

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



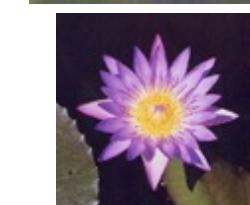
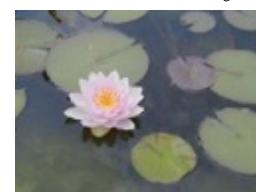
Beaver



Lotus



Water Lily

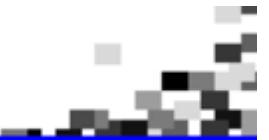


Huge variability
inside classes

Find invariants



High Dimensional Learning



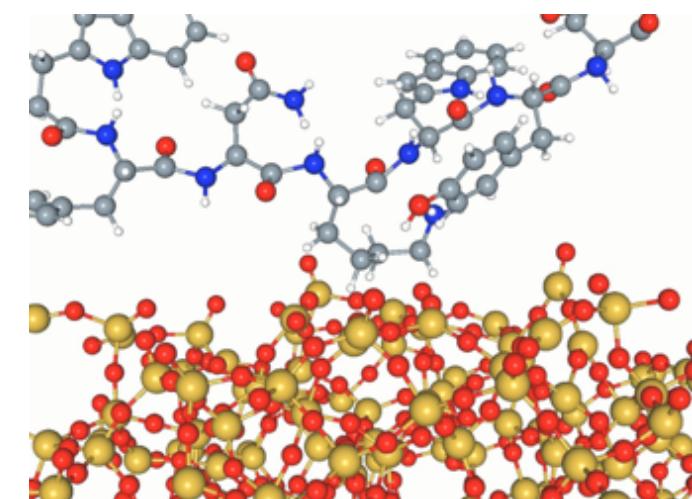
- High-dimensional $x = (x(1), \dots, x(d)) \in \mathbb{R}^d$:
- **Regression:** approximate a *functional* $f(x)$
given n sample values $\{x_i, y_i = f(x_i) \in \mathbb{R}\}_{i \leq n}$

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

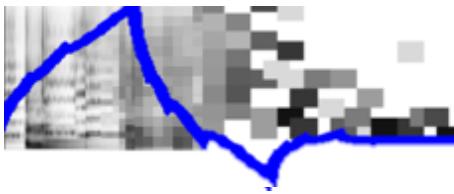
Astronomy



Quantum Chemistry

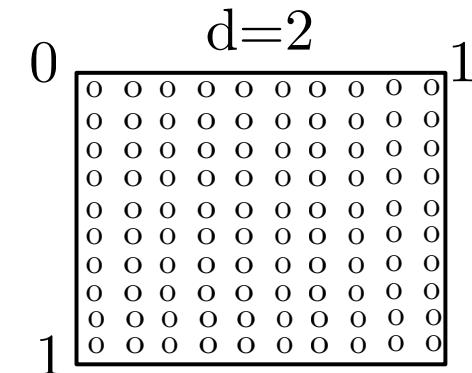
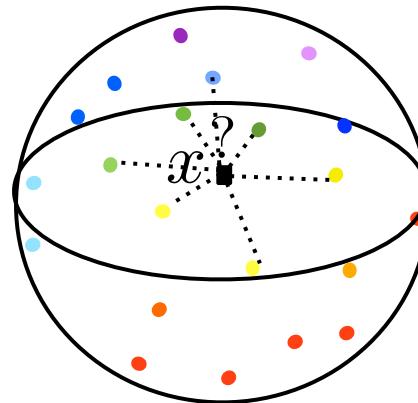


Importance of symmetries.

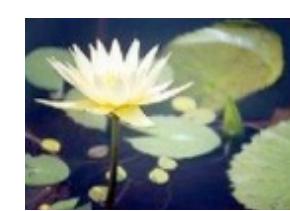


Curse of Dimensionality

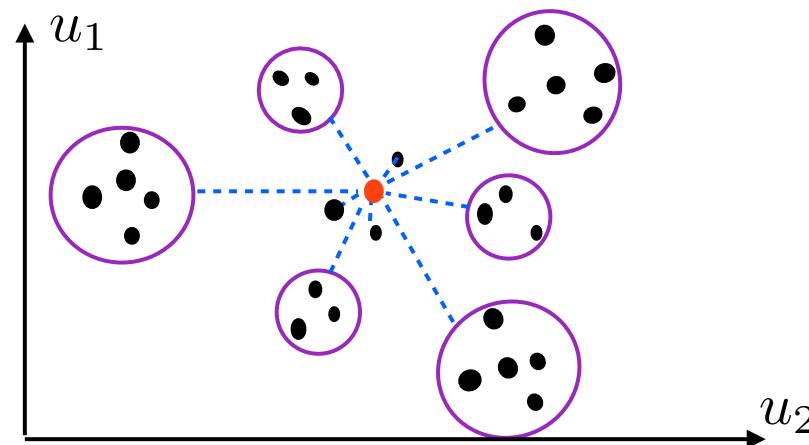
- $f(x)$ can be approximated from examples $\{x_i, f(x_i)\}_i$ by local interpolation if f is regular and there are close examples:



- 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.



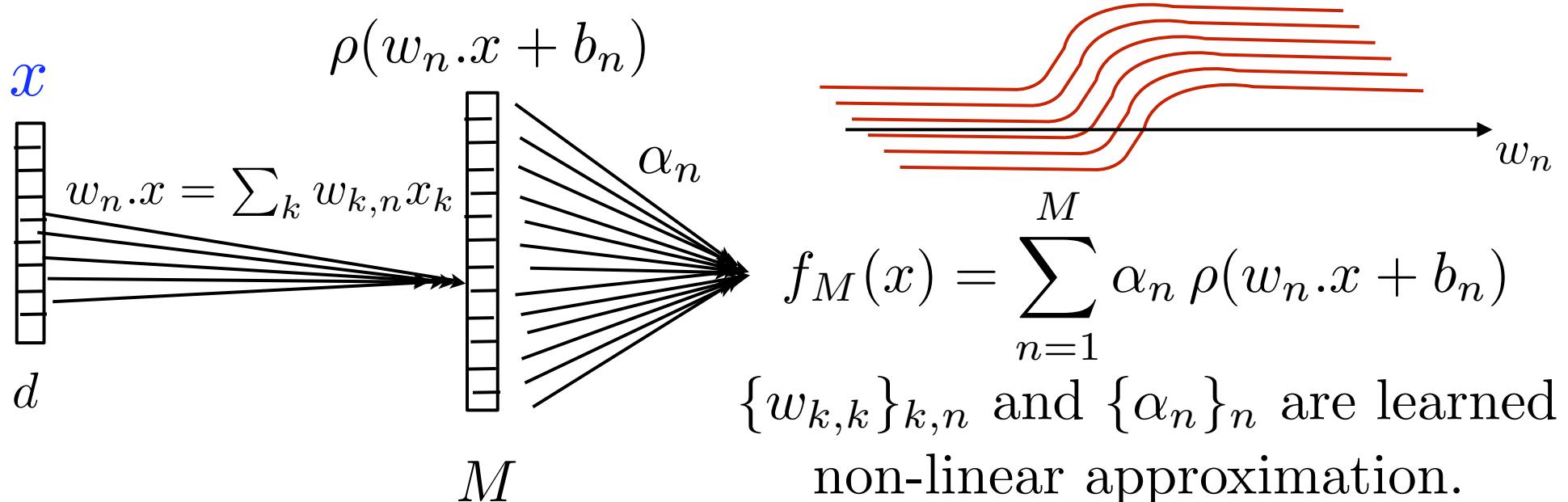
Learning as an Approximation

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



Cybenko, Hornik, Stinchcombe, White

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

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

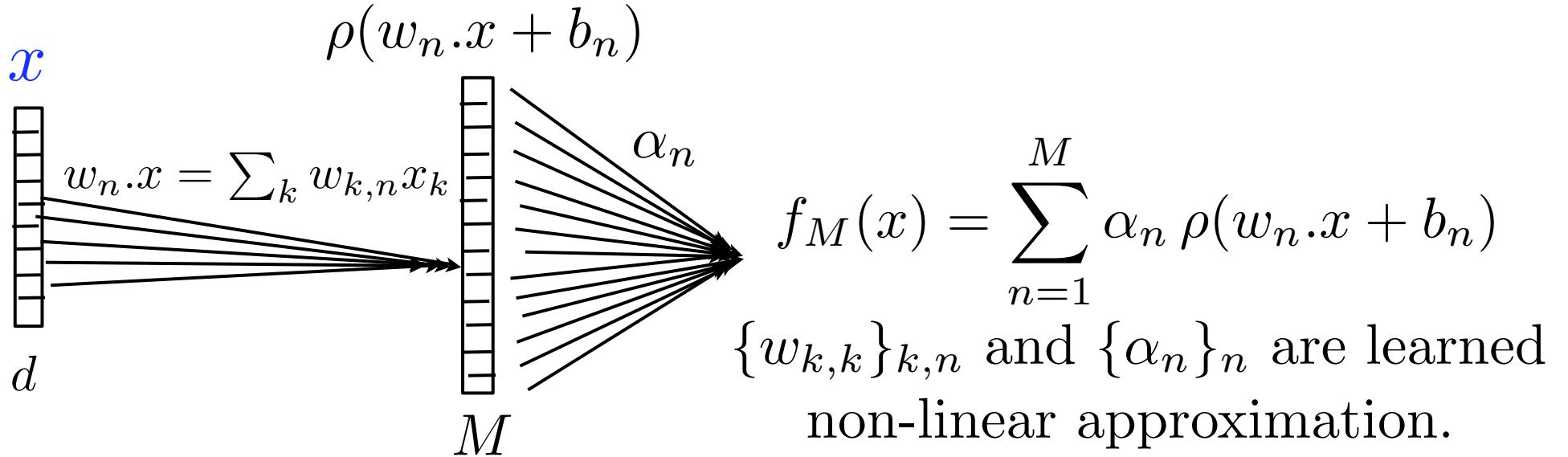
$$\forall f \in \mathbb{L}^2[0, 1]^d \quad \lim_{M \rightarrow \infty} \|f - f_M\| = 0 .$$

