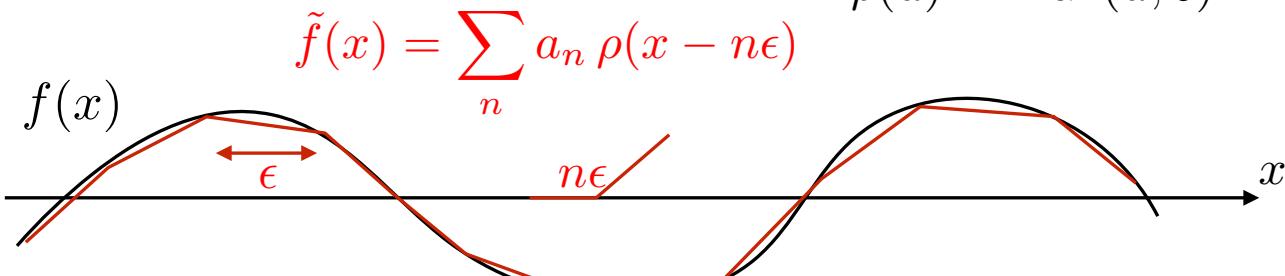
$$f_M(x) = \sum_{n=1}^M \alpha_n e^{iw_n \cdot x}$$

For nearly all  $\rho$ : essentially same approximation results.

# Piecewise Linear Approximation

- Piecewise linear approximation:


$$\rho(u) = \max(u, 0)$$



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

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

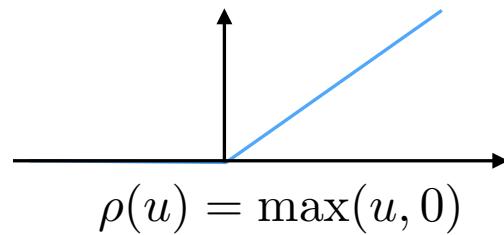
Need  $M = \epsilon^{-1}$  points to cover  $[0, 1]$  at a distance  $\epsilon$

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

# 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  $\epsilon$ :

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

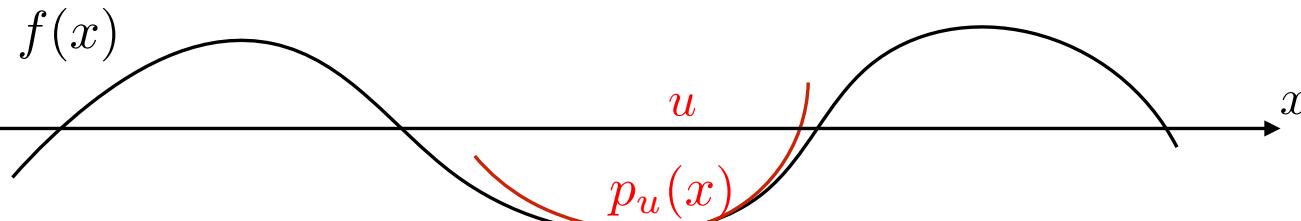
need  $M = \epsilon^{-d}$  points to cover  $[0, 1]^d$  at a distance  $\epsilon$

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

Curse of dimensionality!

# Approximation with Regularity

- What prior condition makes learning possible ?
- Approximation of regular functions in  $\mathbf{C}^s[0, 1]^d$ :  
$$\forall x, u \quad |f(x) - p_u(x)| \leq C |x - u|^s \text{ with } p_u(x) \text{ polynomial}$$



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

Need  $M^{-d/s}$  points 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.**



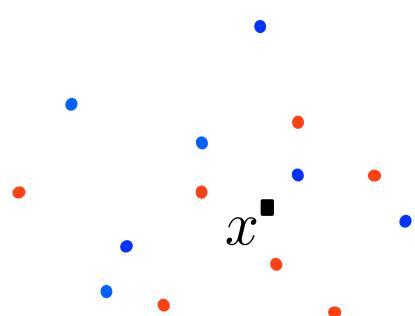
# Kernel Learning

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

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

$$\tilde{f}(x) = \langle \Phi(x), w \rangle = \sum_{\text{1D projection}} w_k \phi_k(x) .$$

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

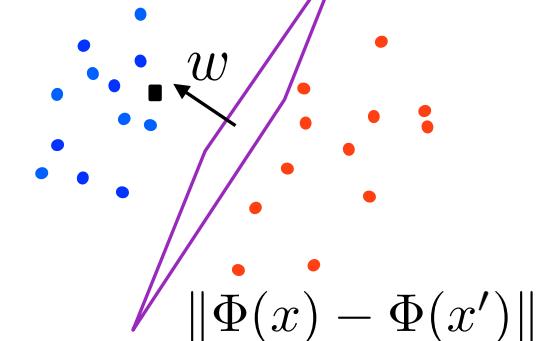


Metric:  $\|x - x'\|$

$$\xrightarrow{\Phi}$$

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

Linear Classifier



- How and when is possible to find such a  $\Phi$  ?
- 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  $< N$ .

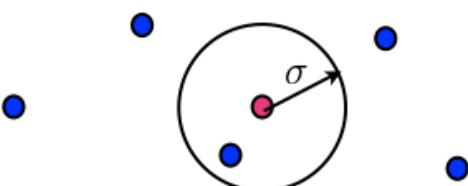
⇒ Choose  $\Phi$  increasing dimensionality !

