



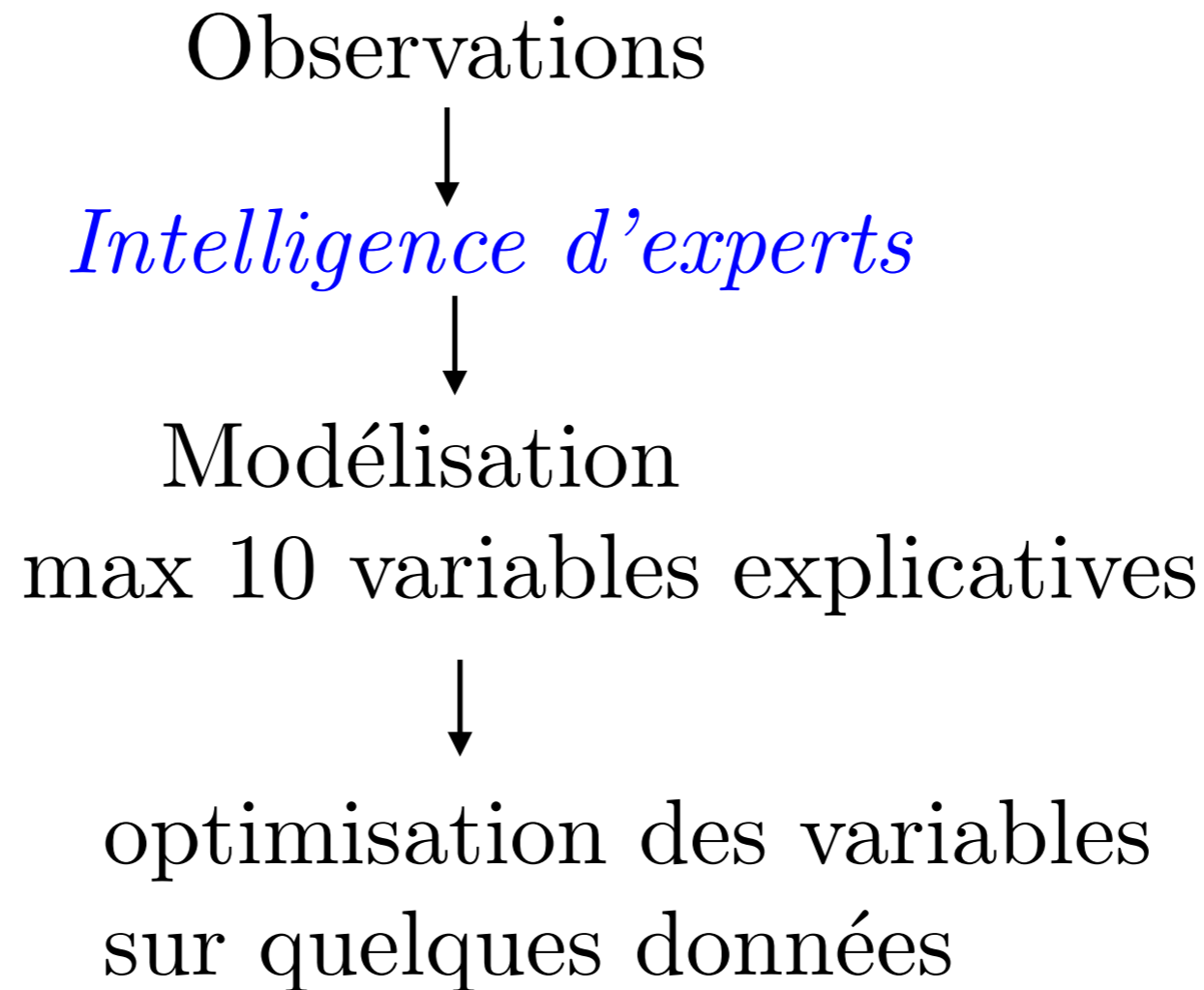
# **La Malédiction Mathématiques des Données en Grande Dimension**

*Stéphane Mallat*

**École Normale Supérieure**

# Un Changement de Paradigme

- Pratique classique de la science et de la connaissance:



- Lent, couteux, et difficile pour des problèmes complexes.

