# Can Robots Learn to See?

Yann LeCun

**Silver Professor of Computer Science and Neural Science**

**The Courant Institute of Mathematical Sciences**

**And Center for Neural Science**

**New York University**

New York University

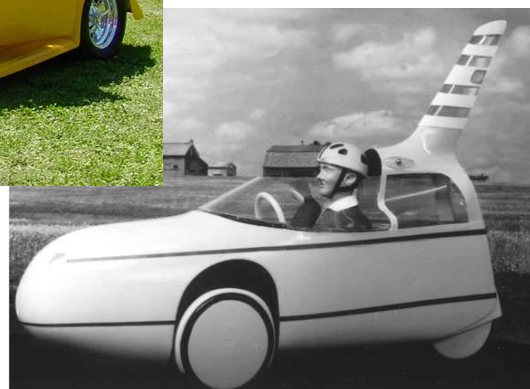# The Next Challenge for AI, Robotics, and Neuroscience

- **How do we learn vision and perception?**
  - From the image of an airplane, how do we extract a representation that is invariant to pose, illumination, background, clutter, object instance....
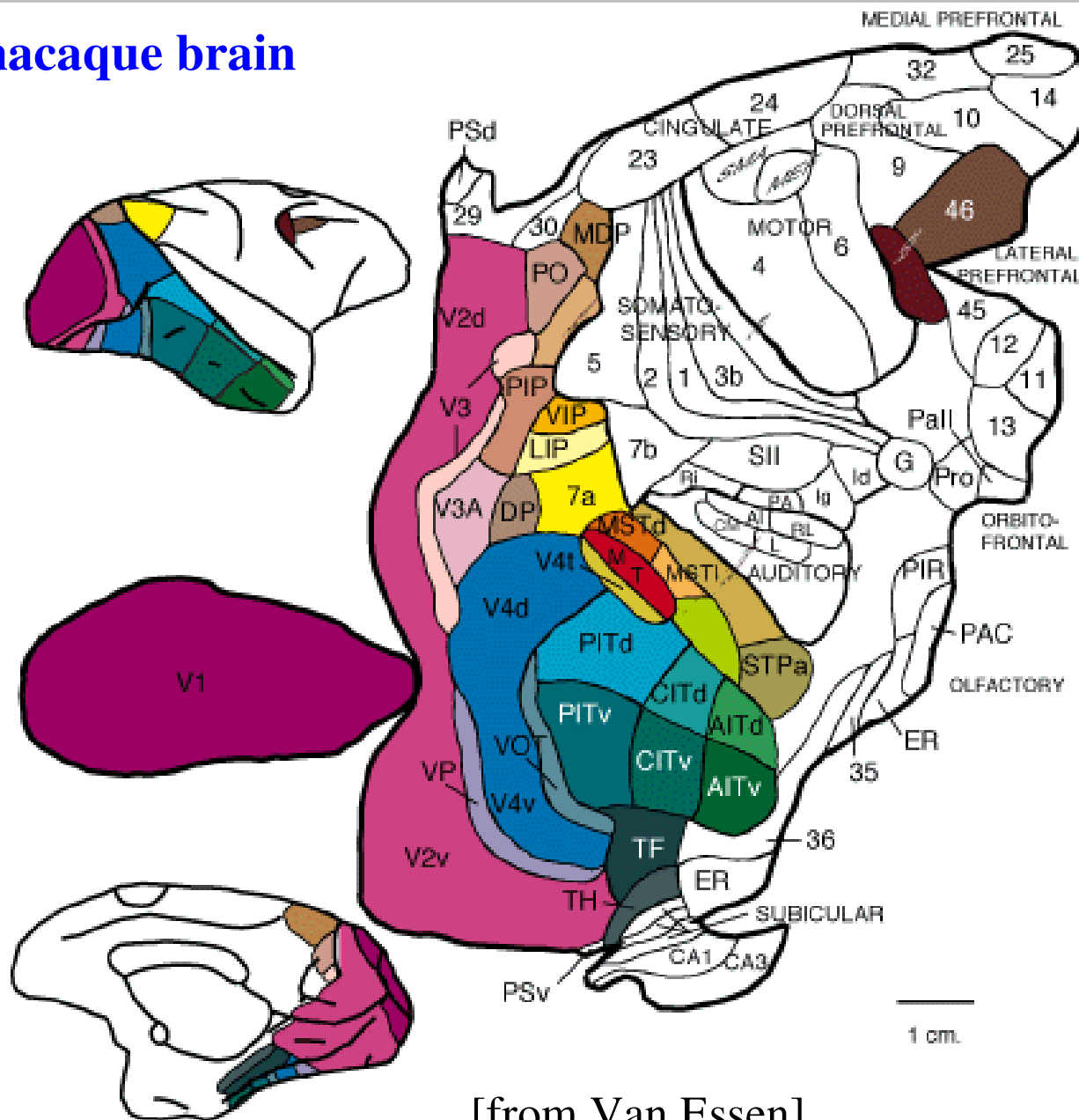  - How can a human (or a machine) learn those representations by just looking at the world?