**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  $\sigma$  is small, nearest neighbor classifier type:



# Spirit in Fisher's Linear Discriminant Analysis

## Reduction of Dimensionality

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

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

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

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

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

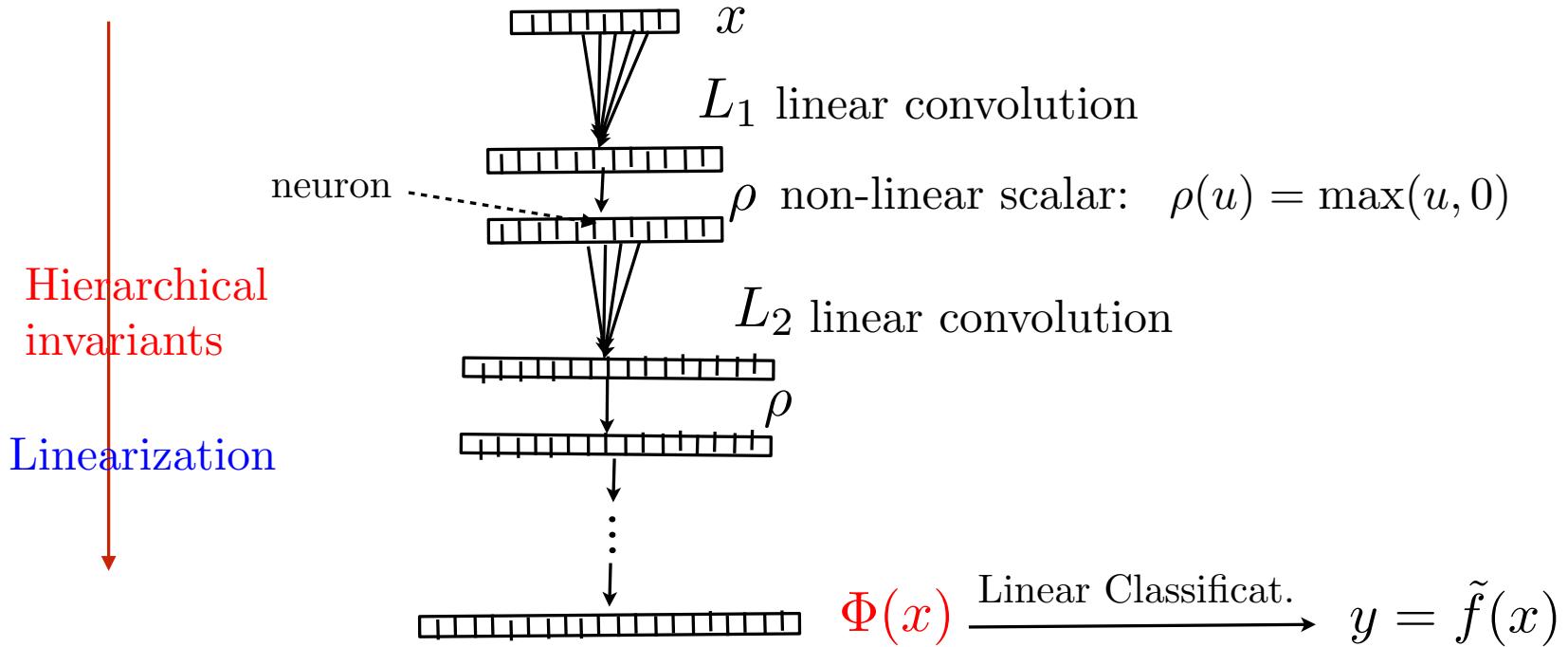
Discriminative:  $\|\Phi(x) - \Phi(x')\| \geq C^{-1} |f(x) - f(x')|$

- 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*



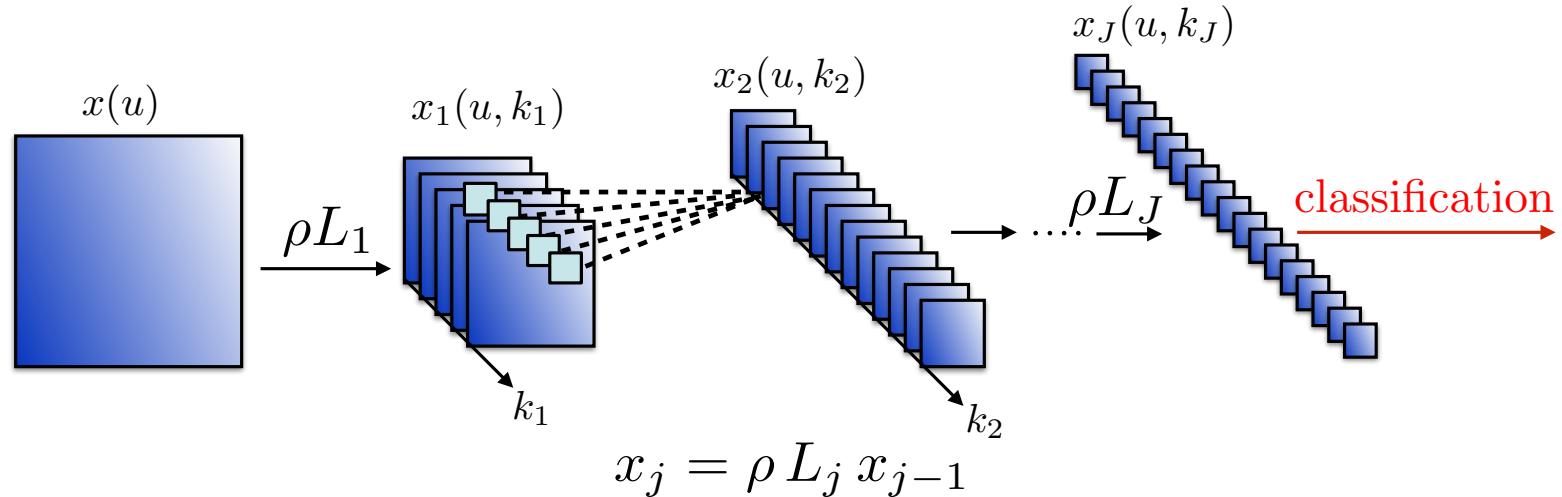
Optimize  $L_j$  with **architecture constraints**: over  $10^9$  parameters

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

Why does it work so well ? **A difficult problem**



# Deep Convolutional Networks



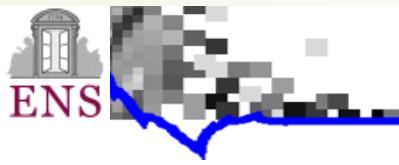
- $L_j$  is a linear combination of convolutions and subsampling:

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

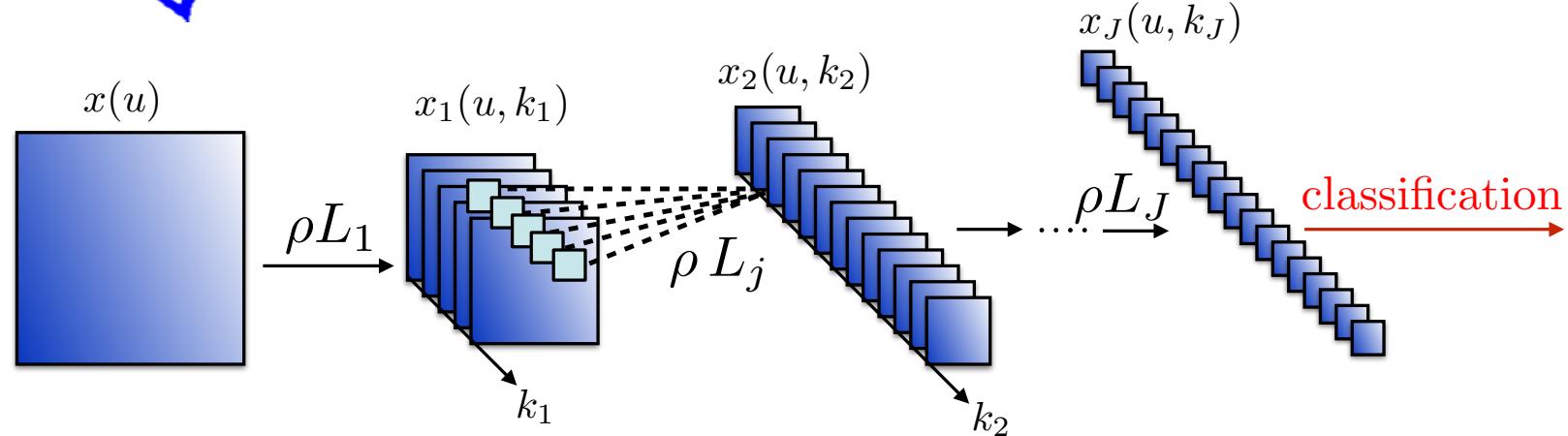
sum across channels

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

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



# Many Questions

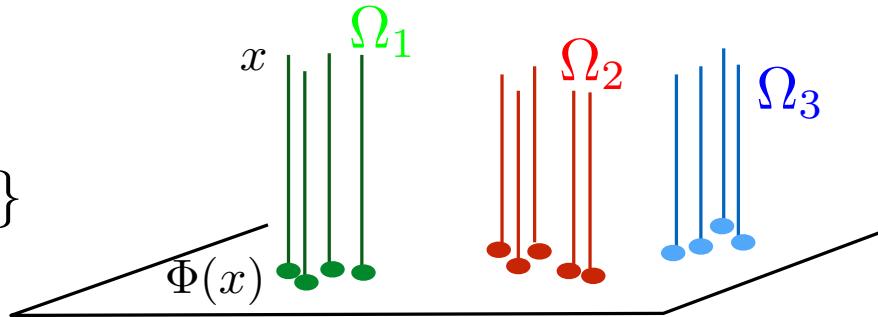


- Why convolutions ? Translation covariance.
- Why no overfitting ? Contractions, dimension reduction
- Why hierarchical cascade ?
- Why introducing non-linearities ?
- How and what to linearise ?
- What are the roles of the multiple channels in each layer ?



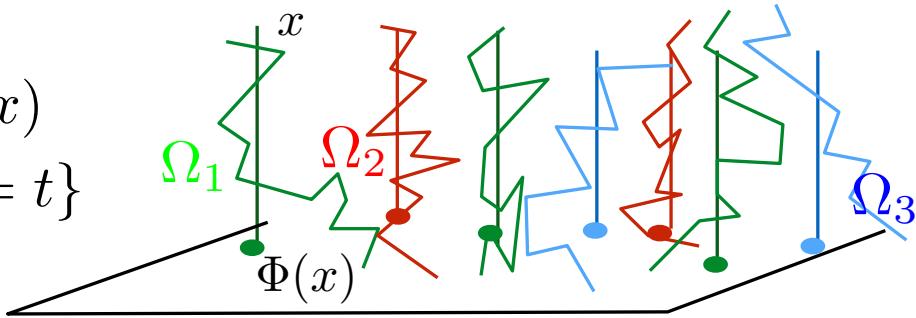
# Linear Dimension Reduction

*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.*

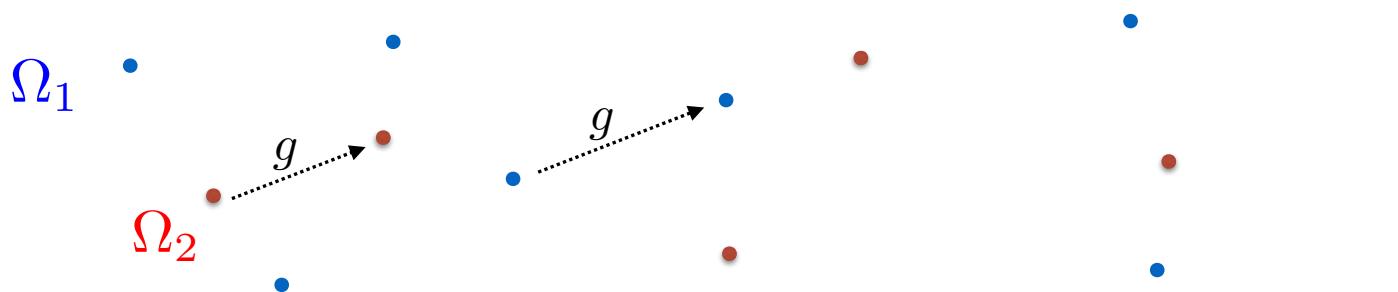
*Classes*  
*Level sets of  $f(x)$*   
 $\Omega_t = \{x : f(x) = t\}$



- If level sets  $\Omega_t$  are not parallel to a linear space
  - Linearise them with a change of variable  $\Phi(x)$
  - Then reduce dimension with linear projections
- Difficult because  $\Omega_t$  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:

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

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:

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

- Lie group: infinitely small generators (Lie Algebra)

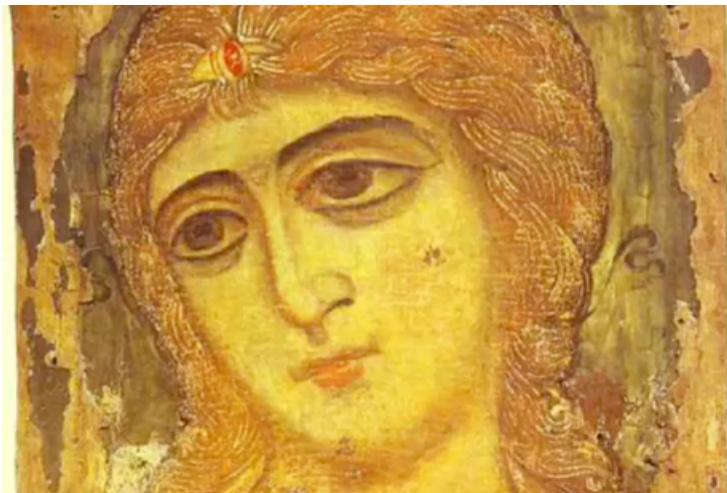
# Translation and Deformations

- Digit classification:

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



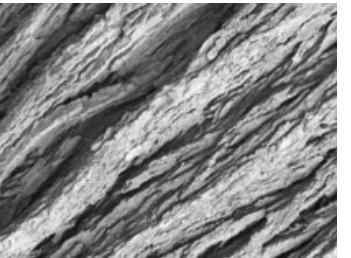
- Globally invariant to the translation group: small
- Locally invariant to small diffeomorphisms: huge group