# Un Changement de Paradigme

Masses de données



*Apprentissage automatique*

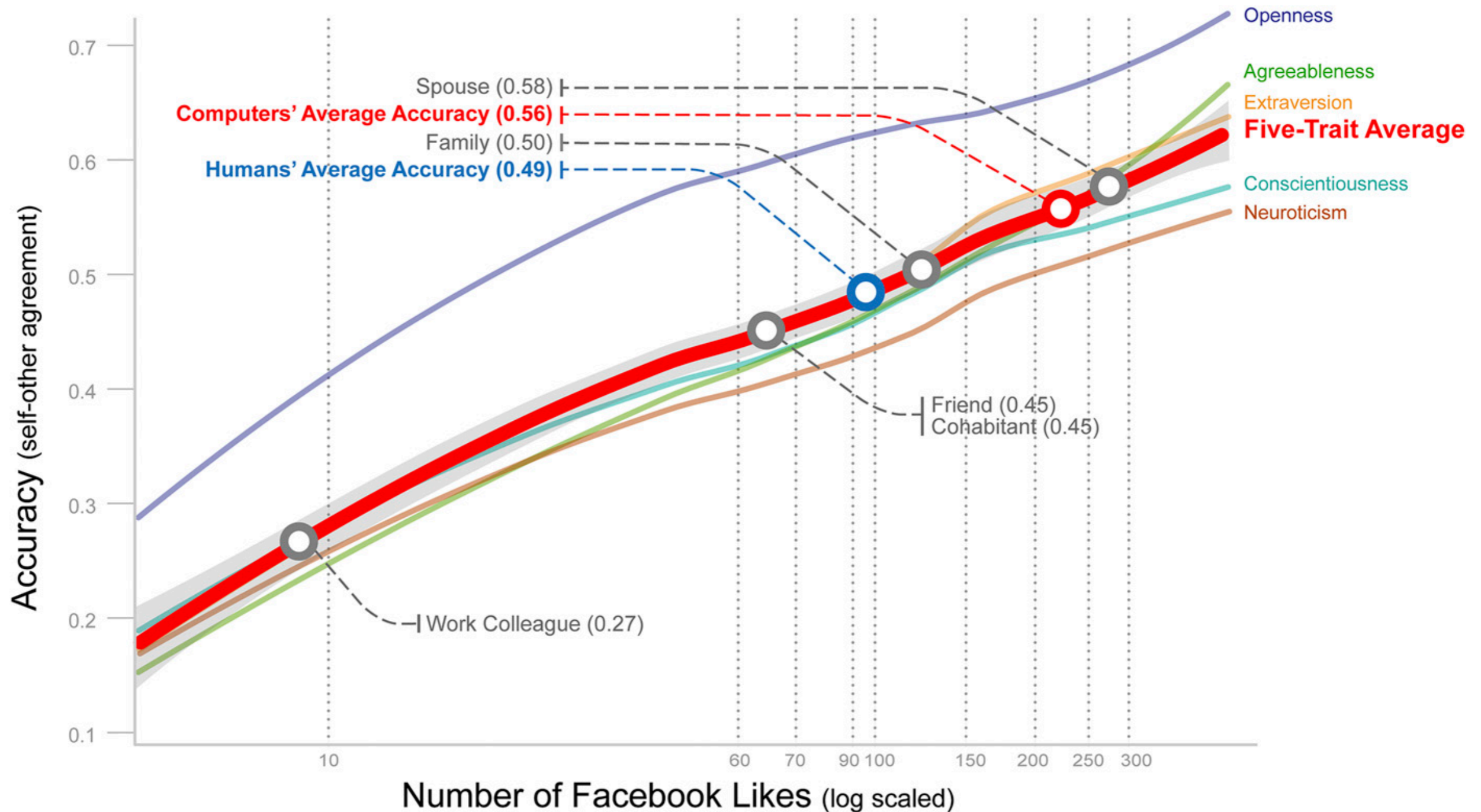


Solution: vote de beaucoup d'indicateurs faibles

- Intelligence: nécessite beaucoup de données et de mémoire.
- Applications:
  - Diagnostiques (médical, industriel)
  - Reconnaissance (images, sons, langage naturel...)
  - Prédiction (sciences sociales, marketing, finance, sciences...)

# Un Changement de Paradigme

Personality Evaluation (Proc. Nat. Acad. Sciences 2014)



# Le Menu

---

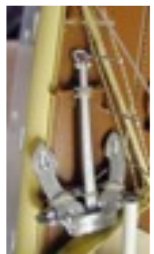
- Malédiction de la dimensionalité
- Enjeux mathématiques de l'apprentissage
- Paysage scientifique et industriel

# High Dimensional Learning

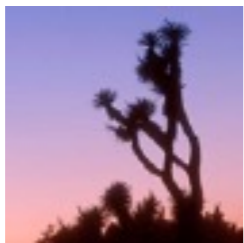
- High-dimensional  $x = (x(1), \dots, x(d)) \in \mathbb{R}^d$ :
- **Classification:** estimate a class label  $f(x)$  given  $n$  sample values  $\{x_i, y_i = f(x_i)\}_{i \leq n}$

Image Classification  $d = 10^6$

Anchor



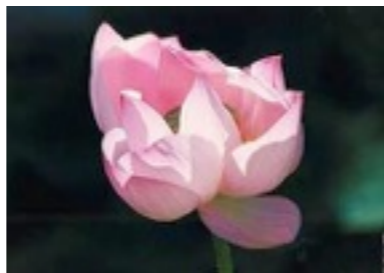
Joshua Tree



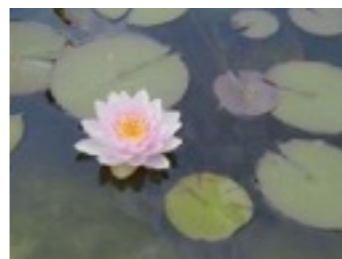
Beaver



Lotus



Water Lily



Huge variability  
inside classes

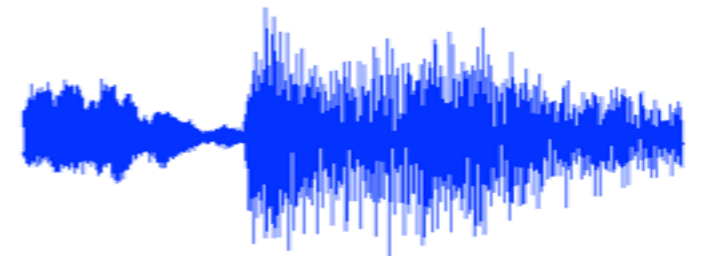
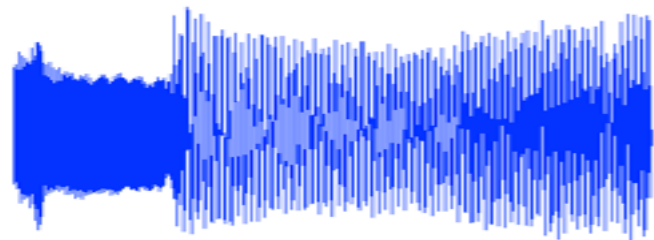
Need to find  
informative  
invariants

# High Dimensional Learning

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

Audio: instrument recognition

Huge variability  
inside classes



# High Dimensional Learning

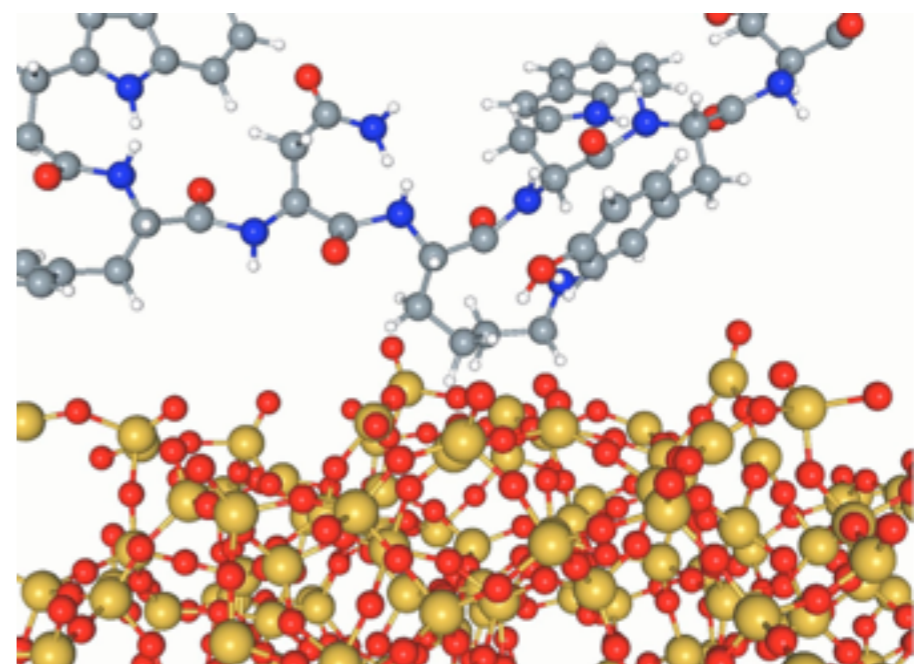
- High-dimensional  $x = (x(1), \dots, x(d)) \in \mathbb{R}^d$ :
- **Regression:** approximate a *functional*  $f(x)$   
given  $n$  sample values  $\{x_i, y_i = f(x_i) \in \mathbb{R}\}_{i \leq n}$