No big deal: curse of dimensionality still there.



1 Hidden Layer Neural Networks

One-hidden layer neural network:



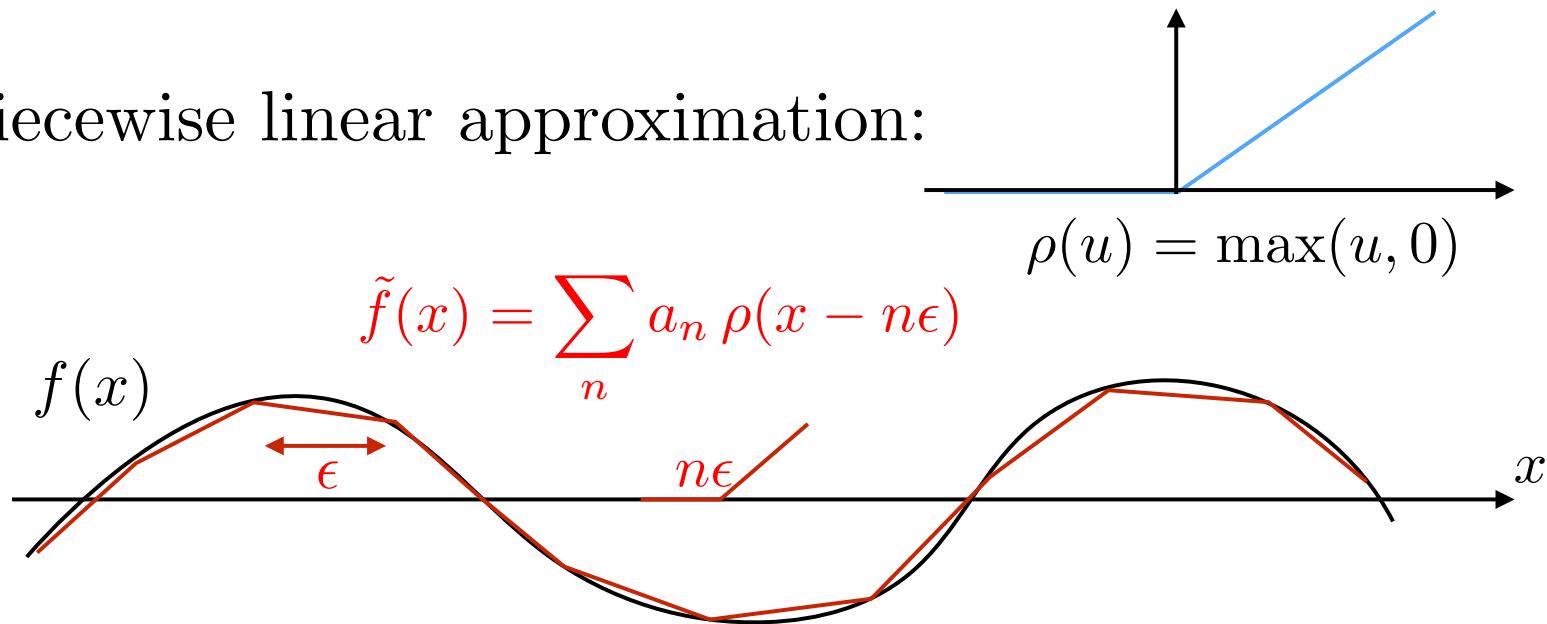
Fourier series: $\rho(u) = e^{iu}$

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

For nearly all ρ : essentially same approximation results.

Piecewise Linear Approximation

- Piecewise linear approximation:



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

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

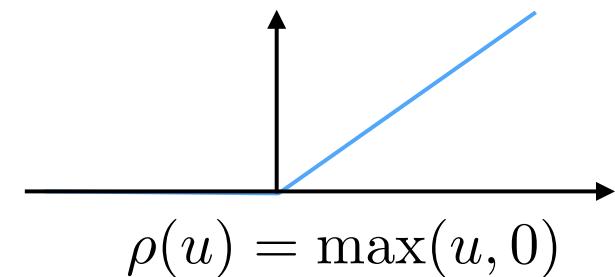
Need $M = \epsilon^{-1}$ points to cover $[0, 1]$ at a distance ϵ

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

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

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

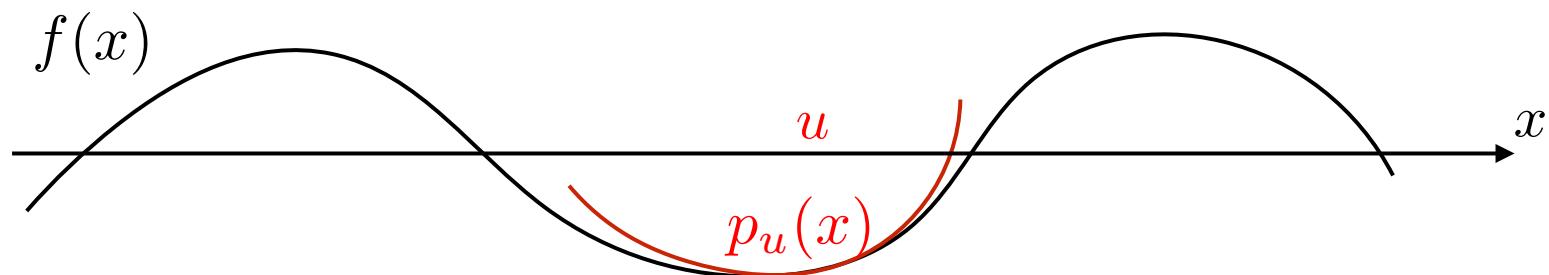
Curse of dimensionality!



Approximation with Regularity

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

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



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

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

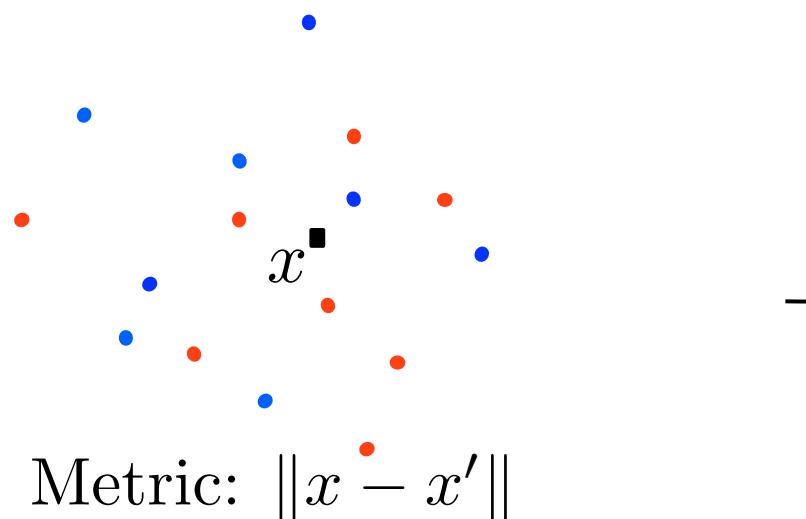
Change of variable $\Phi(x) = \{\phi_k(x)\}_{k \leq d'}$

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

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

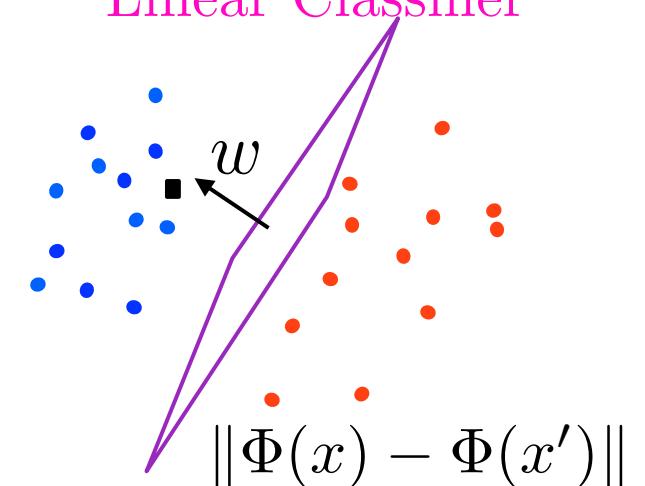
Data: $x \in \mathbb{R}^d$

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



$$\xrightarrow{\Phi}$$

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 $< N$.