*Video of Philipp Scott Johnson*

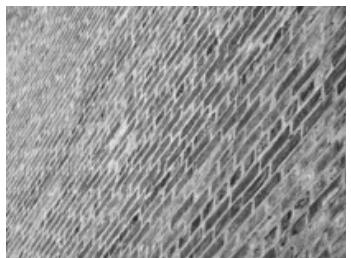
# Rotation and Scaling Variability

- Rotation and **deformations**



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

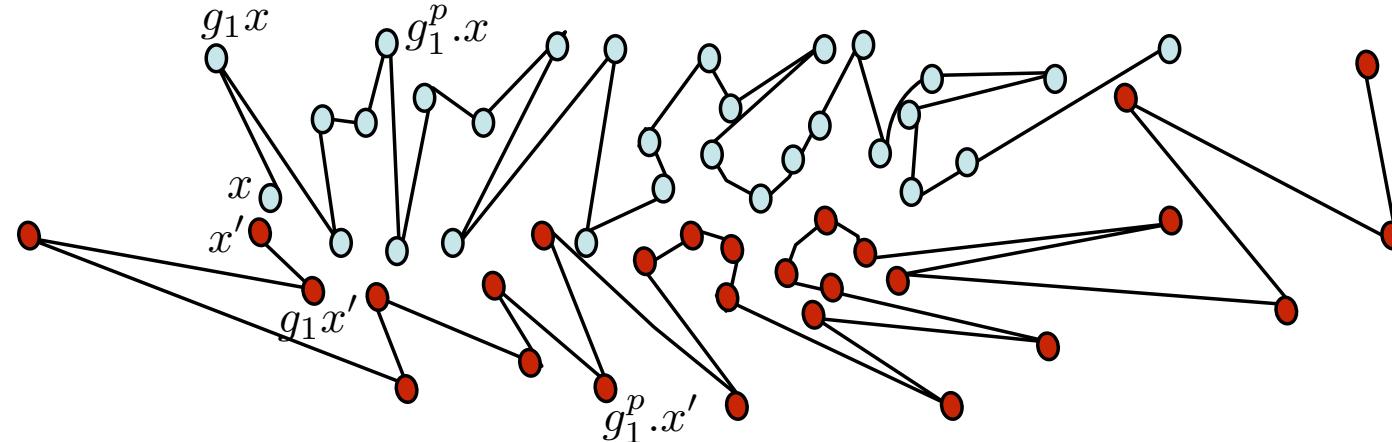
- Scaling and **deformations**



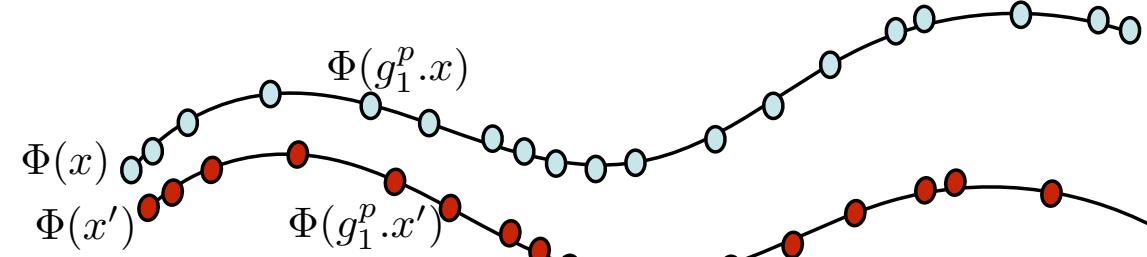
Group:  $\mathbb{R} \times \text{Diff}(\mathbb{R})$

# Linearize Symmetries

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



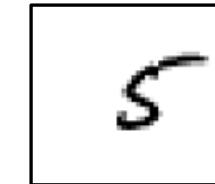
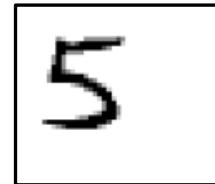
- Linearise symmetries with a change of variable  $\Phi(x)$



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

# Translation and Deformations

- Digit classification:

 $x(u)$  $x'(u)$ 

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

Linearize small  
diffeomorphisms:  
 $\Rightarrow$  Lipschitz regular



*Video of Philipp Scott Johnson*



# Translations and Deformations

- Invariance to translations:

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

- 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 Deformation Instability

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

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

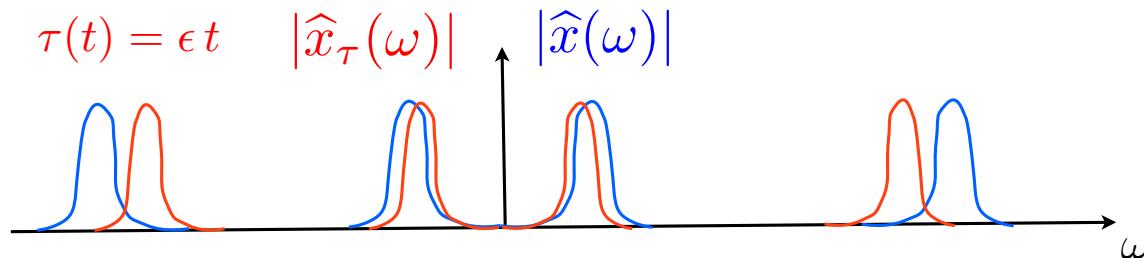
The modulus is invariant to translations:

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

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

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

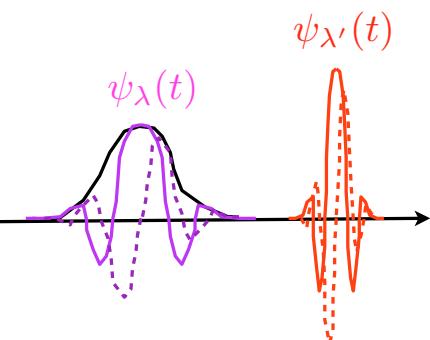
$$\tau(t) = \epsilon t \quad |\hat{x}_\tau(\omega)|$$



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

# Wavelet Transform

- Complex wavelet:  $\psi(t) = \psi^a(t) + i \psi^b(t)$
- Dilated:  $\psi_\lambda(t) = 2^{-j} \psi(2^{-j}t)$  with  $\lambda = 2^{-j}$ .



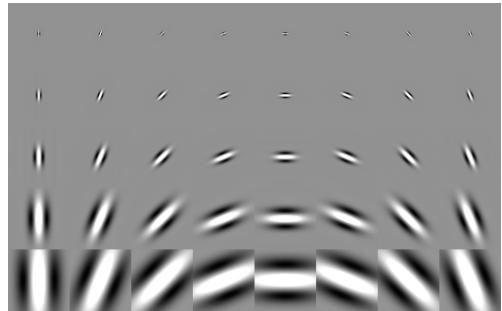
- Wavelet transform:  $x \star \psi_\lambda(t) = \int x(u) \psi_\lambda(t-u) du$   
$$Wx = \begin{pmatrix} x \star \phi(t) \\ x \star \psi_\lambda(t) \end{pmatrix}_{t,\lambda}$$

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

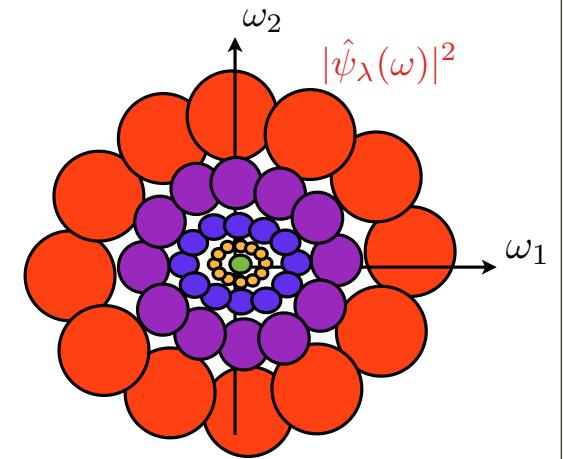
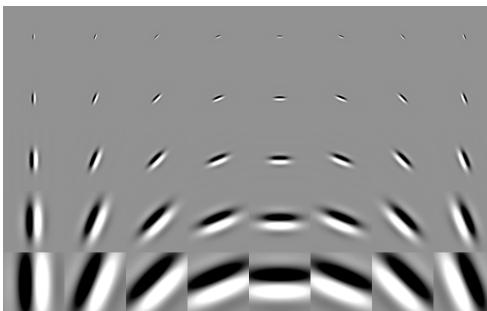
## Image Wavelet Transform