- **How can we learn visual categories from just a few examples?**
  - I don't need to see many airplanes before I can recognize every airplane (even really weird ones)
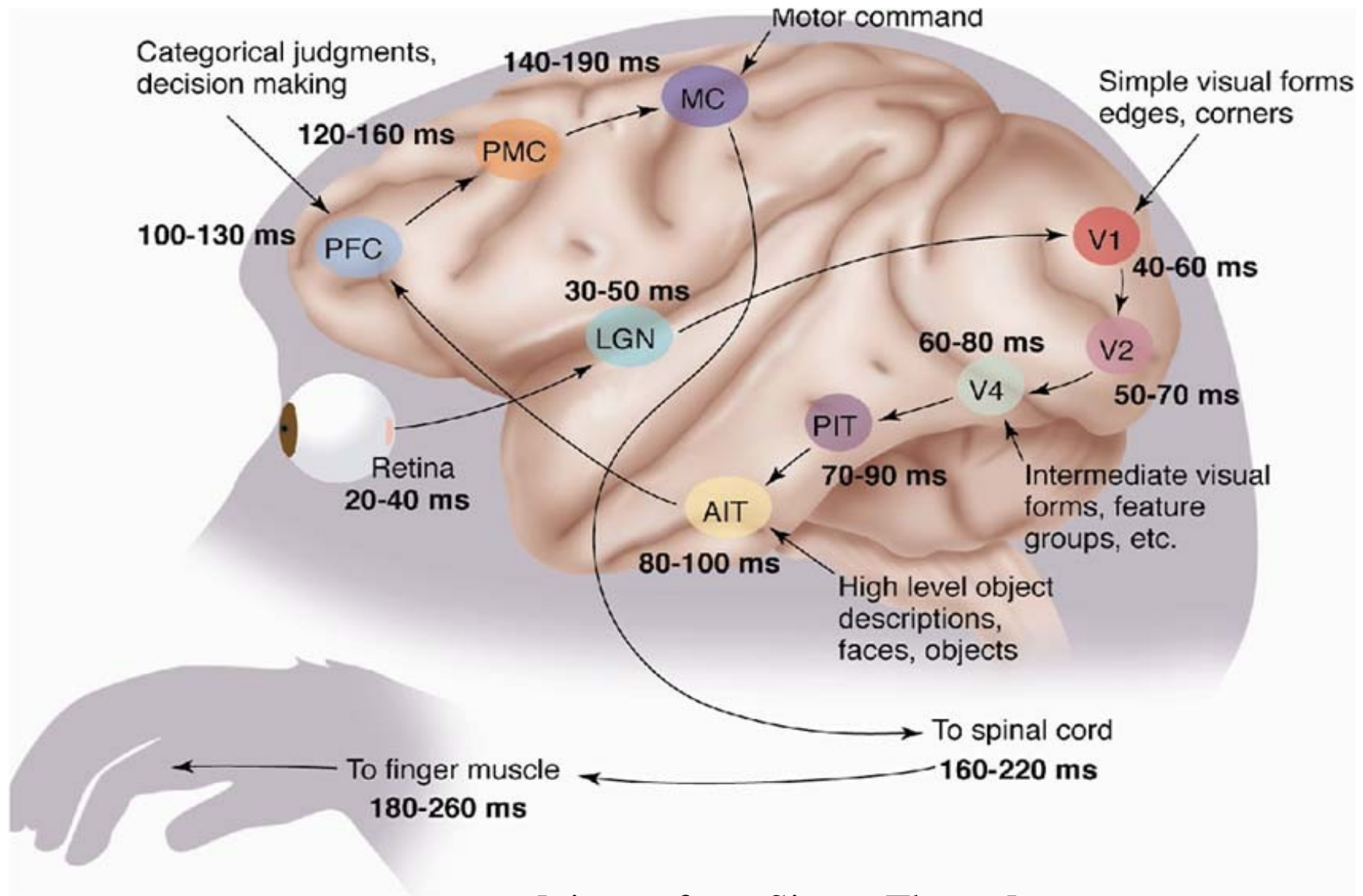
*Yann LeCun*

# Vision occupies a big chunk of our brains

- **1/3 of the macaque brain**



[from Van Essen]

**The ventral (recognition) pathway in the visual cortex**
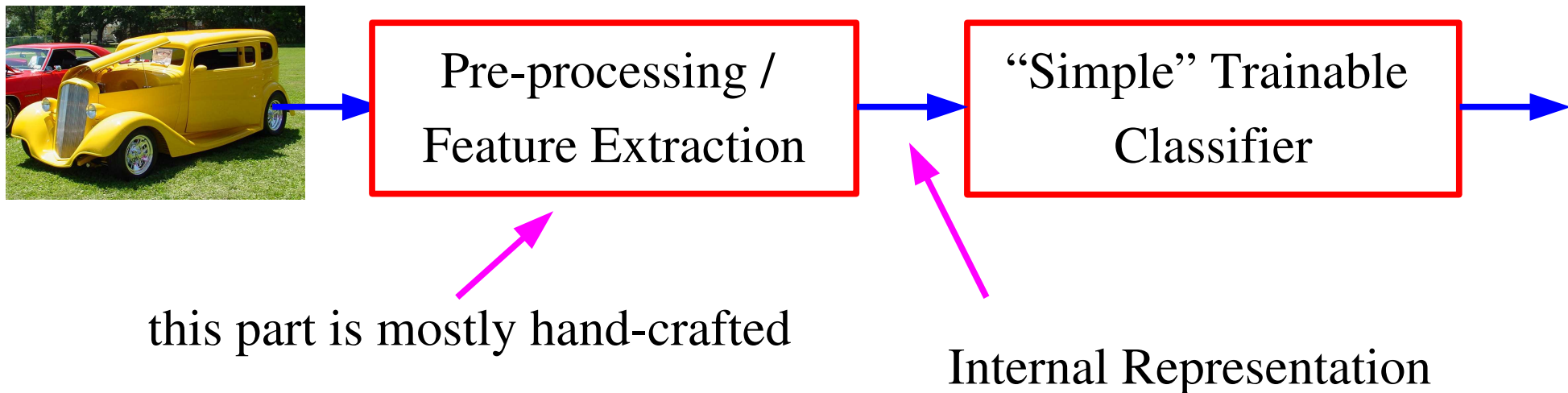


[picture from Simon Thorpe]

# The Primate's Visual System is Deep (LGN->V1->V2->V4->IT)

- **The recognition of everyday objects is a very fast process.**
  - The recognition of common objects is essentially "feed forward."
  - But not all of vision is feed forward.

- **Much of the visual system (all of it?) is the result of learning**
  - How much prior structure is there?

- **If the visual system is deep (around 10 layers) and learned**

- **what is the learning algorithm of the visual cortex?**
  - What learning algorithm can train neural nets as "deep" as the visual system (10 layers?).
  - Unsupervised vs Supervised learning
  - What is the loss function?
  - What is the organizing principle?
  - Broader question (Hinton): what is the learning algorithm of the neo-cortex?