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

Astronomy



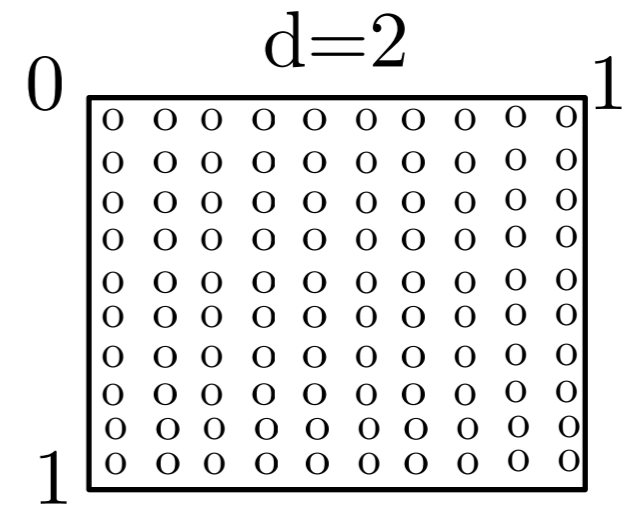
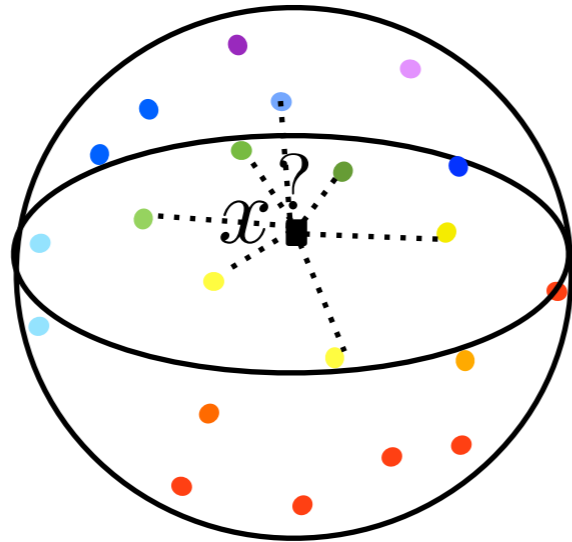
Quantum Chemistry





# Curse of Dimensionality

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



- To cover  $[0, 1]^d$  at a distance  $10^{-1}$  we need  $10^d$  points  
 $\Rightarrow \|x - x_i\|$  is always large