- Complex wavelet:  $\psi(t) = \psi^a(t) + i \psi^b(t)$  ,  $t = (t_1, t_2)$   
rotated and dilated:  $\psi_\lambda(t) = 2^{-j} \psi(2^{-j} rt)$  with  $\lambda = (2^j, r)$

real parts



imaginary parts



- Wavelet transform:  $Wx = \begin{pmatrix} x \star \phi(t) \\ x \star \psi_\lambda(t) \end{pmatrix}_{t,\lambda}$

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



# 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.

# Why Wavelets?

- ▶ Wavelets (complex band limited) are uniformly **stable to deformations**

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

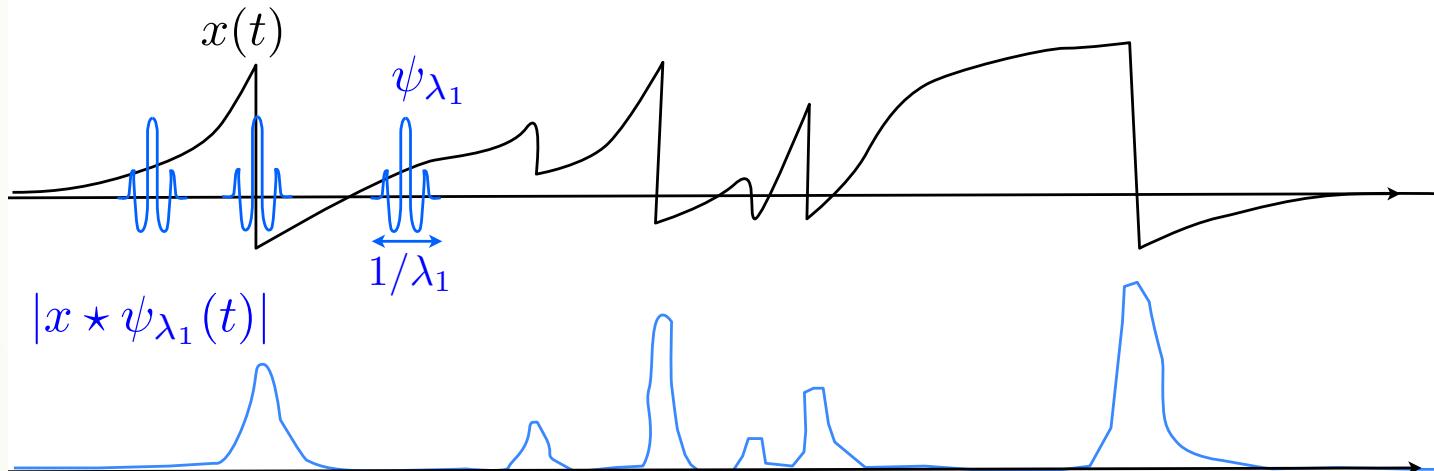
$$\|\psi_\lambda - \psi_{\lambda,\tau}\| \leq C \sup_t |\nabla \tau(t)| .$$

- ▶ Wavelets are **sparse** representations of functions
- ▶ Wavelets separate **multiscale** information
- ▶ Wavelets can be locally **translation invariant**