*Yann LeCun*
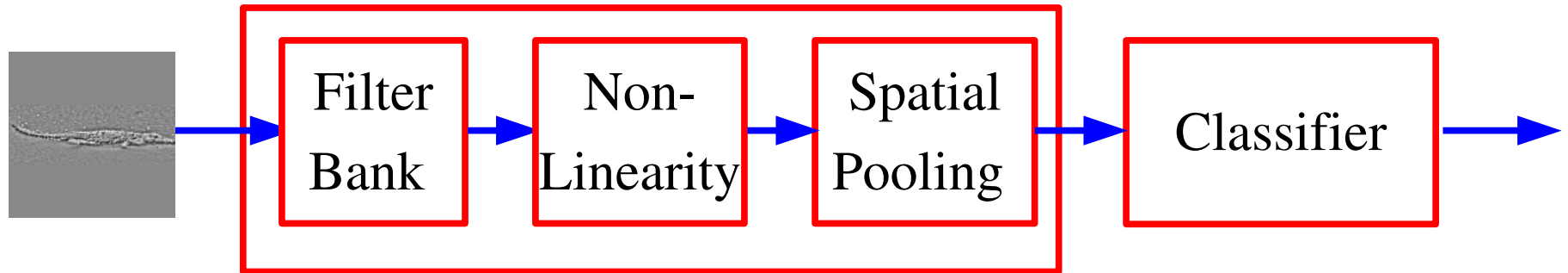
## The Broader Challenge of Machine Learning and AI

- **Can we devise learning algorithms to train a "deep" artificial visual system, and other artificial perception systems.**

- **How can we learn the structure of the world?**
  - How can we build/learn internal representations of the world that allow us to discover its hidden structure?
  - How can we learn internal representations that capture the relevant information and eliminates irrelevant variabilities?

- **How can a human or a machine learn internal representations by just looking at the world?**

- **Can we find learning methods that solve really complex problems end-to-end, such as vision, natural language, speech....?**

# The Traditional "Shallow" Architecture for Recognition

Pre-processing /
Feature Extraction

"Simple" Trainable
Classifier

this part is mostly hand-crafted

Internal Representation

- The raw input is pre-processed through a hand-crafted feature extractor
- **The features are not learned**
- The trainable classifier is often generic (task independent), and "simple" (linear classifier, kernel machine, nearest neighbor,.....)
- The most common Machine Learning architecture: the Kernel Machine

*Yann LeCun*

New York University

# "Modern" Object Recognition Architecture in Computer Vision



**Filter Bank** → **Non-Linearity** → **Spatial Pooling** → **Classifier**

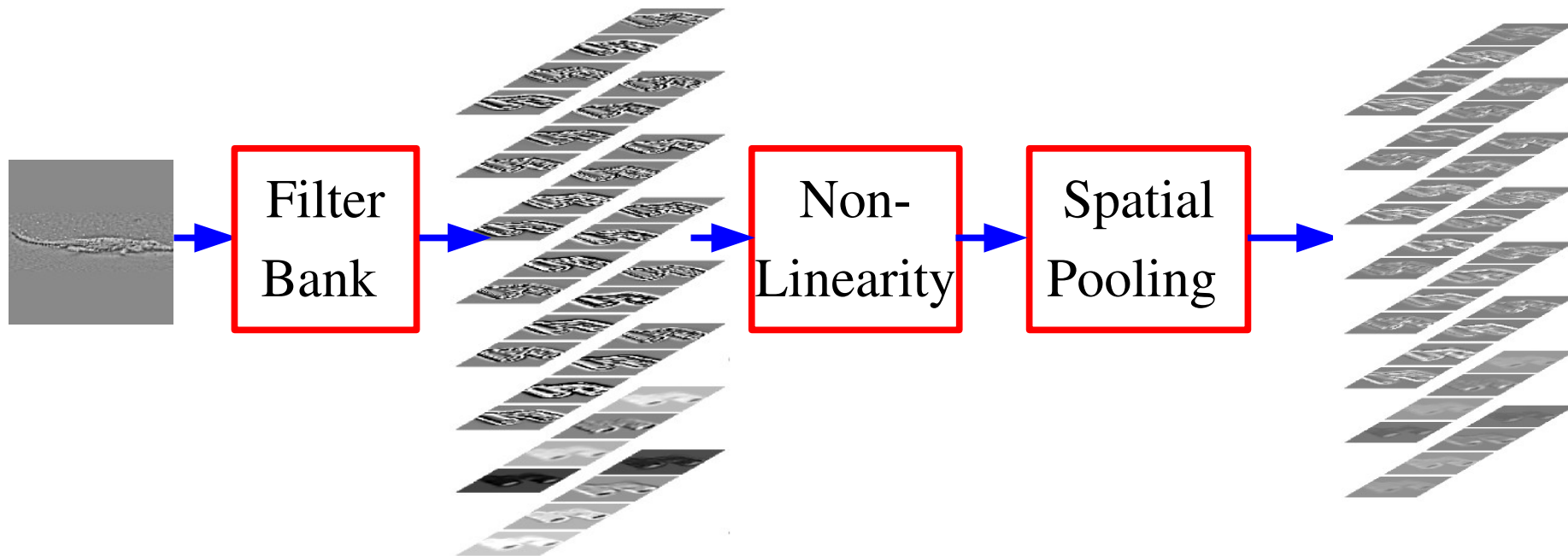| Oriented Edges | Sigmoid | Averaging |
|---|---|---|
| Gabor Wavelets | Rectification | Max pooling |
| Other Filters... | Vector Quant. | VQ+Histogram |
| | Contrast Norm. | Geometric Blurr |