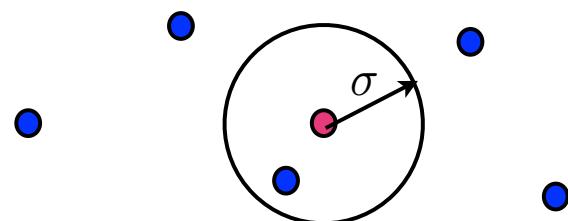
⇒ Choose Φ increasing dimensionality !

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 $\Phi(x)$:

$$\begin{aligned}\Phi(x) \neq \Phi(x') &\text{ if } f(x) \neq f(x') \\ \Rightarrow \exists \tilde{f} \text{ with } f(x) = \tilde{f}(\Phi(x))\end{aligned}$$

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

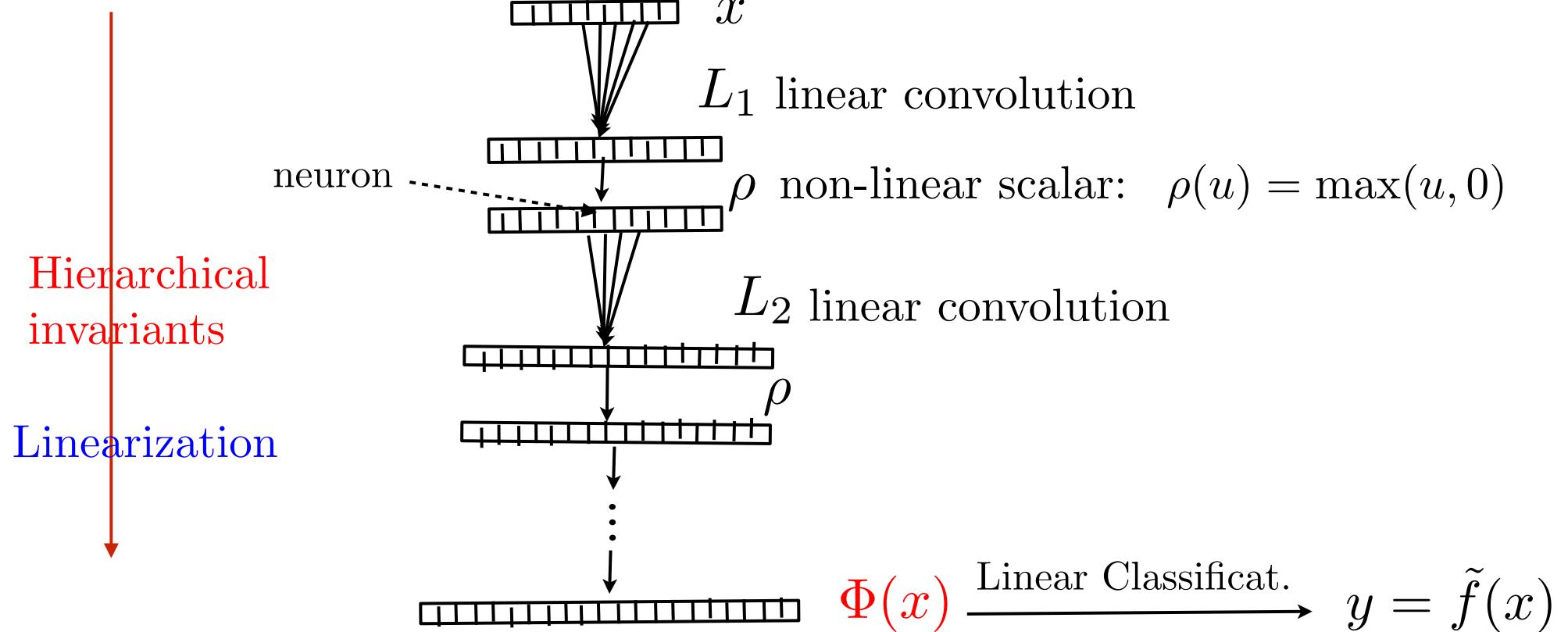
- For $x \in \Omega$, if $\Phi(\Omega)$ is bounded and a low dimension d'

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



Deep Convolution Networks

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



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

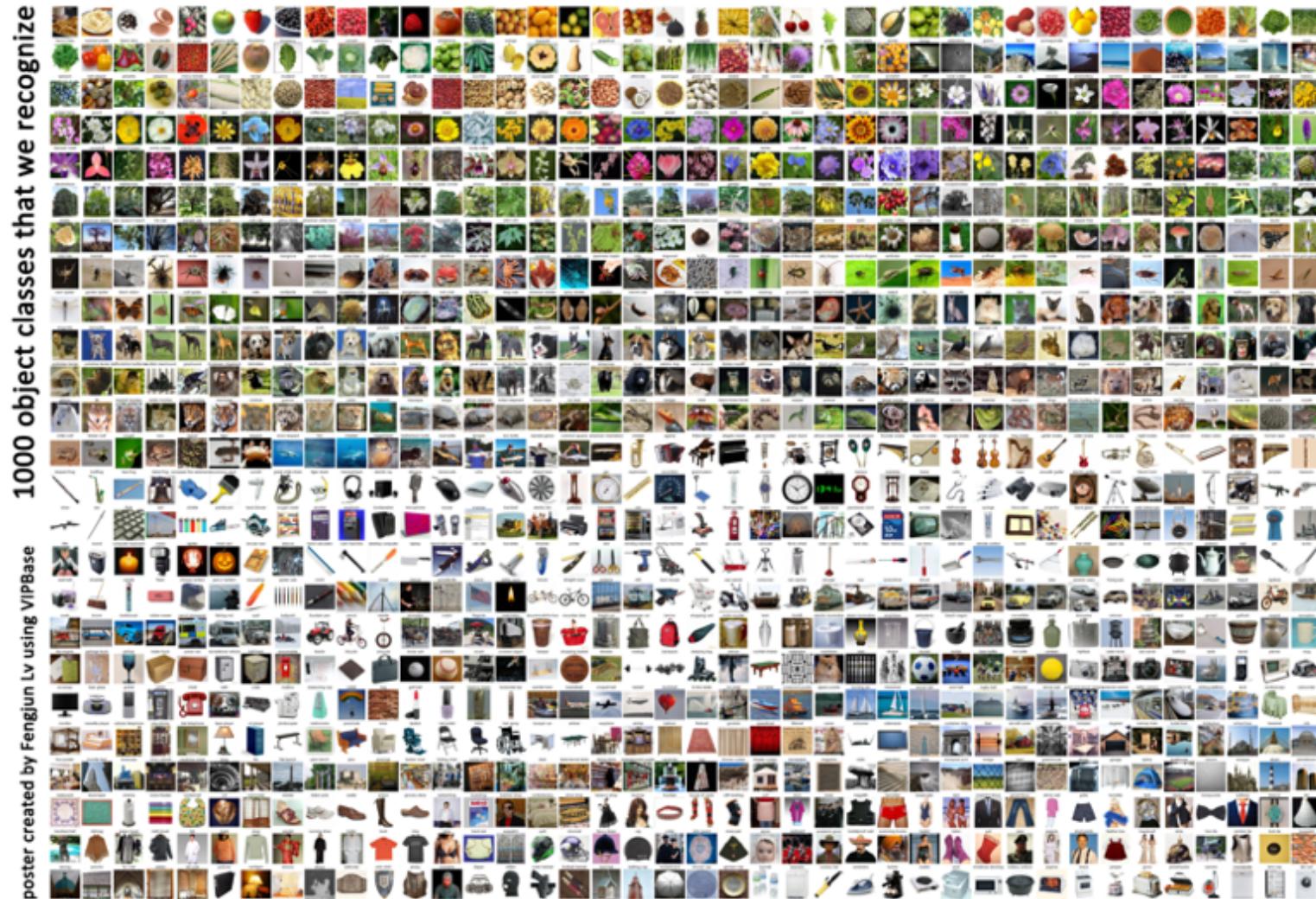
Exceptional results for *images, speech, language, bio-data...*