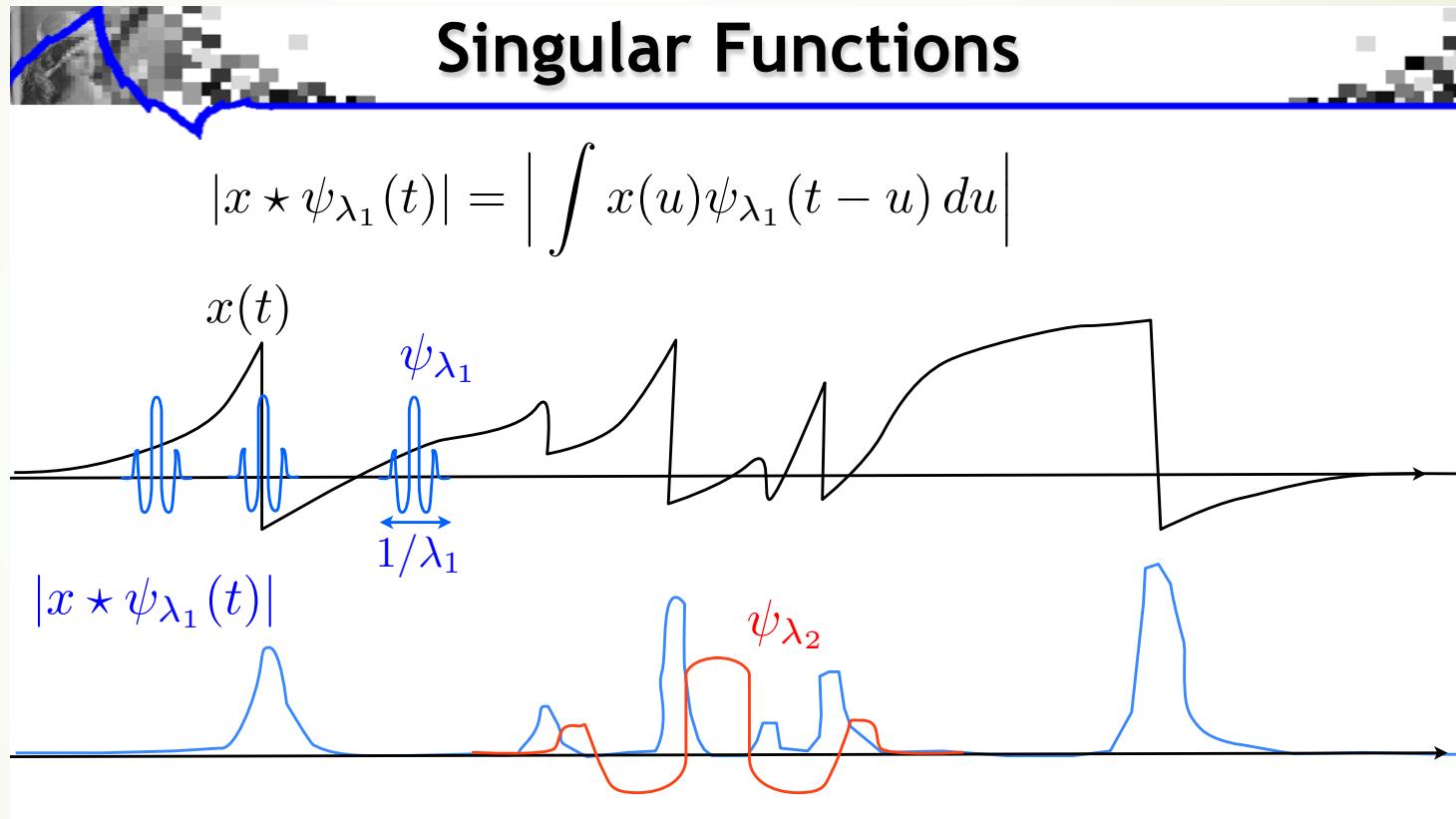
# Sparsity of Wavelet Transforms

## Singular Functions

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



# Singularity is preserved in multiscale transform

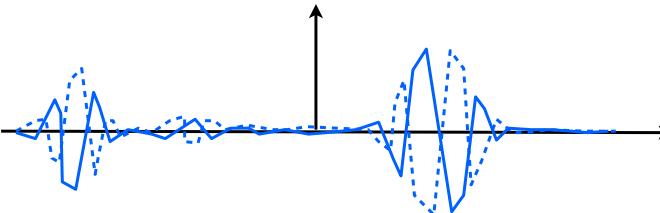


Second wavelet transform modulus

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

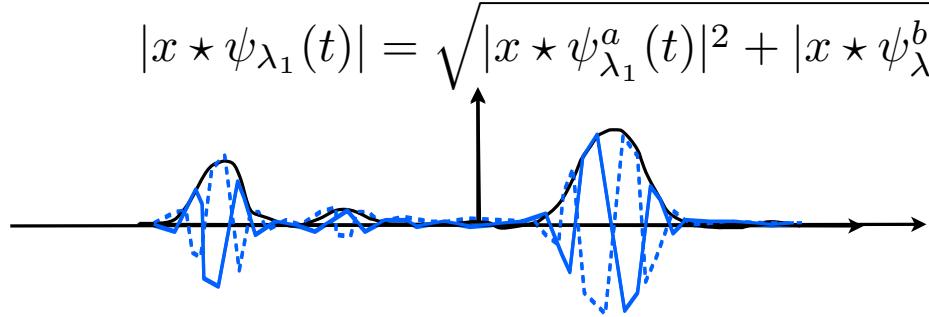
## Wavelet Translation Invariance

$$x \star \psi_{\lambda_1}(t) = x \star \psi_{\lambda_1}^a(t) + i x \star \psi_{\lambda_1}^b(t)$$



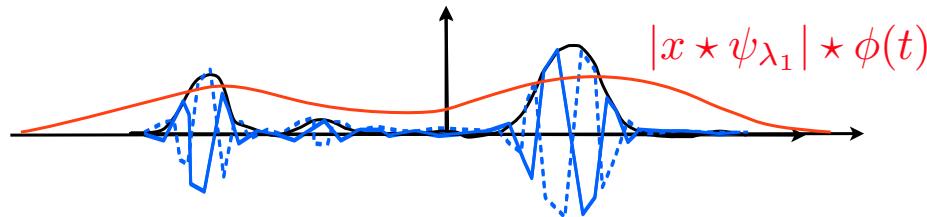
## Wavelet Translation Invariance

$$|x \star \psi_{\lambda_1}(t)| = \sqrt{|x \star \psi_{\lambda_1}^a(t)|^2 + |x \star \psi_{\lambda_1}^b(t)|^2} \text{ pooling}$$



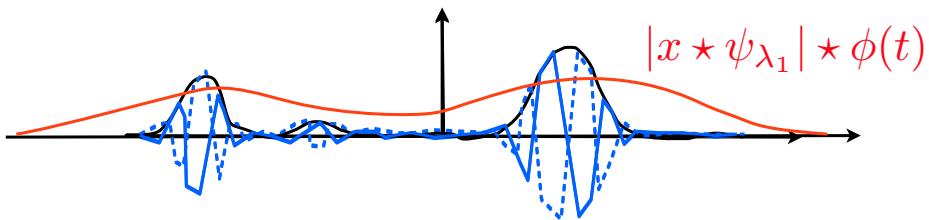
- The modulus  $|x \star \psi_{\lambda_1}|$  is a regular envelop

## Wavelet Translation Invariance



- The modulus  $|x * \psi_{\lambda_1}|$  is a regular envelop
- The average  $|x * \psi_{\lambda_1}| * \phi(t)$  is invariant to small translations relatively to the support of  $\phi$ .

## Wavelet Translation Invariance

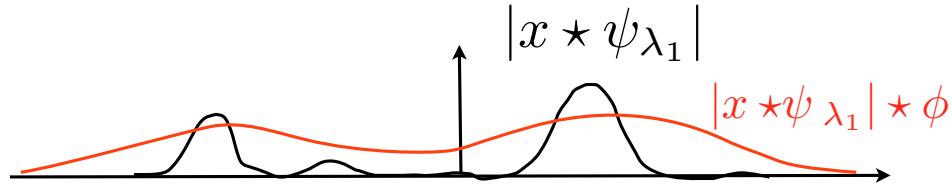


- The modulus  $|x \star \psi_{\lambda_1}|$  is a regular envelop
- The average  $|x \star \psi_{\lambda_1}| \star \phi(t)$  is invariant to small translations relatively to the support of  $\phi$ .
- Full translation invariance at the limit:

$$\lim_{\phi \rightarrow 1} |x \star \psi_{\lambda_1}| \star \phi(t) = \int |x \star \psi_{\lambda_1}(u)| du = \|x \star \psi_{\lambda_1}\|_1$$

but few invariants.

# Recovering Lost Information



- The high frequencies of  $|x * \psi_{\lambda_1}|$  are in wavelet coefficients:

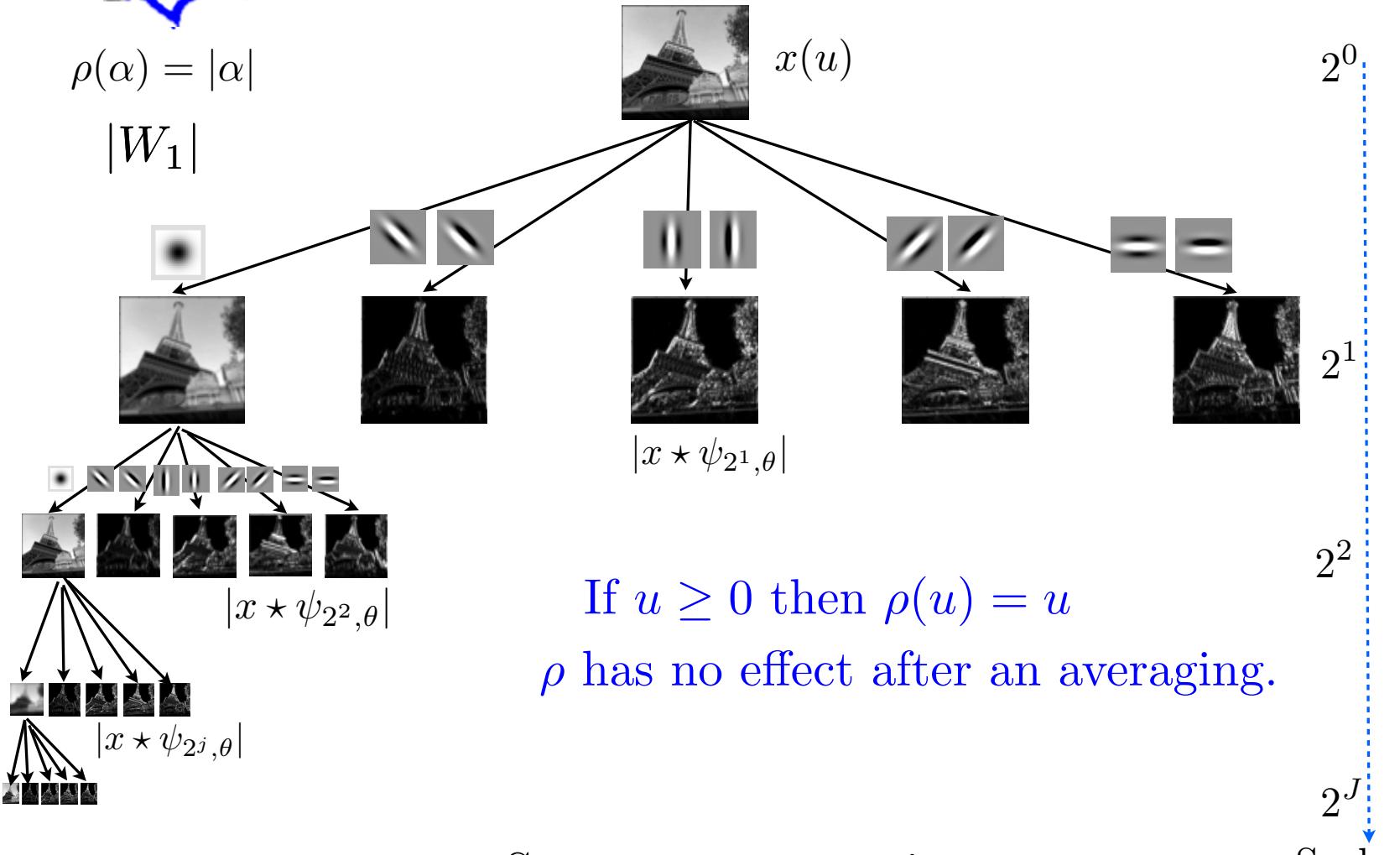
$$W|x * \psi_{\lambda_1}| = \begin{pmatrix} |x * \psi_{\lambda_1}| * \phi(t) \\ |x * \psi_{\lambda_1}| * \psi_{\lambda_2}(t) \end{pmatrix}_{t, \lambda_2}$$

- Translation invariance by time averaging the amplitude:

$$\forall \lambda_1, \lambda_2, \quad ||x * \psi_{\lambda_1}| * \psi_{\lambda_2}| * \phi(t)$$



# Wavelet Filter Bank





# Contraction

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

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

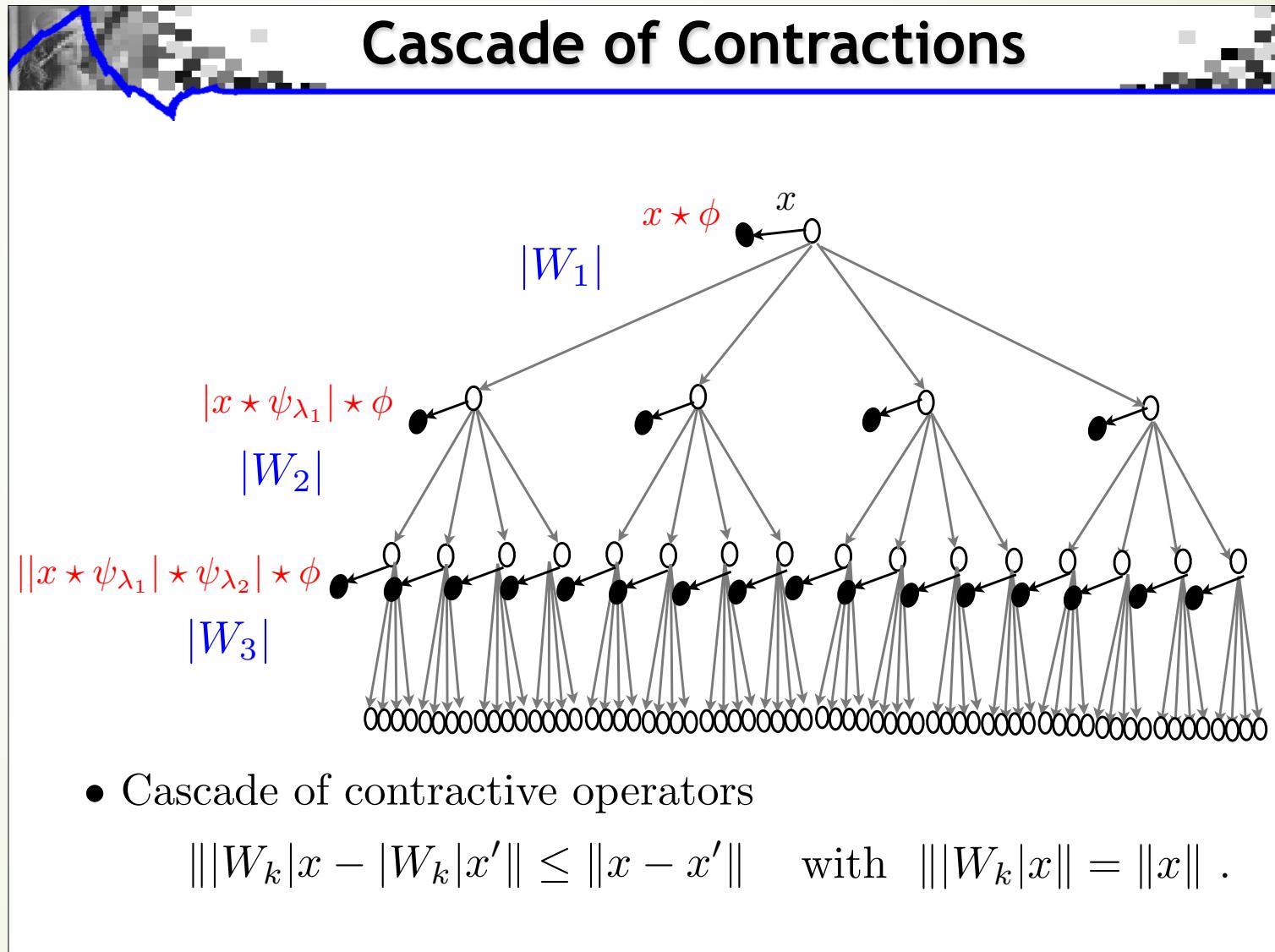
$$|W|x = \begin{pmatrix} x \star \phi(t) \\ |x \star \psi_\lambda(t)| \end{pmatrix}_{t,\lambda} \text{ is non-linear}$$