🔵 **Example:**
- ▶ Edges + Rectification + Histograms + SVM [Dalal & Triggs 2005]
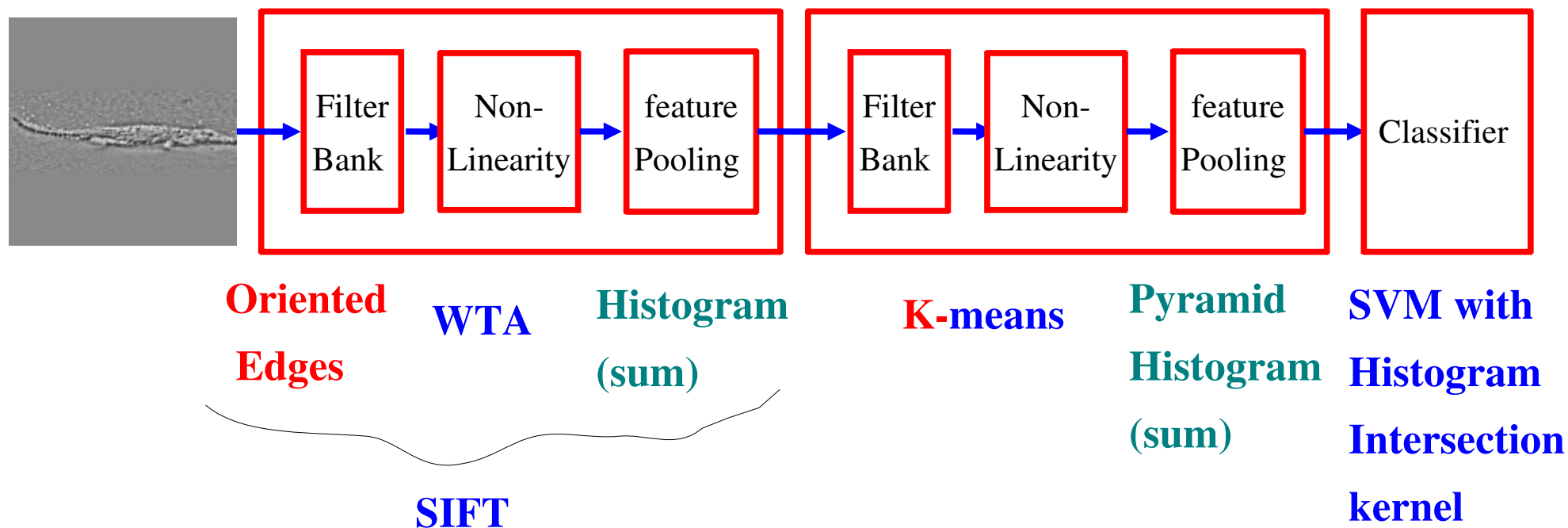- ▶ SIFT + classification

🔵 **Fixed Features + "shallow" classifier**

# Feature Extraction by Filtering and Pooling



🔵 **Biologically-inspired models of low-level feature extraction**

▶ Inspired by [Hubel and Wiesel 1962]

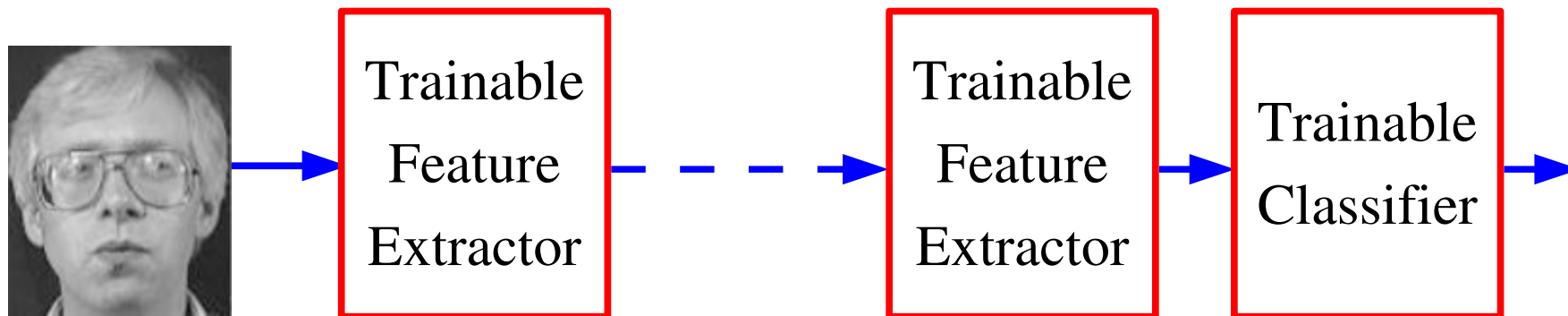# "State of the Art" architecture for object recognition



**Oriented Edges** → **WTA** → **Histogram (sum)** | **K-means** → **Pyramid Histogram (sum)** → **SVM with Histogram Intersection kernel**

SIFT

- 🔵 **Example:**
  - ▶ SIFT features with Spatial Pyramid Match Kernel SVM [Lazebnik et al. 2006]
- 🔵 **Fixed Features + unsupervised features + "shallow" classifier**

New York University

# Good Representations are Hierarchical



**In Language: hierarchy in syntax and semantics**

▶ Words->Parts of Speech->Sentences->Text
▶ Objects,Actions,Attributes...-> Phrases -> Statements -> Stories

**In Vision: part-whole hierarchy**

▶ Pixels->Edges->Textons->Parts->Objects->Scenes

# "Deep" Learning: Learning Hierarchical Representations



Learned Internal Representation

- **Deep Learning**: learning a hierarchy of internal representations

- From low-level features to mid-level invariant representations, to object identities

- Representations are increasingly invariant as we go up the layers

- **using multiple stages gets around the specificity/invariance dilemma**

New York University

# Do we really need deep architectures?

- **We can approximate any function as close as we want with shallow architecture (e.g. a kernel machine). Why would we need deep ones?**

$$y = \sum_{i=1}^{P} \alpha_i K(X, X^i) \qquad y = F(W^1 . F(W^0 . X))$$

  ▶ kernel machines and 2-layer neural net are "universal".

- **Deep learning machines**

$$y = F(W^K . F(W^{K-1} . F(.... F(W^0 . X)...)))$$

- **Deep machines are more efficient for representing certain classes of functions, particularly those involved in visual recognition**

  ▶ they can represent more complex functions with less "hardware"

- **We need an efficient parameterization of the class of functions that are useful for "AI" tasks.**

# Why are Deep Architectures More Efficient?

[Bengio & LeCun 2007 "Scaling Learning Algorithms Towards AI"]

- **A deep architecture trades space for time (or breadth for depth)**
  - more layers (more sequential computation),
  - but less hardware (less parallel computation).
  - Depth-Breadth tradoff

- **Example1: N-bit parity**
  - requires N-1 XOR gates in a tree of depth log(N).
  - requires an exponential number of gates of we restrict ourselves to 2 layers (DNF formula with exponential number of minterms).