# Learning by Euclidean Embedding

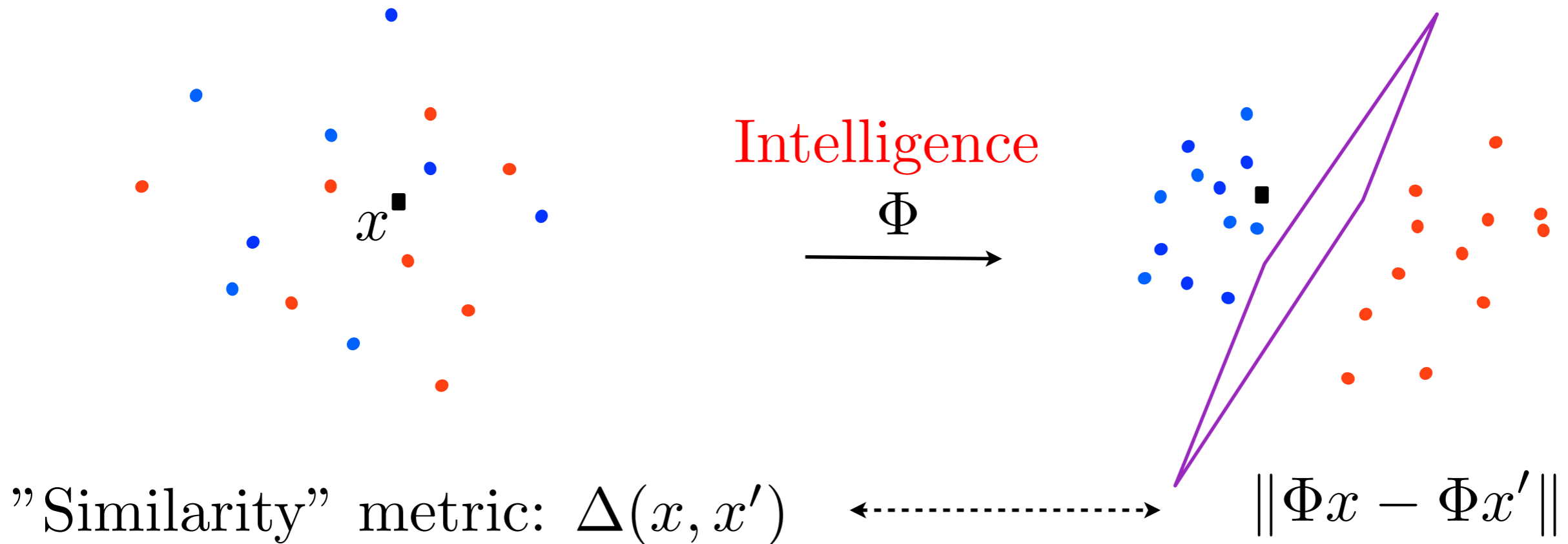
Data:  $x \in \mathbb{R}^d$

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

Representation

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

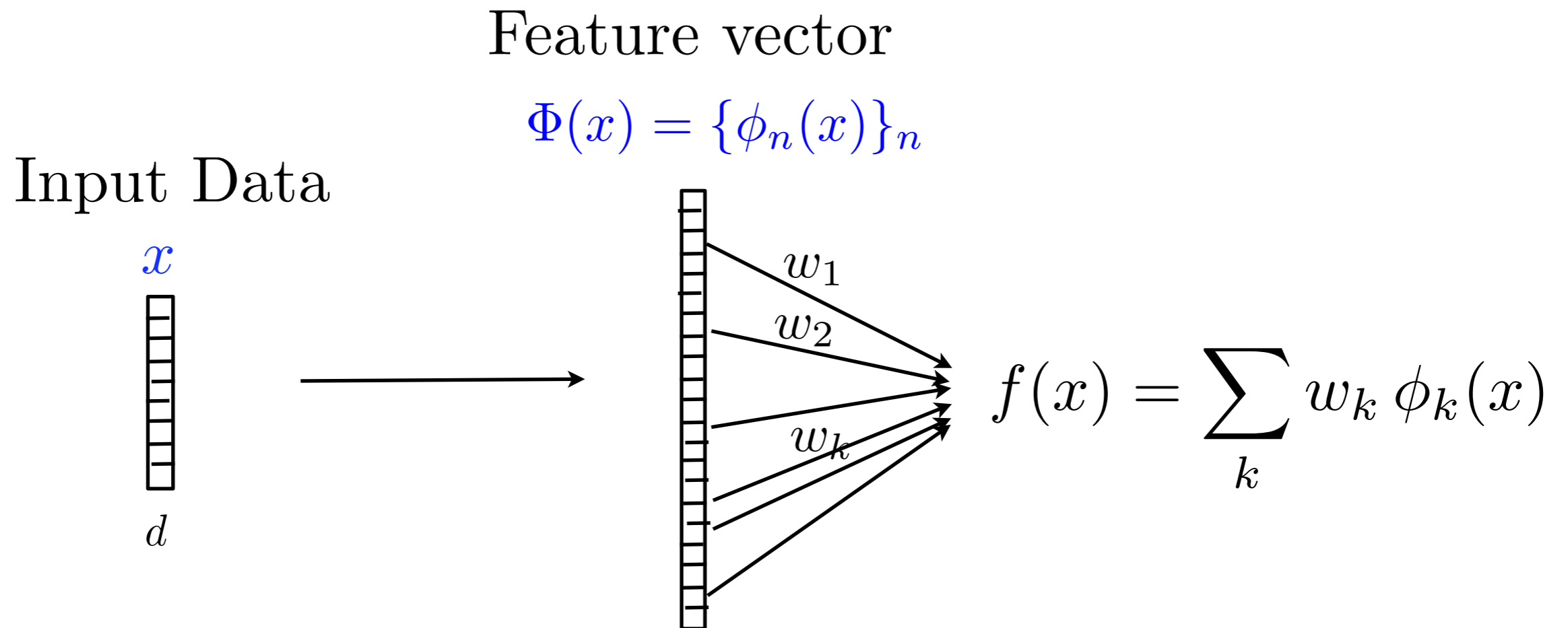
Linear Classifier



Equivalent Euclidean metric:

$$\Delta(x, x') \approx \|\Phi x - \Phi x'\|$$

# Aggregation of Weak Features



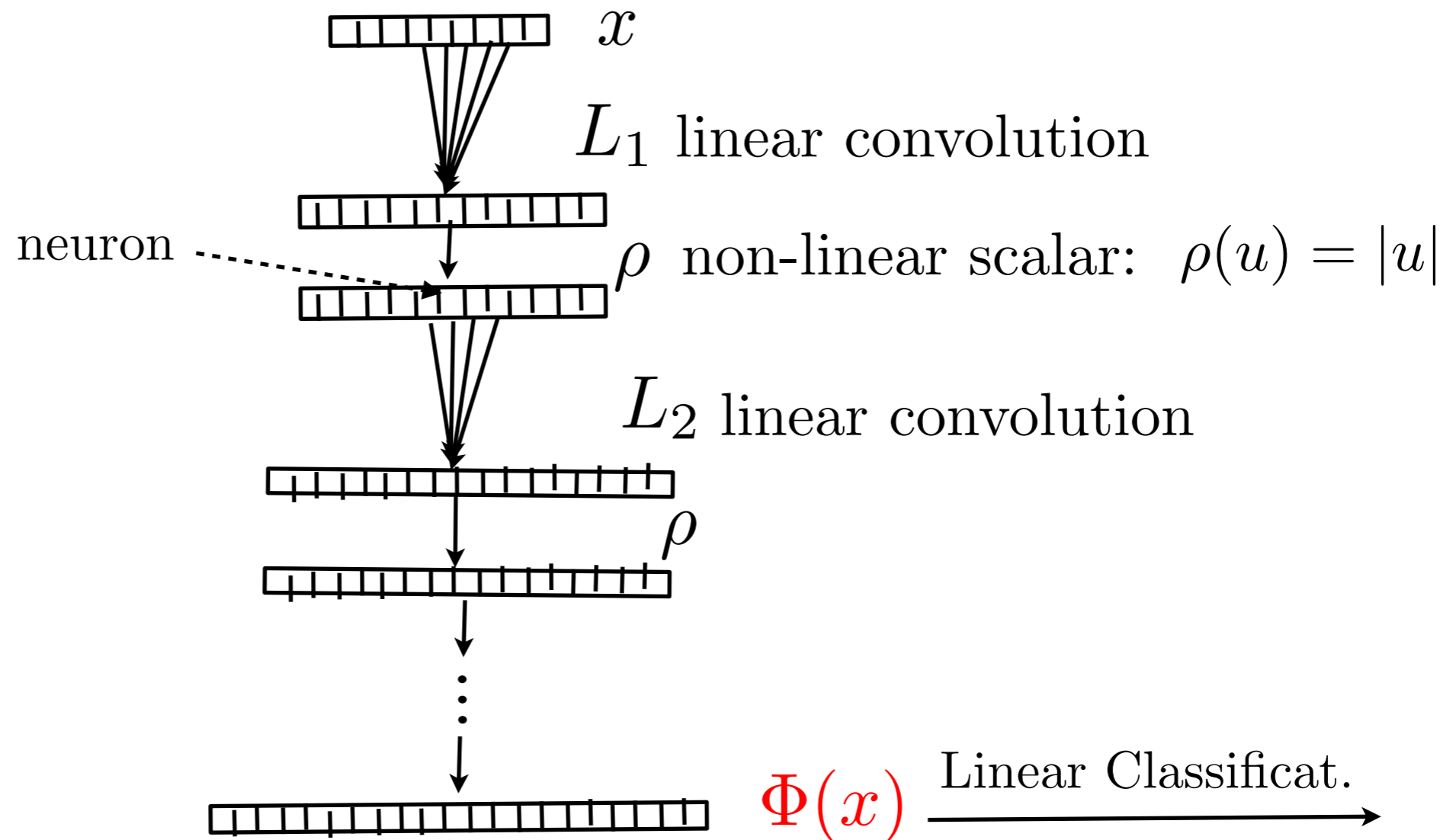
- Selection and vote of many weak features.