- 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  $\||W|x\| = \|x\|$

# Wavelet Scattering Network



# Stability of Wavelet Scattering Transform



$$Sx = \begin{pmatrix} x * \phi(u) \\ |x * \psi_{\lambda_1}| * \phi(u) \\ ||x * \psi_{\lambda_1}| * \psi_{\lambda_2}| * \phi(u) \\ |||x * \psi_{\lambda_1}| * \psi_{\lambda_2}| * \psi_{\lambda_3}| * \phi(u) \\ \dots \\ u, \lambda_1, \lambda_2, \lambda_3, \dots \end{pmatrix}$$

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

contractive  $\|Sx - Sy\| \leq \|x - y\|$

preserves norms  $\|Sx\| = \|x\|$

stable to deformations  $x_\tau(t) = x(t - \tau(t))$

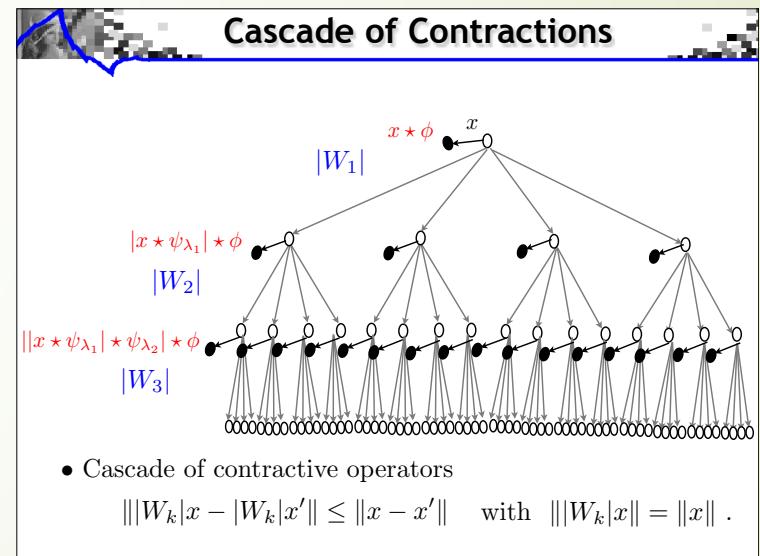
$$\|Sx - Sx_\tau\| \leq C \sup_t |\nabla \tau(t)| \|x\|$$

$\Rightarrow$  linear discriminative classification from  $\Phi x = Sx$

# Summary: Wavelet Scattering Net

- ▶ Architecture:
  - ▶ Convolutional filters: band-limited wavelets
  - ▶ Nonlinear activation: modulus (Lipschitz)
  - ▶ Pooling: L1 norm as averaging
- ▶ Properties:
  - ▶ A Multiscale Sparse Representation
  - ▶ Norm Preservation (Parseval's identity):
 
$$\|Sx\| = \|x\|$$
- ▶ Contraction:
 
$$\|Sx - Sy\| \leq \|x - y\|$$

$$Sx = \begin{pmatrix} x * \phi(u) \\ |x * \psi_{\lambda_1}| * \phi(u) \\ ||x * \psi_{\lambda_1}| * \psi_{\lambda_2}| * \phi(u) \\ |||x * \psi_{\lambda_2}| * \psi_{\lambda_2}| * \psi_{\lambda_3}| * \phi(u) \\ \dots \end{pmatrix}_{u, \lambda_1, \lambda_2, \lambda_3, \dots}$$



# Invariants/Stability of Scattering Net

## ► Translation Invariance:

- The average  $|x \star \psi_{\lambda_1}| \star \phi(t)$  is invariant to small translations relatively to the support of  $\phi$ .
- Full translation invariance at the limit:

$$\lim_{\phi \rightarrow 1} |x \star \psi_{\lambda_1}| \star \phi(t) = \int |x \star \psi_{\lambda_1}(u)| du = \|x \star \psi_{\lambda_1}\|_1$$

## ► Stable Small Deformations:

*stable to deformations*  $x_\tau(t) = x(t - \tau(t))$

$$\|Sx - Sx_\tau\| \leq C \sup_t |\nabla \tau(t)| \|x\|$$

# Feature Extraction

## Linearized Classification

Joan Bruna

- Each class  $X_k$  is represented by a scattering centroid  $E(SX_k)$   
Affine space model  $\mathbf{A}_k = E(SX_k) + \mathbf{V}_k$ . computed with PCA.

MNIST data basis:

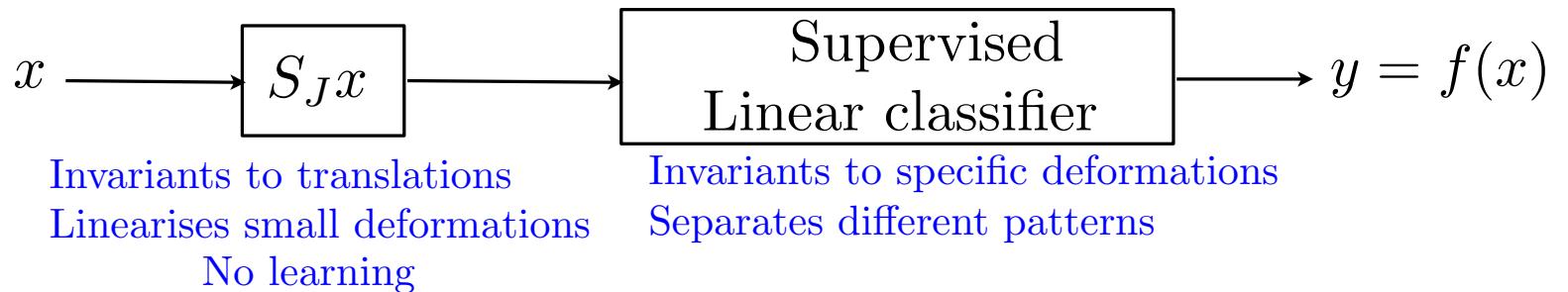
|   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|
| 3 | 6 | 8 | 1 | 7 | 9 | 6 | 6 | 9 | 1 |
| 6 | 7 | 5 | 7 | 8 | 6 | 3 | 4 | 8 | 5 |
| 2 | 1 | 7 | 9 | 7 | 1 | 2 | 8 | 4 | 6 |
| 4 | 8 | 1 | 9 | 0 | 1 | 8 | 8 | 9 | 4 |

# Digit Classification: MNIST



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 6  
4 8 1 9 0 1 8 8 9 4

Joan Bruna



## Classification Errors

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

LeCun et. al.