- **Example2: circuit for addition of 2 N-bit binary numbers**
  - Requires O(N) gates, and O(N) layers using N one-bit adders with ripple carry propagation.
  - Requires lots of gates (some polynomial in N) if we restrict ourselves to two layers (e.g. Disjunctive Normal Form).
  - Bad news: almost all boolean functions have a DNF formula with an exponential number of minterms O(2^N).....

# Deep Supervised Learning is Hard

- **The loss surface is non-convex, ill-conditioned, has saddle points, has flat spots.....**

- **For large networks, it will be horrible! (not really, actually)**

- **Back-prop doesn't work well with networks that are tall and skinny.**
  - Lots of layers with few hidden units.

- **Back-prop works fine with short and fat networks**
  - But over-parameterization becomes a problem without regularization
  - Short and fat nets with fixed first layers aren't very different from SVMs.

- **For reasons that are not well understood theoretically, back-prop works well when they are highly structured**
  - e.g. convolutional networks.

# Can't we train multi-stage vision architectures?



- Stacking multiple stages of feature extraction/pooling.

- Creates a hierarchy of features

New York University

# Convolutional Network



- **Hierarchical/multilayer:** features get progressively more global, invariant, and numerous

- **dense features:** features detectors applied everywhere (no interest point)

- **broadly tuned (possibly invariant) features:** sigmoid units are on half the time.

- **Global discriminative training:** The whole system is trained "end-to-end" with a gradient-based method to minimize a global loss function

- **Integrates segmentation, feature extraction, and invariant classification in one fell swoop.**

New York University

# An Old Idea for Local Shift Invariance

- **[Hubel & Wiesel 1962]:**
  - simple cells detect local features
  - complex cells "pool" the outputs of simple cells within a retinotopic neighborhood.



"Simple cells"

"Complex cells"

pooling subsampling

Multiple convolutions

Retinotopic Feature Maps

# The Multistage Hubel-Wiesel Architecture

**Building a complete artificial vision system:**

- Stack multiple stages of simple cells / complex cells layers
- Higher stages compute more global, more invariant features
- Stick a classification layer on top
- [Fukushima 1971-1982]
  - neocognitron
- [LeCun 1988-2007]
  - convolutional net
- [Poggio 2002-2006]
  - HMAX
- [Ullman 2002-2006]
  - fragment hierarchy
- [Lowe 2006]
  - HMAX

**QUESTION: How do we find (or learn) the filters?**

# Convolutional Net Architecture



input
1@32x32

Layer 1
6@28x28

Layer 2
6@14x14

Layer 3
12@10x10

Layer 4
12@5x5

Layer 5
100@1x1

Layer 6: 10

10

5x5
convolution

2x2
pooling/
subsampling

5x5
convolution

2x2
pooling/
subsampling

5x5
convolution

- **Convolutional net for handwriting recognition** (400,000 synapses)
- Convolutional layers (simple cells): all units in a feature plane share the same weights
- Pooling/subsampling layers (complex cells): for invariance to small distortions.
- **Supervised gradient-descent learning using back-propagation**
- **The entire network is trained end-to-end. All the layers are trained simultaneously.**

# Face Detection and Pose Estimation with Convolutional Nets

- **Training:** 52,850, 32x32 grey-level images of faces, 52,850 non-faces.

- **Each sample:** used 5 times with random variation in scale, in-plane rotation, brightness and contrast.

- **2nd phase:** half of the initial negative set was replaced by false positives of the initial version of the detector .

# Face Detection: Results

| Data Set-> | TILTED | | PROFILE | | MIT+CMU | |
|---|---|---|---|---|---|---|
| False positives per image-> | 4.42 | 26.9 | 0.47 | 3.36 | 0.5 | 1.28 |
| **Our Detector** | 90% | 97% | 67% | 83% | 83% | 88% |
| **Jones & Viola (tilted)** | 90% | 95% | x | | x | |
| **Jones & Viola (profile)** | x | | 70% | 83% | x | |

University

*Yann LeCun*

# Face Detection with a ConvNet



- Demo produced with EBLearn open source package

- http://eblearn.sf.net

New York University

# Generic Object Detection and Recognition with Invariance to Pose and Illumination

- **50** toys belonging to 5 categories: **animal, human figure, airplane, truck, car**
- **10** instance per category: 5 instances used for training, 5 instances for testing
- **Raw dataset:** **972** stereo pair of each object instance. **48,600** image pairs total.

- **For each instance:**
- **18 azimuths**
  - 0 to 350 degrees every 20 degrees
- **9 elevations**
  - 30 to 70 degrees from horizontal every 5 degrees
- **6 illuminations**
  - on/off combinations of 4 lights
- **2 cameras (stereo)**
  - 7.5 cm apart
  - 40 cm from the object



**Training instances**    **Test instances**

*Yann LeCun*

# Convolutional Network



- **90,857 free parameters, 3,901,162 connections.**
- The architecture alternates convolutional layers (feature detectors) and subsampling layers (local feature pooling for invariance to small distortions).
- **The entire network is trained end-to-end** (all the layers are trained simultaneously).
- A gradient-based algorithm is used to minimize a supervised loss function.

*Yann LeCun*

New York University

# Normalized-Uniform Set: Error Rates

- **Linear Classifier on raw stereo images:** **30.2% error.**
- **K-Nearest-Neighbors on raw stereo images:** **18.4% error.**
- **K-Nearest-Neighbors on PCA-95:** **16.6% error.**
- **Pairwise SVM on 96x96 stereo images:** **11.6% error**
- **Pairwise SVM on 95 Principal Components:** **13.3% error.**
- **Convolutional Net on 96x96 stereo images:** **5.8% error.**



**Training instances**   **Test instances**

*Yann LeCun*

# Jittered-Cluttered Dataset



- **Jittered-Cluttered Dataset:**

- **291,600** tereo pairs for training, **58,320** for testing

- Objects are jittered: position, scale, in-plane rotation, contrast, brightness, backgrounds, distractor objects,...