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



ImageNet Data Basis

- Data basis with 1 million images and 2000 classes



Alex Deep Convolution Network

A. Krizhevsky, Sutskever, Hinton

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

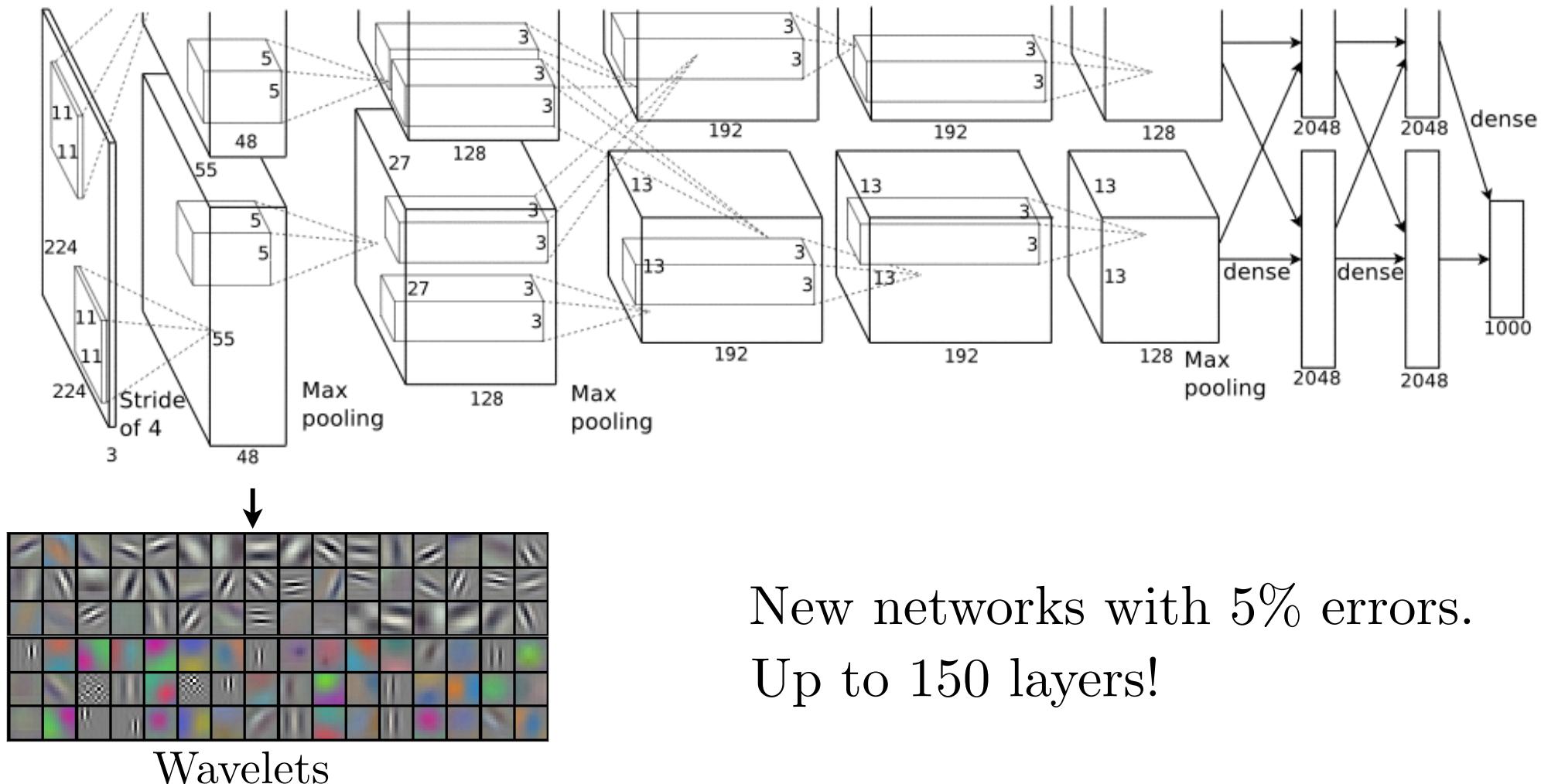
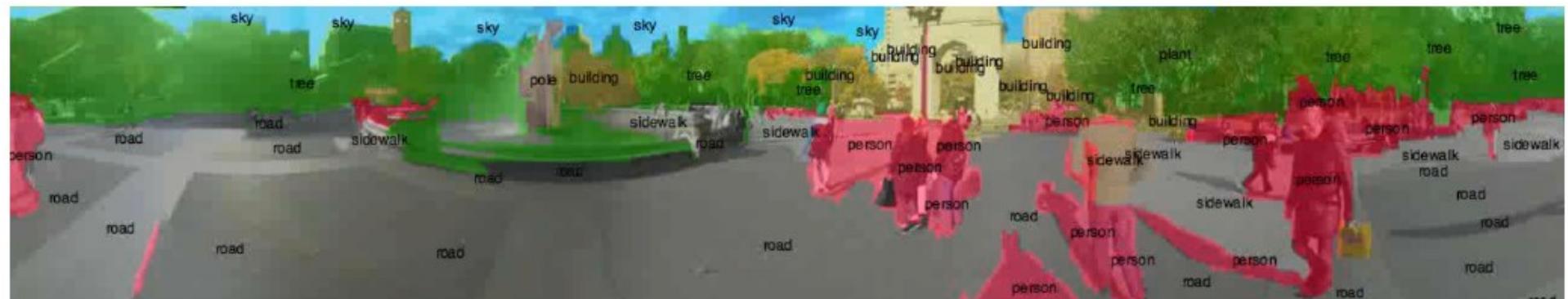


Image Classification

			
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 convertible grille pickup beach wagon fire engine	mushroom agaric mushroom jelly fungus gill fungus dead-man's-fingers	cherry dalmatian grape elderberry ffordshire bullterrier currant	Madagascar cat squirrel monkey spider monkey titi indri howler monkey

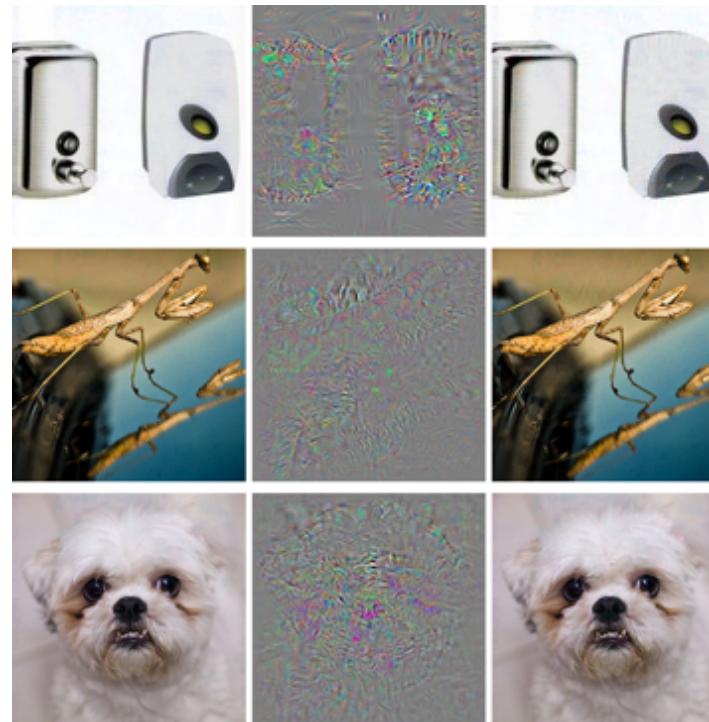
Scene Labeling / Car Driving



Why Understanding ?

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

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

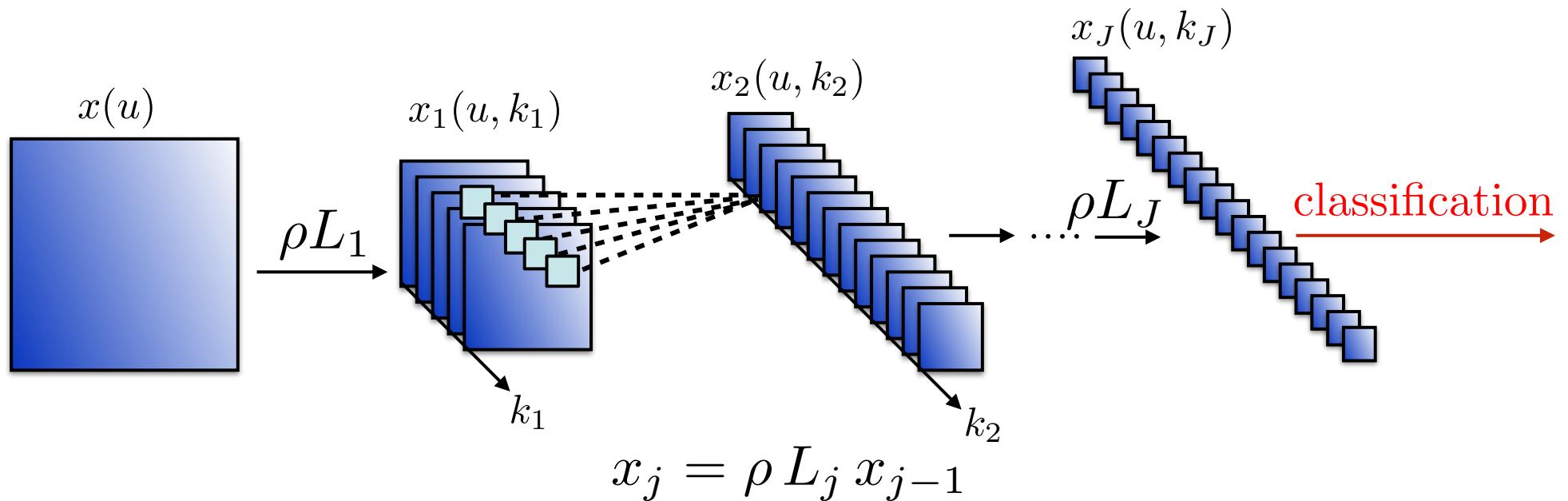


correctly
classified

classified as
ostrich

- Trial and error testing can not guarantee reliability.