How to define  $\Phi$  ?

# Deep Convolution Networks

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

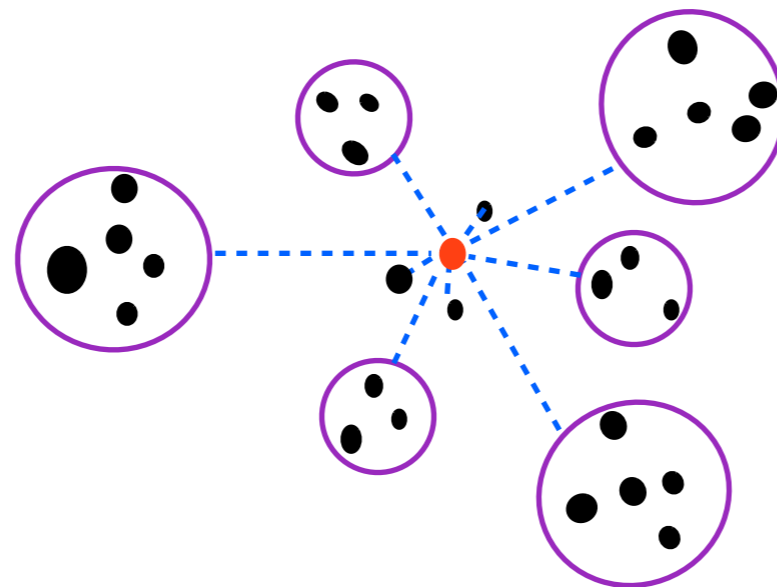
Hierarchical  
Invariants



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

Why does it work so well ?

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

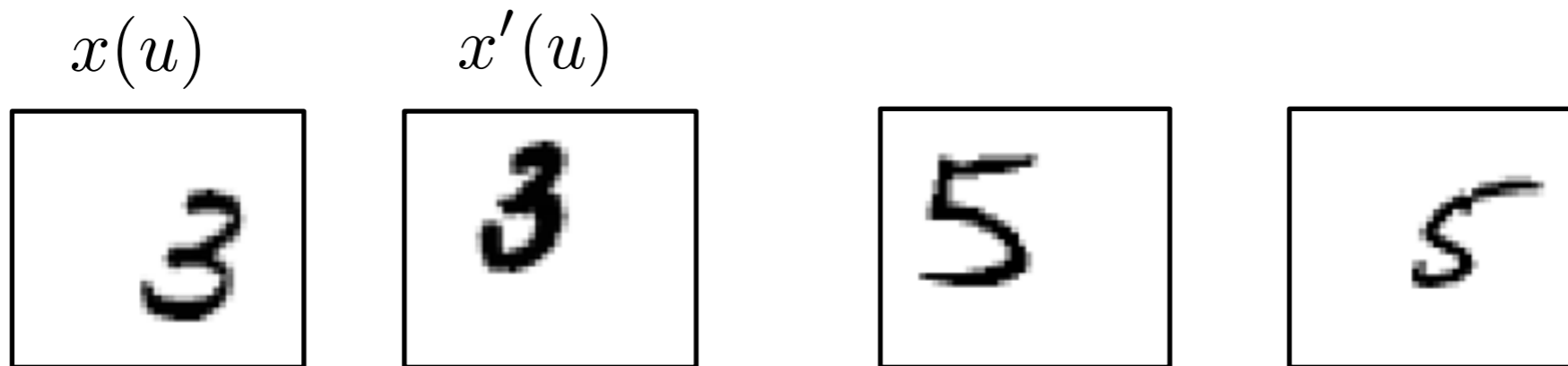


- Multiscale analysis: in mathematics and physics.

# Learning Metric Transformations

- Patterns are too diverse to memorize them all
- We learn interactions and transformations (forces in physics)

Geometric shapes: deformation metric

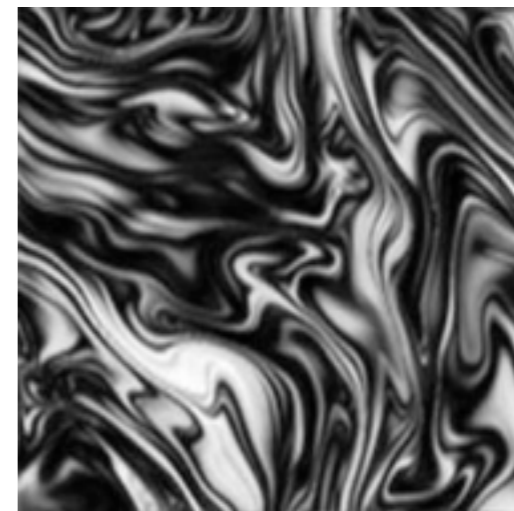
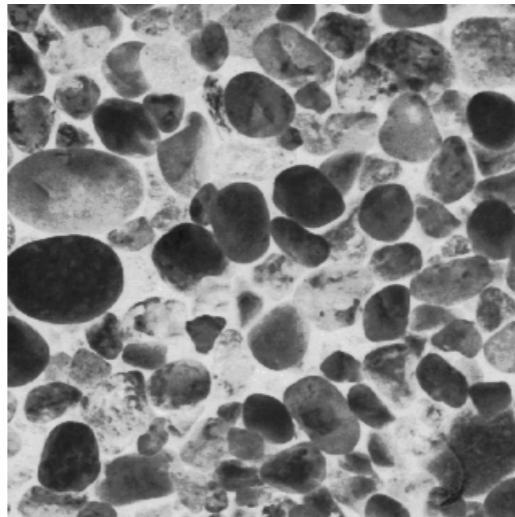
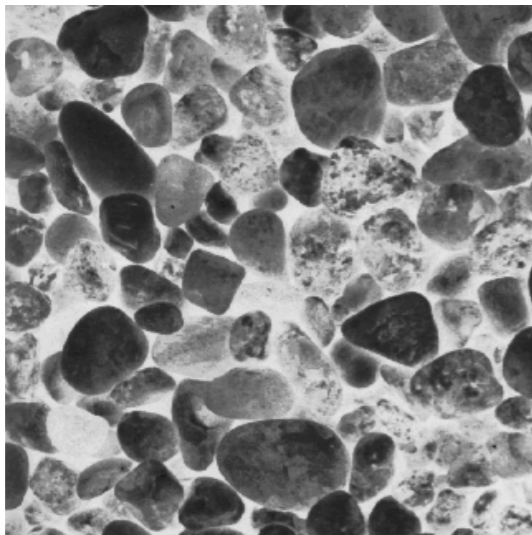