- Input dimension: 98x98x2 (approx 18,000)

# Experiment 2: Jittered-Cluttered Dataset



- **291,600** training samples, **58,320** test samples

- SVM with Gaussian kernel                                    **43.3% error**

- Convolutional Net with **binocular** input:                 **7.8% error**

- **Convolutional Net + SVM on top:**                         **5.9% error**

- Convolutional Net with **monocular** input:                 **20.8% error**

- Smaller **mono** net (DEMO):                                **26.0% error**

- **Dataset available from http://www.cs.nyu.edu/~yann**

# Examples (Monocular Mode)

# Examples (Monocular Mode)

# Visual Navigation for a Mobile Robot

[LeCun et al. NIPS 2005]

- Mobile robot with two cameras
- The convolutional net is trained to emulate a human driver from recorded sequences of video + human-provided steering angles.
- The network maps stereo images to steering angles for obstacle avoidance

# Convolutional Nets For Brain Imaging and Biology

**Brain tissue reconstruction from slice images [Jain,....,Denk, Seung 2007]**

- Sebastian Seung's lab at MIT.
- 3D convolutional net for image segmentation
- ConvNets Outperform MRF, Conditional Random Fields, Mean Shift, Diffusion,...[ICCV'07]

# Industrial Applications of ConvNets

- **AT&T/Lucent/NCR**
  - Check reading, OCR, handwriting recognition (deployed 1996)

- **Vidient Inc**
  - Vidient Inc's "SmartCatch" system deployed in several airports and facilities around the US for detecting intrusions, tailgating, and abandoned objects (Vidient is a spin-off of NEC)

- **NEC Labs**
  - Cancer cell detection, automotive applications, kiosks

- **Google**
  - OCR, face and license plate removal from StreetView

- **Microsoft**
  - OCR, handwriting recognition, speech detection

- **France Telecom**
  - Face detection, HCI, cell phone-based applications

- **Other projects: HRL (3D vision)....**

# FPGA Custom Board: NYU ConvNet Processor

🔵 **Xilinx Virtex 4 FPGA, 8x5 cm board**       [Farabet et al. 2009]

  ▶ Dual camera port, Fast dual QDR RAM,

🔵 **New version being developed with Eugenio Culurciello (Yale EE)**

  ▶ Full custom chip
  ▶ Version for Virtex 6 FPGA

New York University

**Reconfigurable Dataflow Architecture**

[Farabet et al. 2010]

# FPGA Performance

**Seconds per frame for a robot vision task (log scale)** [Farabet et al. 2010]

# Problem: supervised ConvNets don't work with few labeled samples

- **On recognition tasks with few labeled samples, deep supervised architectures don't do so well**

- **Example: Caltech-101 Object Recognition Dataset**
  - 101 categories of objects (gathered from the web)
  - Only 30 training samples per category!

- **Recognition rates (OUCH!):**
  - Supervised ConvNet:                                    **29.0%**
  - SIFT features + Pyramid Match Kernel SVM:    **64.6%**
    - [Lazebnik et al. 2006]

- **When learning the features, there are simply too many parameters to learn in purely supervised mode (or so we thought).**

*face*

*beaver*

*wild cat*

*lotus*

*ant*

*dollar*

*w. chair*

*minaret*

*cellphone*

*joshua t.*

*cougar body*

*background*

*metronome*

*Yann LeCun*

New York University

# Unsupervised Deep Learning: Leveraging Unlabeled Data

[Hinton 05, Bengio 06, LeCun 06, Ng 07]

- **Unlabeled data is usually available in large quantity**

- **A lot can be learned about the world by just looking at it**

- **Unsupervised learning captures underlying regularities about the data**

- **The best way to capture underlying regularities is to learn good representations of the data**

- **The main idea of Unsupervised Deep Learning**
  - ▶ Learn each layer one at a time in unsupervised mode
  - ▶ Stick a supervised classifier on top
  - ▶ Optionally: refine the entire system in supervised mode

- **Unsupervised Learning view as Energy-Based Learning**

# Unsupervised Feature Learning with Sparse Coding

[Olshausen & field 1997]

- **Find a dictionary of basis functions such that any input can be reconstructed of a sparse linear combination of them.**



- **Energy:**
$$E(Y^i, Z; W_d) = \|Y^i - W_d Z\|^2 + \lambda \sum_j |z_j|$$

- **Optimal Code**
$$Z^i = argmin_z E(Y^i, z; W_d)$$

- **Free Energy:**
$$F(Y^i; W_d) = F(Z^i) = min_z E(Y^i, z; W_d)$$

*Yann LeCun*

# Unsupervised Feature Learning with Sparse Coding

🔹 **The learning algorithm minimizes the loss function:**

$$L(W_d) = \sum_i F(Y^i; W_d) = \sum_i \left( min_Z \, E(Y^i, Z; W_d) \right)$$

🔹 **The columns of Wd are normalized**



🔹 **Energy:**
$$E(Y^i, Z; W_d) = \|Y^i - W_d Z\|^2 + \lambda \sum_j |z_j|$$

🔹 **Free Energy:**
$$F(Y^i; W_d) = F(Z^i) = min_z \, E(Y^i, z; W_d)$$

**Inference: find Z that minimizes the energy for a given Y**

$$E(Y^i, Z^i; W_d) = \|Y^i - W_d Z^i\|^2 + \lambda \sum_j |z^i_j|$$

$$Z^i = argmin_z E(Y^i, z; W_d)$$

▶ **For each new Y, an optimization algorithm must be run to find the corresponding optimal Z**

▶ **This would be very slow for large scale vision tasks**

▶ **Also, the optimal Z are very unstable:**

- A small change in Y can cause a large change in the optimal Z

# Solution: Predictive Sparse Decomposition (PSD)

[Kavukcuoglu, Ranzato, LeCun, 2009]