Deep Convolutional Networks



- L_j is a linear combination of convolutions and subsampling:

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

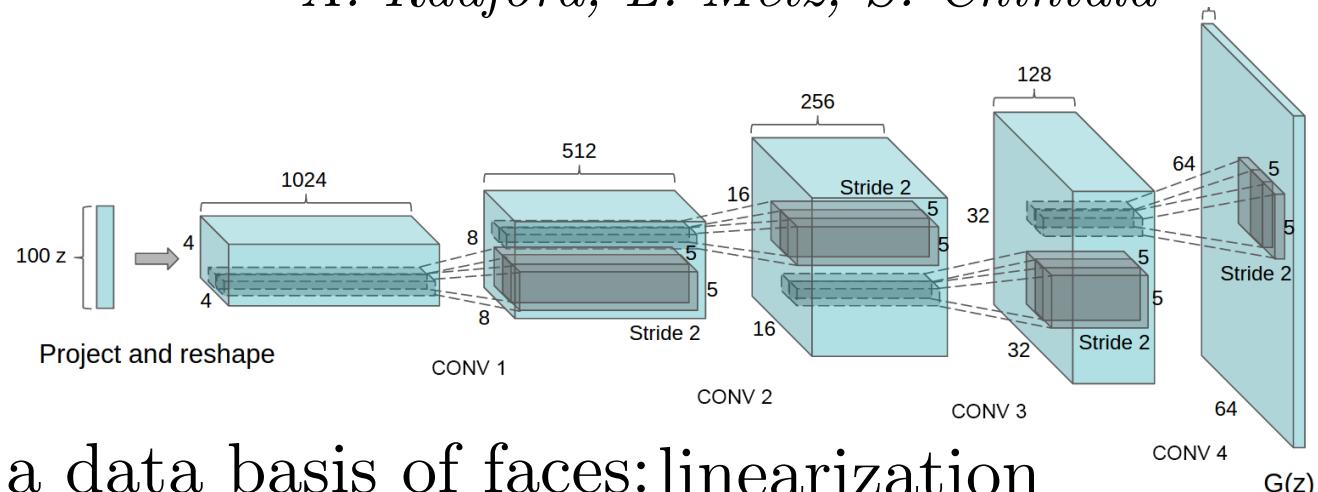
sum across channels

- 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



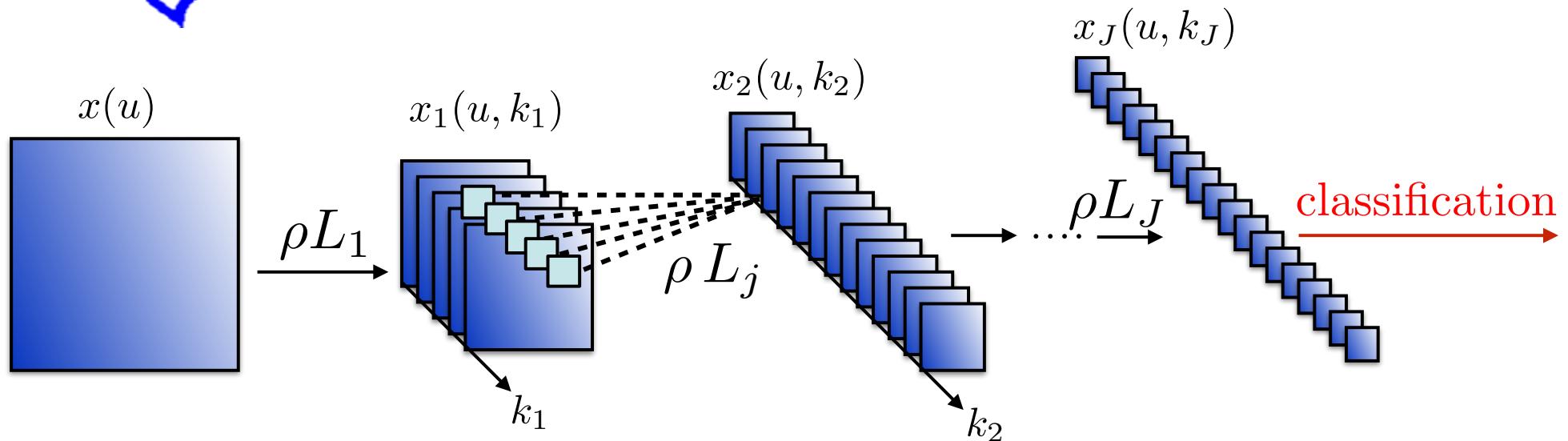
- Trained on a data basis of faces: linearization



- On a data basis including bedrooms: interpolations



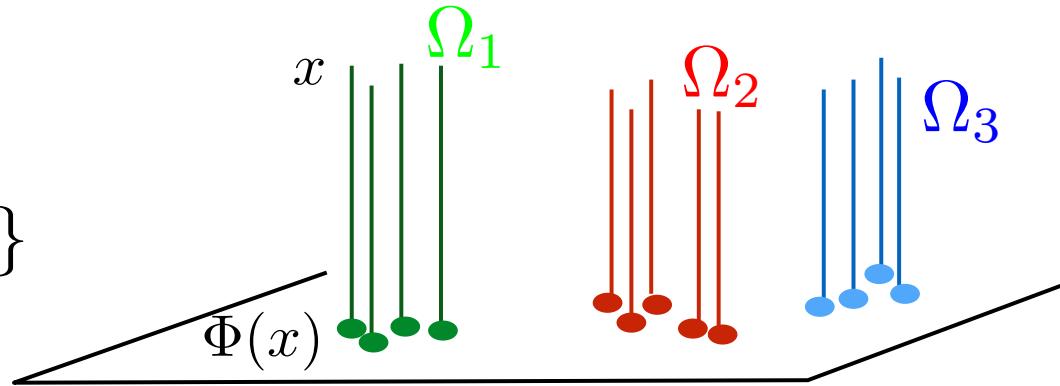
Many Questions



- 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 ?

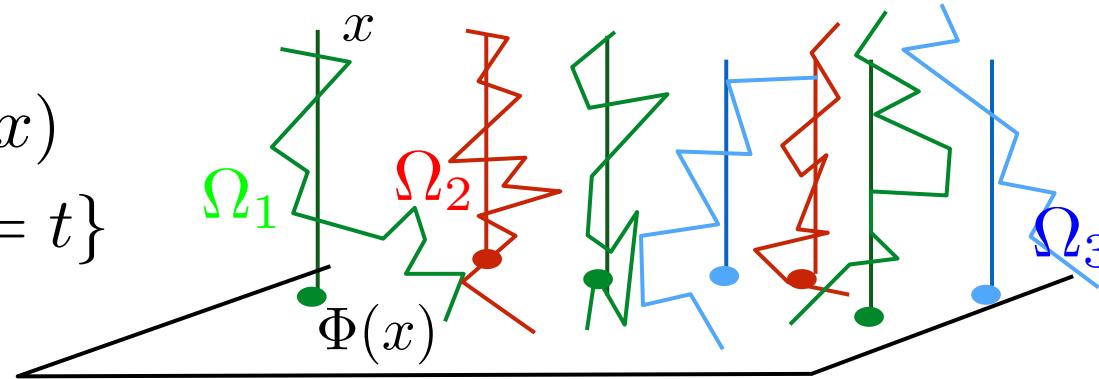
Linear Dimension Reduction

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*.