Learn groups of operators

# Learning Metric Transformations

- Patterns are too diverse to memorize them all
- We learn interactions and transformations (forces in physics)

## Stationary Textures



2D Turbulence

Learn probabilistic metrics

# ImageNet Data Basis

- Data basis with 1 million images and 2000 classes

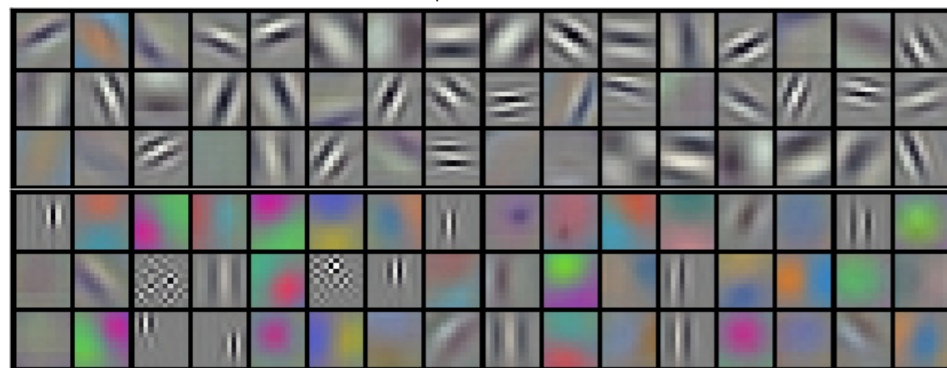
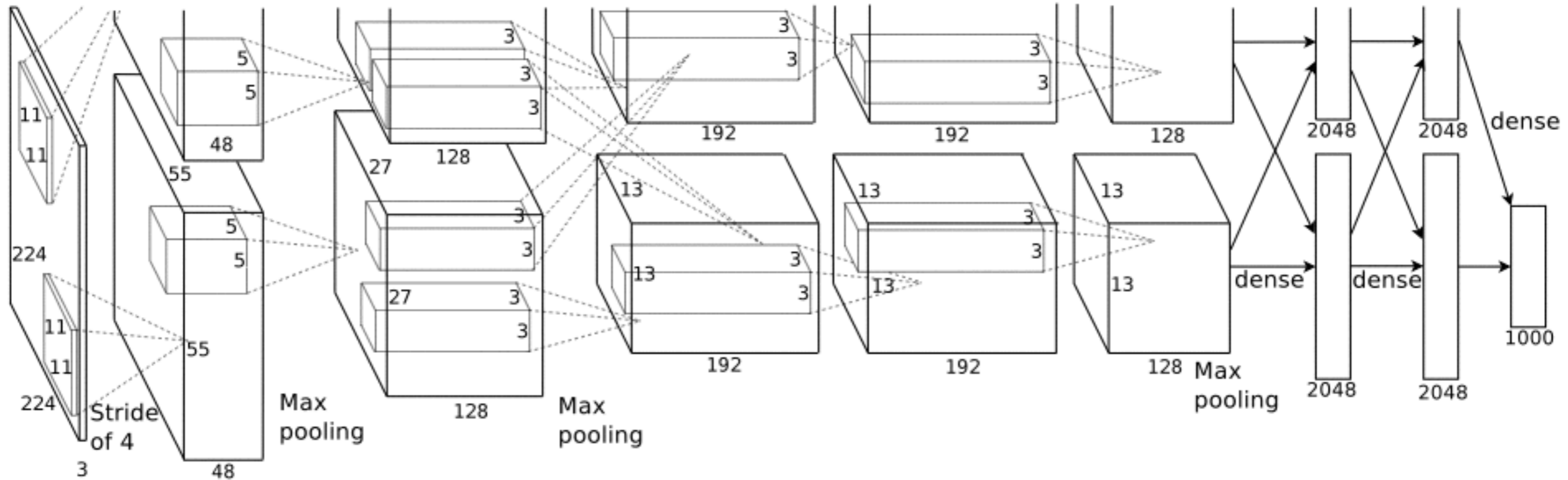