- **Prediction the optimal code with a trained encoder**

- **Energy = reconstruction_error + code_prediction_error + code_sparsity**

$$E(Y^i, Z) = \|Y^i - W_d Z\|^2 + \|Z - g_e(W_e, Y^i)\|^2 + \lambda \sum_j |z_j|$$

$$g_e(W_e, Y^i) = D \tanh(W_e Y)$$



**INPUT** Y → $\|Y^i - \tilde{Y}\|^2$ ← $W_d Z$ ← Z → $|z_j|$ → $\lambda \sum_j .$ **FEATURES**

$g_e(W_e, Y^i)$ → $\|Z - \tilde{Z}\|^2$

# PSD: Inference

🔵 **Inference by gradient descent starting from the encoder output**

$$E(Y^i, Z) = \|Y^i - W_d Z\|^2 + \|Z - g_e(W_e, Y^i)\|^2 + \lambda \sum_j |z_j|$$

$$Z^i = argmin_z E(Y^i, z; W)$$

# PSD: Learning [Kavukcuoglu et al. 2009]

- **Learning by minimizing the average energy of the training data with respect to Wd and We.**

- **Loss function:**

$$L(W_d, W_e) = \sum_i F(Y^i; W_d, W_e)$$

$$F(Y^i; W_d, W_e) = min_z E(Y^i, z; W_d, W_e)$$

$$\|Y^i - \tilde{Y}\|^2 \qquad W_d Z \qquad \lambda \sum_j .$$

$$Y \qquad Z \qquad |z_j|$$

INPUT

FEATURES

$$g_e(W_e, Y^i) \qquad \|Z - \tilde{z}\|^2$$

🔵 **1. Initialize Z = Encoder(Y)**

🔵 **2. Find Z that minimizes the energy function**

🔵 **3. Update the Decoder basis functions to reduce reconstruction error**

🔵 **4. Update Encoder parameters to reduce prediction error**

🔵 **Repeat with next training sample**

# Decoder Basis Functions on MNIST

▶ **PSD trained on handwritten digits: decoder filters are "parts" (strokes).**

- Any digit can be reconstructed as a linear combination of a small number of these "parts".

# PSD Training on Natural Image Patches

- **Basis functions are like Gabor filters (like receptive fields in V1 neurons)**

- **256 filters of size 12x12**

- **Trained on natural image patches from the Berkeley dataset**

- **Encoder is linear-tanh-diagonal**



iteration no 0

New York University

# Learned Features on natural patches: V1-like receptive fields

# Learned Features: V1-like receptive fields

- **12x12 filters**

- **1024 filters**

*Yann LeCun*

# Classification Error Rate on MNIST

**Supervised Linear Classifier trained on 200 trained sparse features**

▶ Red: linear-tanh-diagonal encoder; Blue: linear encoder

🔵 **Phase 1: train first layer using PSD**



$$\|Y^i - \tilde{Y}\|^2 \quad W_d Z \quad Z \quad |z_j| \quad \lambda \sum_j \cdot$$

$$Y \quad g_e(W_e, Y^i) \quad \|z - \tilde{z}\|^2$$

**FEATURES**

# Using PSD to Train a Hierarchy of Features

- **Phase 1: train first layer using PSD**

- **Phase 2: use encoder + absolute value as feature extractor**

$$Y \quad \rightarrow \quad g_e(W_e, Y^i) \quad \rightarrow \quad |z_j| \quad \rightarrow$$

**FEATURES**

# Using PSD to Train a Hierarchy of Features

- **Phase 1: train first layer using PSD**

- **Phase 2: use encoder + absolute value as feature extractor**

- **Phase 3: train the second layer using PSD**



**FEATURES**

# Using PSD to Train a Hierarchy of Features

- **Phase 1: train first layer using PSD**

- **Phase 2: use encoder + absolute value as feature extractor**

- **Phase 3: train the second layer using PSD**

- **Phase 4: use encoder + absolute value as 2ⁿᵈ feature extractor**

$$Y \longrightarrow g_e(W_e, Y^i) \longrightarrow |z_j| \longrightarrow g_e(W_e, Y^i) \longrightarrow |z_j| \longrightarrow$$

**FEATURES**

*Yann LeCun*

# Using PSD to Train a Hierarchy of Features

- **Phase 1: train first layer using PSD**

- **Phase 2: use encoder + absolute value as feature extractor**

- **Phase 3: train the second layer using PSD**

- **Phase 4: use encoder + absolute value as 2nd feature extractor**

- **Phase 5: train a supervised classifier on top**

- **Phase 6 (optional): train the entire system with supervised back-propagation**



$$Y \rightarrow g_e(W_e, Y^i) \rightarrow |z_j| \rightarrow g_e(W_e, Y^i) \rightarrow |z_j| \rightarrow \text{classifier}$$

**FEATURES**

# "Deep Learning"

[Hinton 05, Bengio 06, LeCun 06, Ng 07]

- **The "deep learning" method was popularized by Hinton for training "deep belief networks".**
  - DBN use a special kind of encoder-decoder architecture called Restricted Boltzmann Machines (RBM)

- **1. Train each layer in an unsupervised fashion, layer by layer**

- **2. Stick a supervised classifier on top, and refine the entire system with gradient descent (back-prop) on a supervised criterion.**

# Unsupervised Learning: Capturing Dependencies Between Variables

🔵 **Energy function: viewed as a negative log probability density**

🔵 **Probabilistic View:**

▶ Produce a probability density function that:

▶ has high value in regions of high sample density

▶ has low value everywhere else (integral = 1).