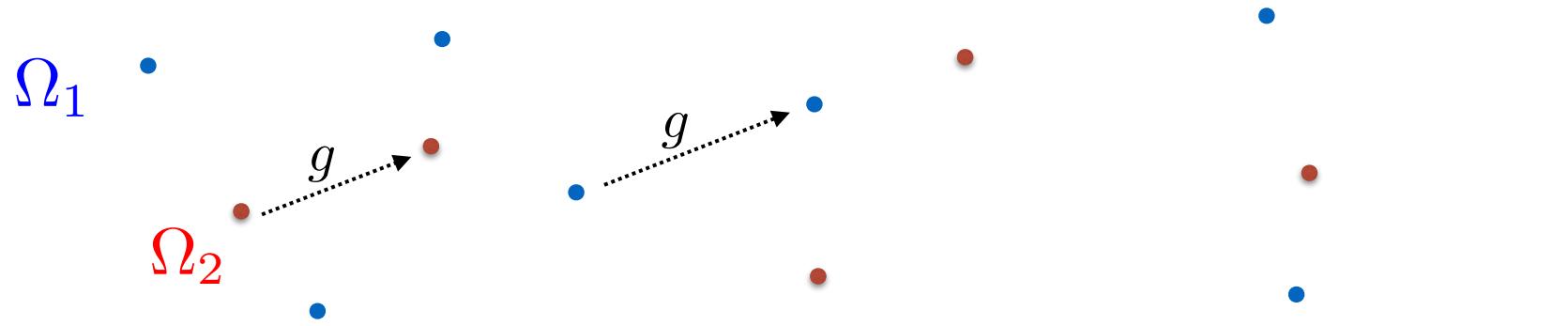
Classes
Level sets of $f(x)$
 $\Omega_t = \{x : f(x) = t\}$



- If level sets 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:

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

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:

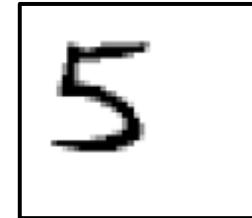
$$g = g_1^{p_1} g_2^{p_2} \dots g_n^{p_n}$$

- Lie group: infinitely small generators (Lie Algebra)

Translation and Deformations

- Digit classification:

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

 Ω_3  Ω_5

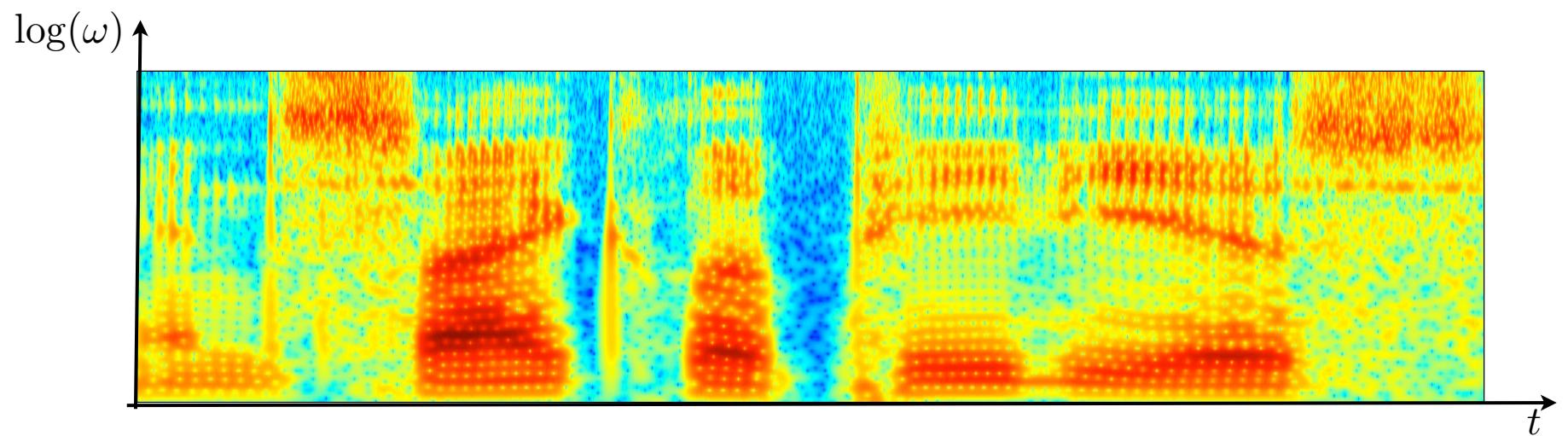
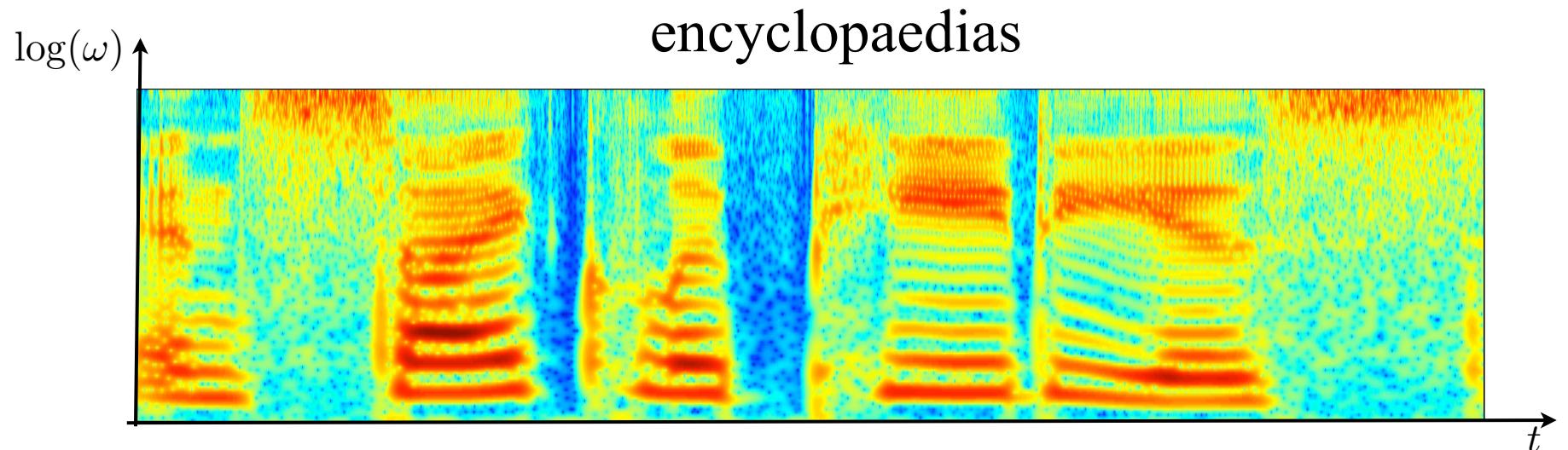
- Globally invariant to the translation group: small
- Locally invariant to small diffeomorphisms: huge group



Video of Philipp Scott Johnson



Frequency Transpositions

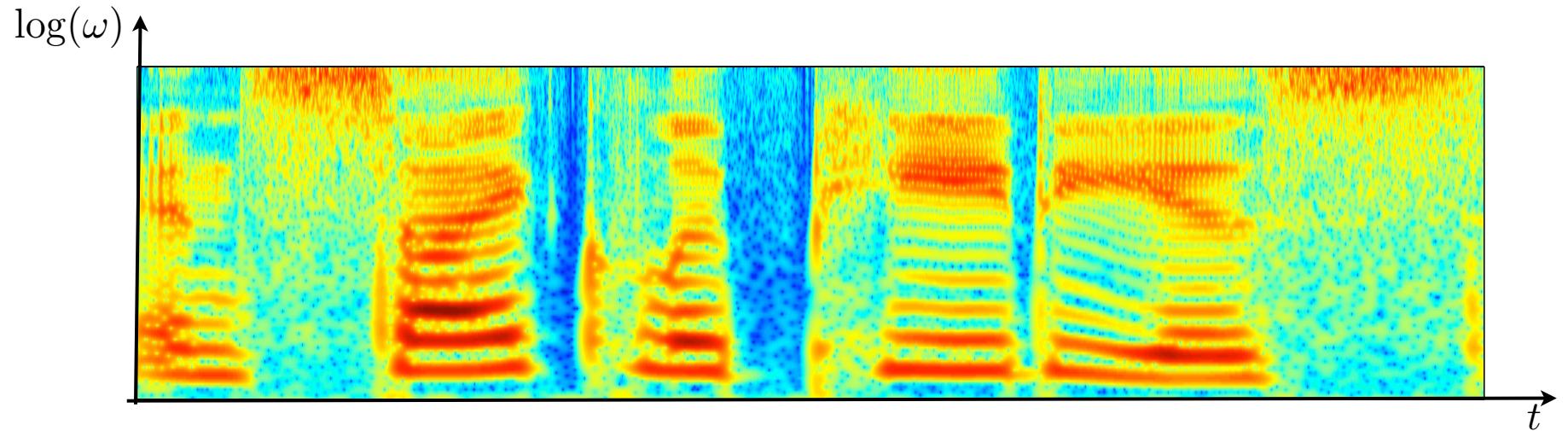


H : Heisenberg group of "time-frequency" translations



Frequency Transpositions

Time and frequency translations and deformations:

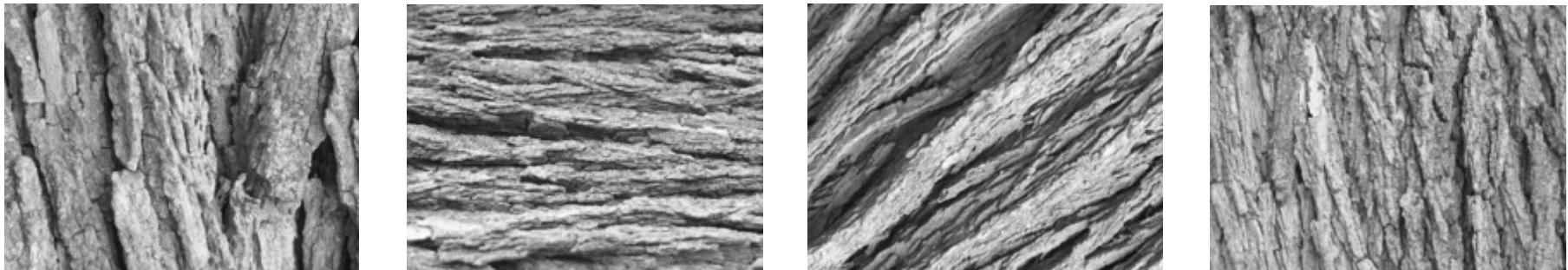


- Frequency transposition invariance is needed for speech recognition not for locutor recognition.



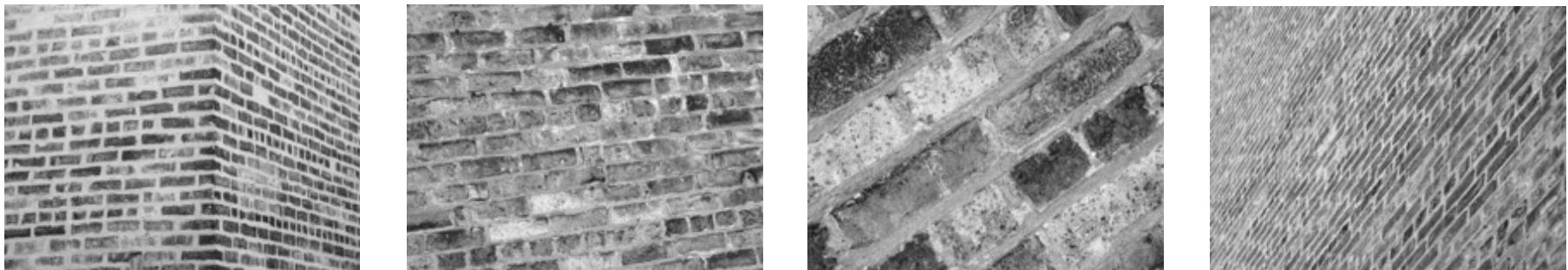
Rotation and Scaling Variability

- Rotation and **deformations**



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

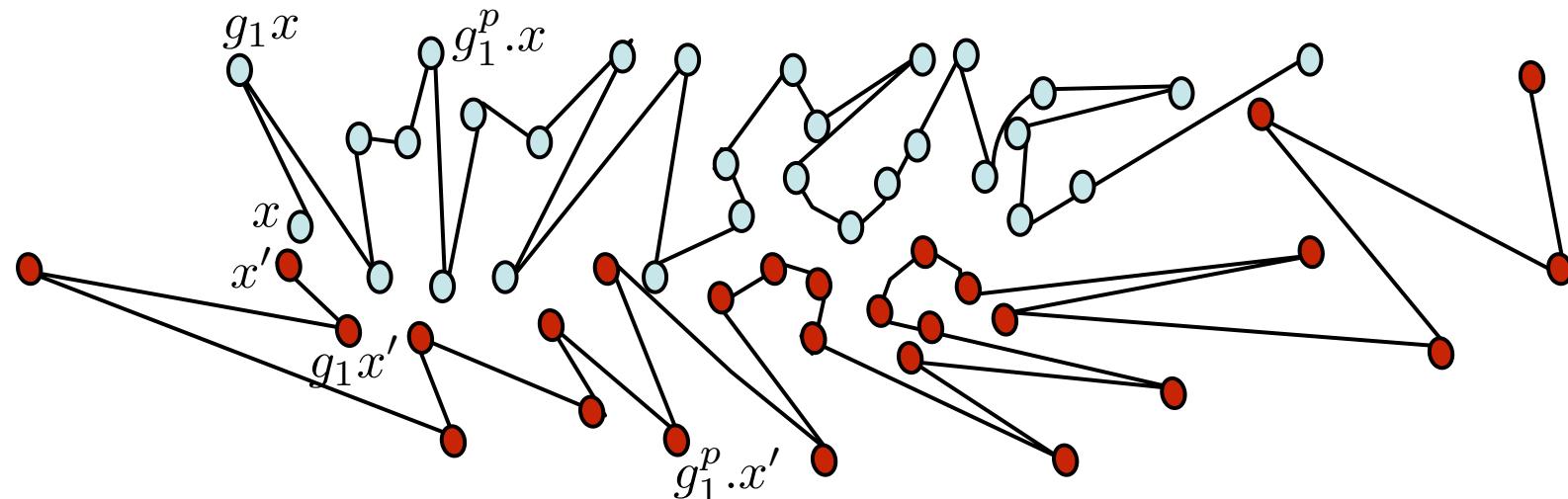
- Scaling and **deformations**



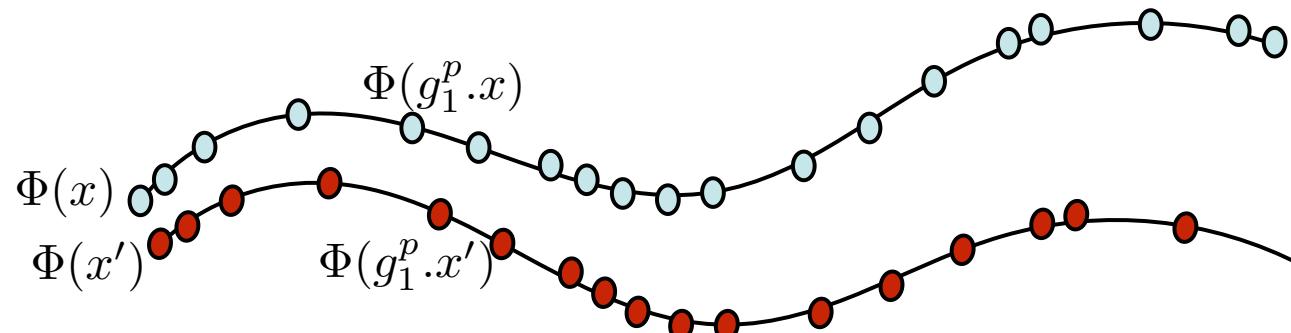
Group: $\mathbb{R} \times \text{Diff}(\mathbb{R})$

Linearize Symmetries

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



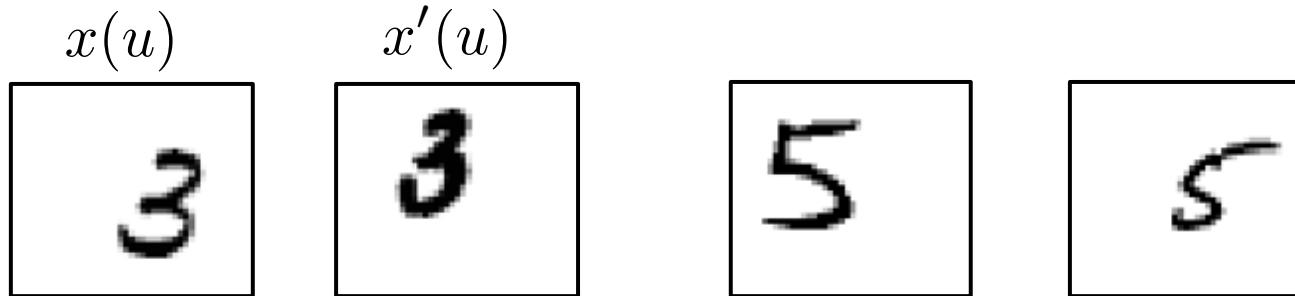
- Linearise symmetries with a change of variable $\Phi(x)$



- 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:

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

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

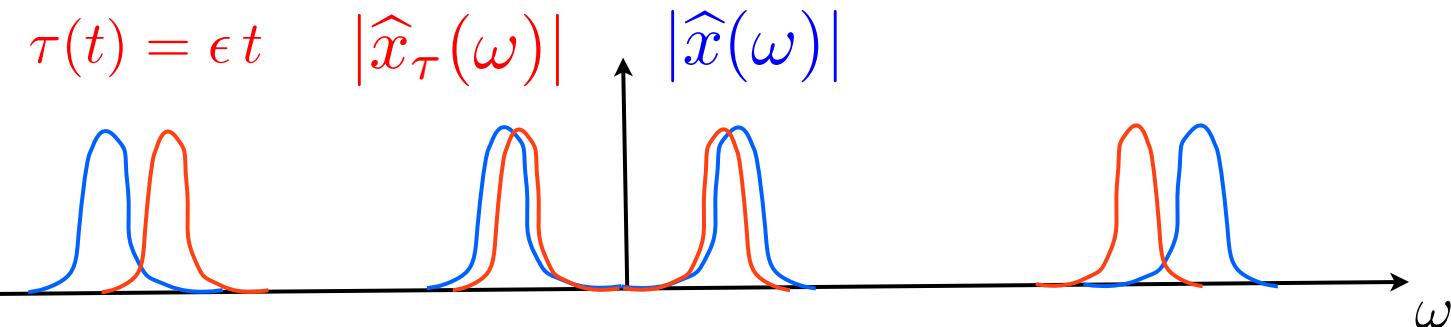
$$x_c(t) = x(t - c) \Rightarrow \hat{x}_c(\omega) = e^{-ic\omega} \hat{x}(\omega)$$

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



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

Deep Neural Network

Mathematical Mysteries

for High Dimensional Learning

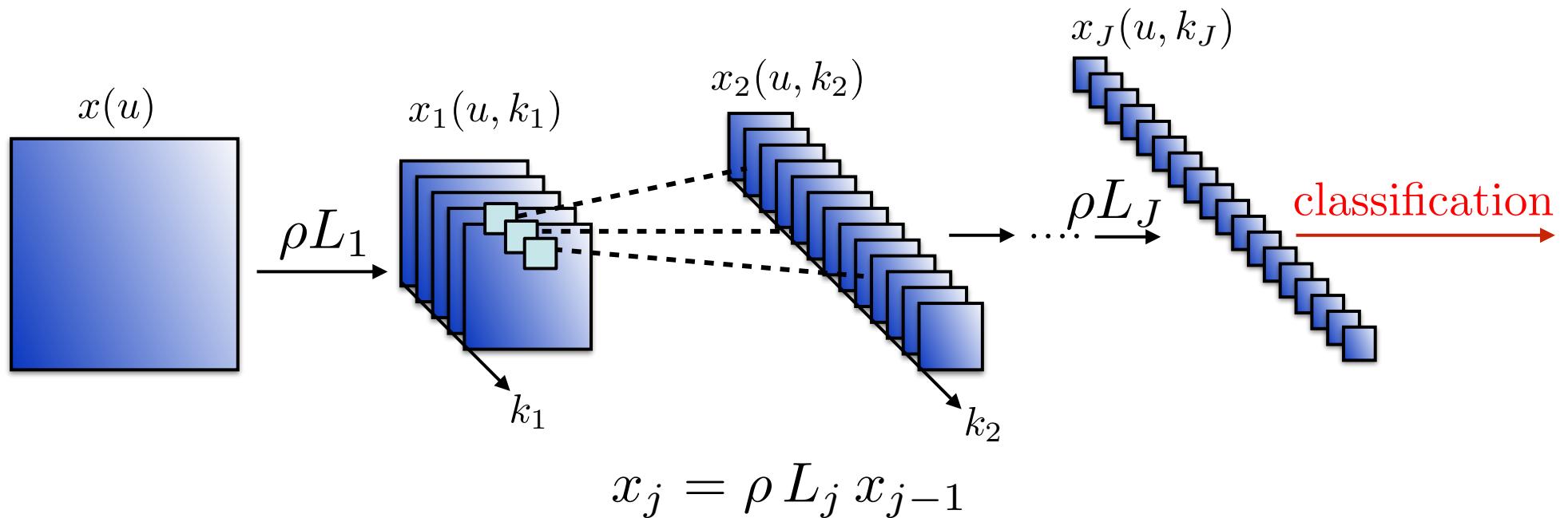


Stéphane Mallat

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:

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

- 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:

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

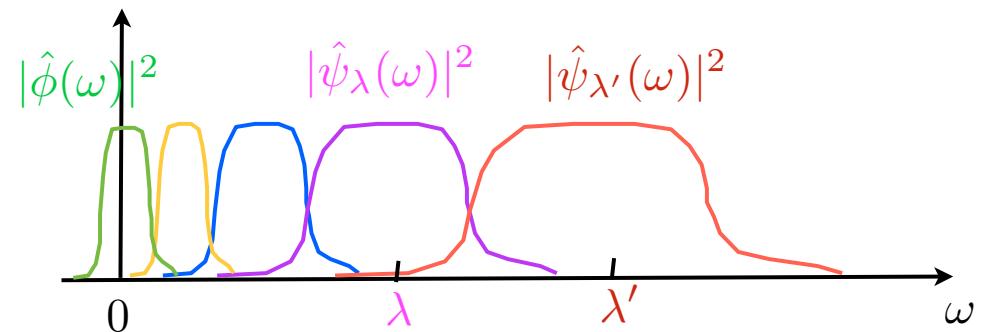
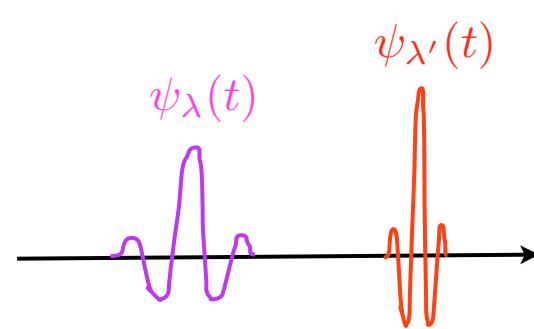
- Wavelet Scattering transform along translations
- Generation of textures and random processes
- Channel connections for more general groups
- Image and audio classification with small training sets
- Quantum chemistry
- Open problems

Understanding Deep Convolutional Networks, arXiv 2016.



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 $\hat{\psi}_\lambda$

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

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

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



Why Wavelets ?

- Wavelets are uniformly stable to deformations:

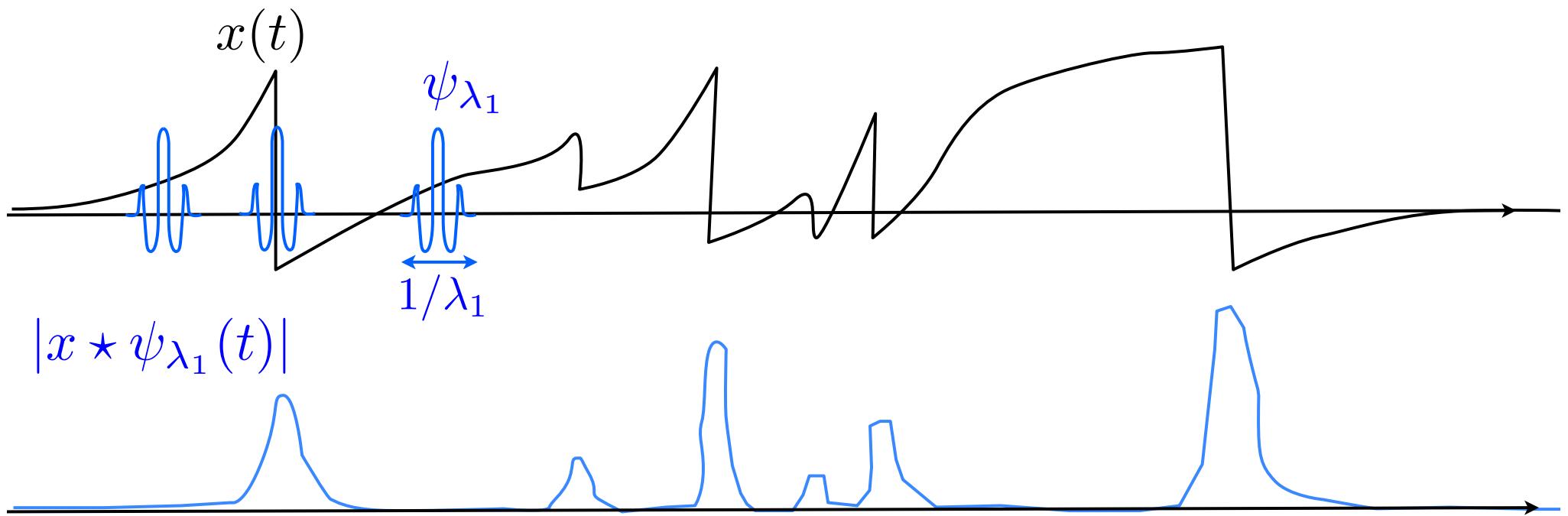
if $\psi_{\lambda,\tau}(t) = \psi_\lambda(t - \tau(t))$ then

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

- 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: $|W|$