# Alex Deep Convolution Network

*A. Krizhevsky, Sutssever, Hinton*

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



Wavelets

# Image Classification



**mite**

█	mite
█	black widow
█	cockroach
█	tick
█	starfish



**container ship**

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



**motor scooter**

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



**leopard**

█	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

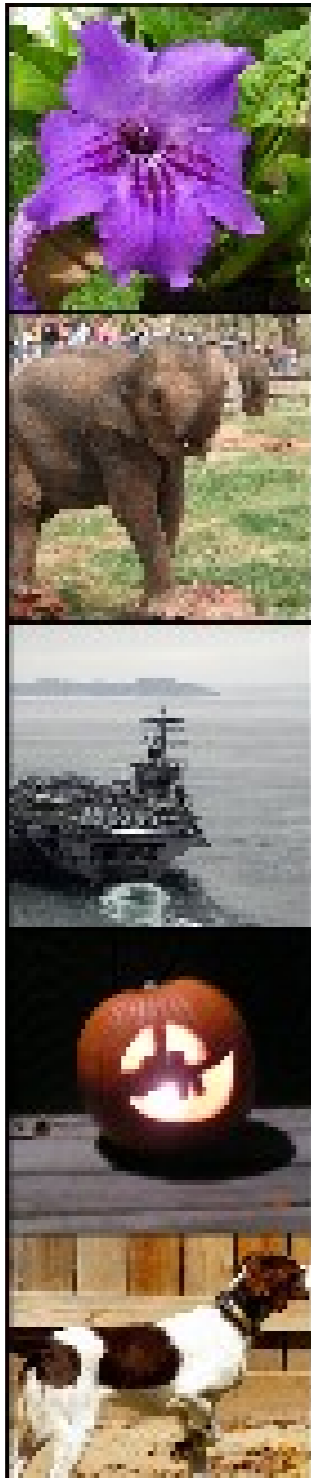


**Madagascar cat**

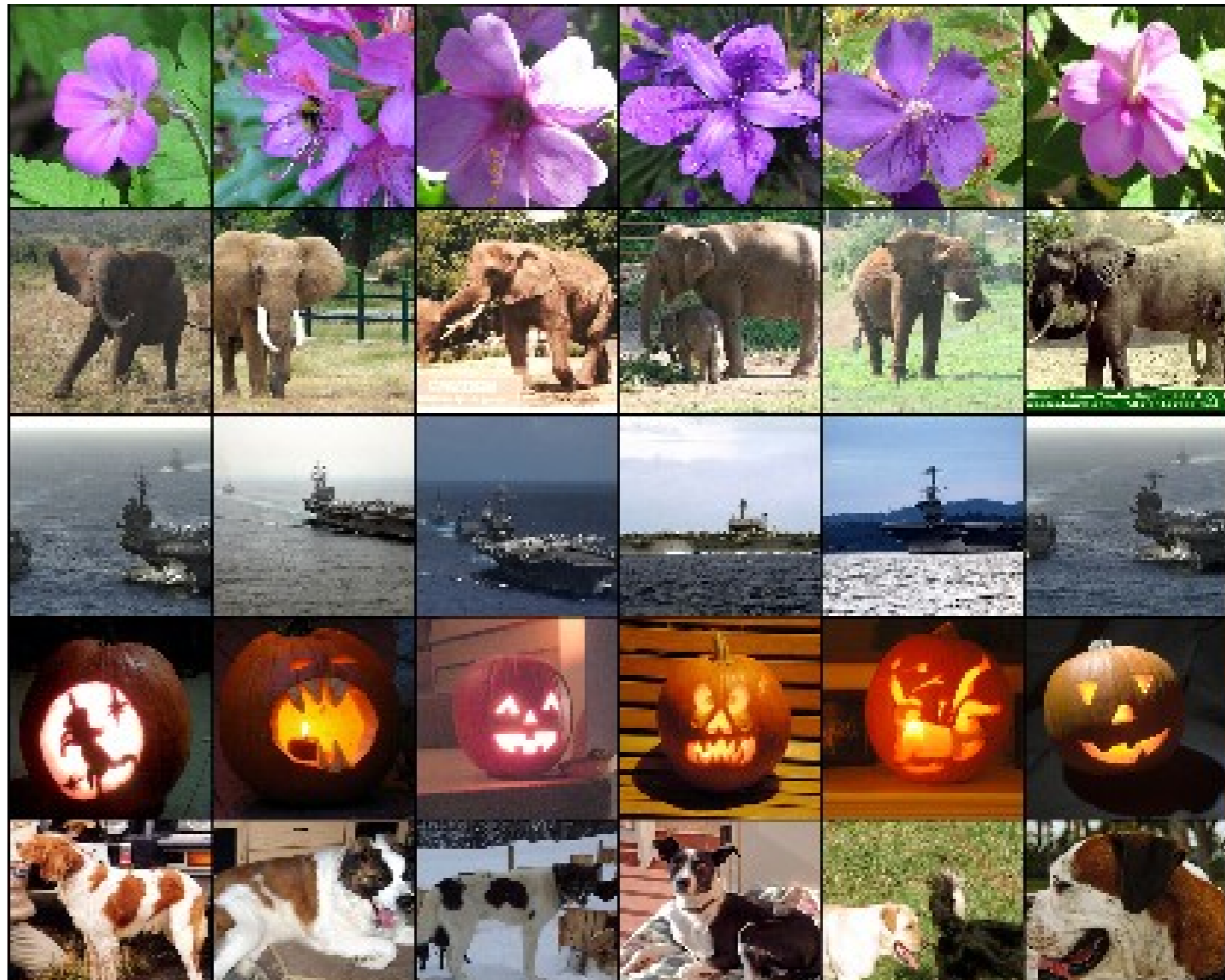
█	squirrel monkey
█	spider monkey
█	titi
█	indri
█	howler monkey

# Image Retrieval

**TEST  
IMAGE**



**RETRIEVED IMAGES**



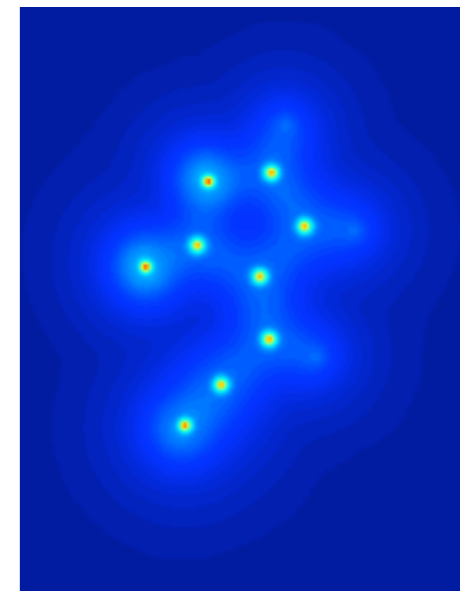
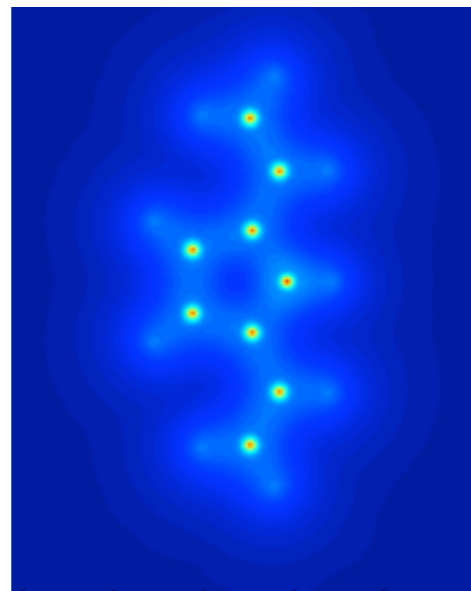
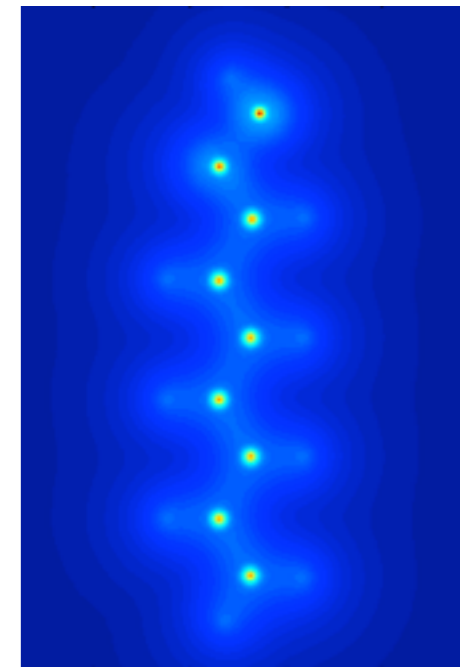
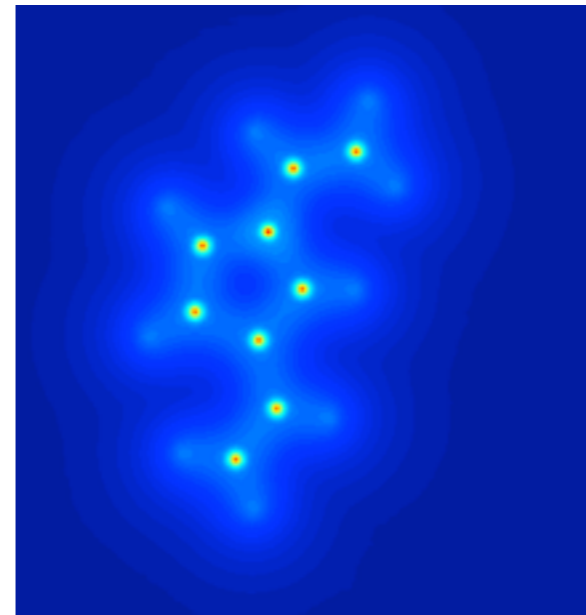
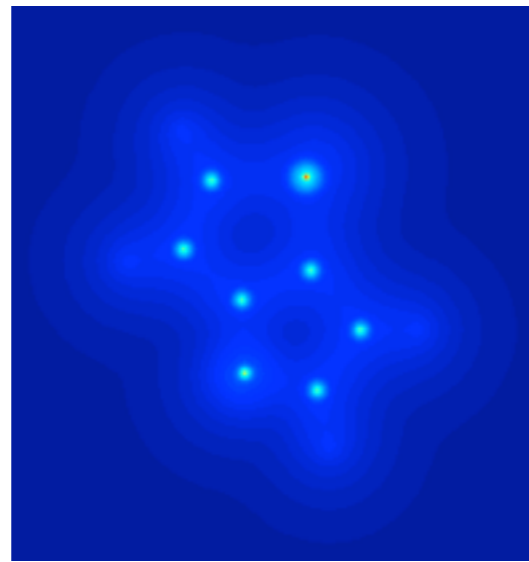
# Scene Labeling



# Quantum Chemistry

Regression of molecule energy from atomic positions  
without solving Schroedinger equation

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



# Paysage Scientifique

Informatique  
Algorithmiciens

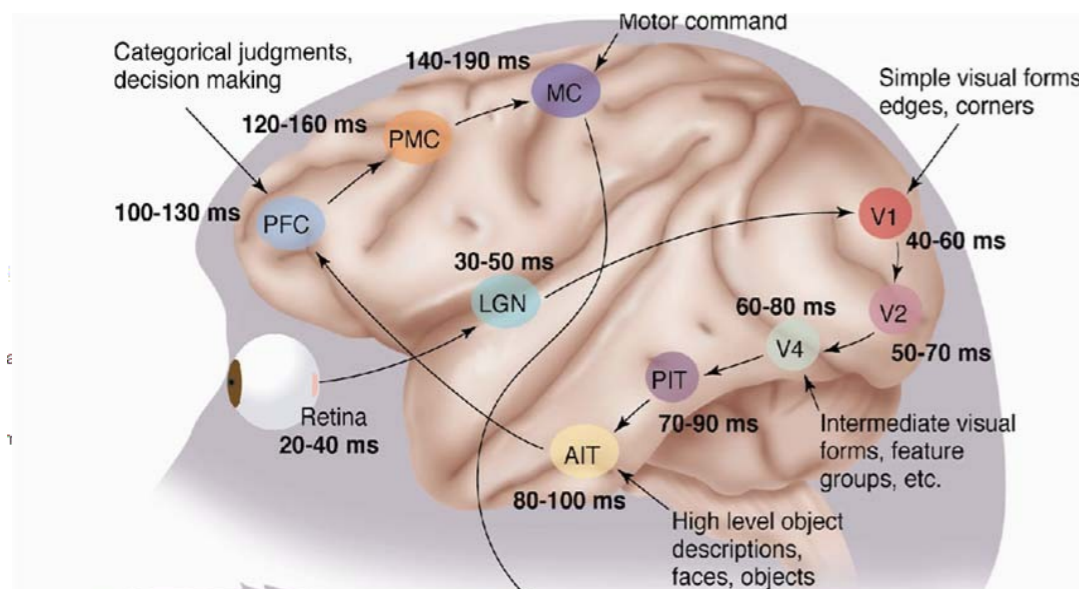
Mathématiques

Physique  
Chimie

Applications  
Médicales  
Industrielles  
Services...

Neurosciences

Sciences Sociales  
Linguistique  
Philosophie



- L'ensemble des sciences vont être profondément impactés

# Paysage Industriel

**Utilisateurs:**

Sociétés

Individus

**Réseaux**

**Terminaux:**

Ordinateurs

Téléphones

Objets connectés

Apple, Samsung

Huawei, Google

**Internet**

Google, FaceBook

Amazon, Microsoft

**Applications/Services**

Start-up, IBM,...

- A terme, l'intelligence ira dans les terminaux
- Pas d'acteur Français important: prise de conscience lente