🔵 **Energy-Based View:**

▶ produce an energy function $E(Y,W)$ that:

▶ has low value in regions of high sample density

▶ has high(er) value everywhere else

**Energy function viewed as a negative log density**

▶ **Example: y = x^2**



Y

X

# Energy <-> Probability

$$P(Y|W) = \frac{e^{-\beta E(Y,W)}}{\int_y e^{-\beta E(y,W)}}$$

$$E(Y,W) \propto -\log P(Y|W)$$



P(Y|W)

Y

E(Y,W)

Y

# Training an Energy-Based Model

- Make the energy around training samples low

- Make the energy everywhere else higher



$$P(Y, W) = \frac{e^{-\beta E(Y,W)}}{\int_y e^{-\beta E(y,W)}}$$

# Training an Energy-Based Model to Approximate a Density

Maximizing P(Y|W) on training samples

make this big

$$P(Y|W) = \frac{e^{-\beta E(Y,W)}}{\int_y e^{-\beta E(y,W)}}$$

make this small

P(Y)

Y

Minimizing -log P(Y,W) on training samples

$$L(Y,W) = E(Y,W) + \frac{1}{\beta} \log \int_y e^{-\beta E(y,W)}$$

make this small

make this big

E(Y)

Y

# Training an Energy-Based Model with Gradient Descent

- **Gradient of the negative log-likelihood loss for one sample Y:**

$$\frac{\partial L(Y, W)}{\partial W} = \frac{\partial E(Y, W)}{\partial W} - \int_y P(y|W) \frac{\partial E(y, W)}{\partial W}$$

- **Gradient descent:**

$$W \leftarrow W - \eta \frac{\partial L(Y, W)}{\partial W}$$

Pushes down on the energy of the samples

Pulls up on the energy of low-energy Y's

$$W \leftarrow W - \eta \frac{\partial E(Y, W)}{\partial W} + \eta \int_y P(y|W) \frac{\partial E(y, W)}{\partial W}$$



*Yann LeCun*

New York University

# How do we push up on the energy of everything else?

- **Solution 1: contrastive divergence [Hinton 2000]**
  - Move away from a training sample a bit
  - Push up on that

- **Solution 2: score matching**
  - On the training samples: minimize the gradient of the energy, and maximize the trace of its Hessian.

- **Solution 3: denoising auto-encoder (not really energy-based)**
  - Train the inference dynamics to map noisy samples to clean samples

- **Solution 4: MAIN INSIGHT! [Ranzato, ..., LeCun AI-Stat 2007]**
  - **Restrict the information content of the code (features) Z**
  - If the code Z can only take a few different configurations, only a correspondingly small number of Ys can be perfectly reconstructed
  - Idea: impose a sparsity prior on Z
  - This is reminiscent of sparse coding [Olshausen & Field 1997]

# Encoder-Decoder with Sparsity (PSD)

[Kavukcuoglu, Ranzato, LeCun, 2009]

🔵 **Prediction the optimal code with a trained encoder**

🔵 **Energy = reconstruction_error + code_prediction_error + code_sparsity**

$$E(Y^i, Z) = \|Y^i - W_d Z\|^2 + \|Z - g_e(W_e, Y^i)\|^2 + \lambda \sum_j |z_j|$$

$$g_e(W_e, Y^i) = D \tanh(W_e Y)$$



$\|Y^i - \tilde{Y}\|^2$    $W_d Z$    $\lambda \sum_j \cdot$

**INPUT**   Y     Z    $|z_j|$   **FEATURES**

$g_e(W_e, Y^i)$    $\|Z - \tilde{Z}\|^2$

New York University

# The Main Insight [Ranzato et al. AISTATS 2007]

- **If the information content of the feature vector is limited (e.g. by imposing sparsity constraints), the energy MUST be large in most of the space.**
  - ▶ pulling down on the energy of the training samples will necessarily make a groove

- **The volume of the space over which the energy is low is limited by the entropy of the feature vector**
  - ▶ Input vectors are reconstructed from feature vectors.
  - ▶ If few feature configurations are possible, few input vectors can be reconstructed properly

New York University

# Why Limit the Information Content of the Code?

- Training sample

- Input vector which is NOT a training sample

- Feature vector

# Why Limit the Information Content of the Code?

- **Training sample**

- **Input vector which is NOT a training sample**

- **Feature vector**

*BAD: machine does not learn structure from training data!!*

*It just copies the data.*



INPUT SPACE

FEATURE SPACE

# Why Limit the Information Content of the Code?

- **Training sample**

- **Input vector which is NOT a training sample**

- **Feature vector**

*IDEA: reduce number of available codes.*



**INPUT SPACE**

**FEATURE SPACE**

Yann LeCun

New York University

# Why Limit the Information Content of the Code?

- **Training sample**

- **Input vector which is NOT a training sample**

- **Feature vector**

*IDEA: reduce number of available codes.*

# Sparsity Penalty to Restrict the Code

- **We are going to impose a sparsity penalty on the code to restrict its information content.**

- **We will allow the code to have higher dimension than the input**

- **Categories are more easily separable in high-dim sparse feature spaces**
  - This is a trick that SVM use: they have one dimension per sample

- **Sparse features are optimal when an active feature costs more than an inactive one (zero).**
  - e.g. neurons that spike consume more energy
  - The brain is about 2% active on average.

- 2 dimensional toy dataset
  - Mixture of 3 Cauchy distrib.
- Visualizing energy surface
(black = low, white = high)

[Ranzato 's PhD thesis 2009]

| | **PCA** (1 code unit) | **autoencoder** (3 code units) | **sparse coding** (3 code units) | **K-Means** (3 code units) |
|---|---|---|---|---|
| encoder | $W^T Y$ | $\sigma(W_e Y)$ | $-$ | $-$ |
| decoder | $WZ$ | $W_d Z$ | $WZ$ | $WZ$ |
| energy | $\|Y - WZ\|^2$ | $\|Y - WZ\|^2$ | $\|Y - WZ\|^2 + \lambda |Z|$ | $\|Y - WZ\|^2$ |
| loss | $F(Y)$ | $F(Y) + \log \Gamma$ | $F(Y)$ | $F(Y)$ |
| **pull-up** | dimens. | part. func. | sparsity | 1-of-N code |

- 2 dimensional toy dataset
  - spiral

- Visualizing energy surface
(black = low, white = high)

|  | **PCA** (1 code unit) | **autoencoder** (1 code unit) | **sparse coding** (20 code units) | **K-Means** (20 code units) |
|---|---|---|---|---|
| encoder | $W'Y$ | $\sigma(W_e Y)$ | $\sigma(W_e Z)$ | – |
| decoder | $WZ$ | $W_d Z$ | $W_d Z$ | $WZ$ |
| energy | $\|Y - WZ\|^2$ | $\|Y - WZ\|^2$ | $\|Y - WZ\|^2$ | $\|Y - WZ\|^2$ |
| loss | $F(Y)$ | $F(Y)$ | $F(Y)$ | $F(Y)$ |
| **pull-up** | dimens. | dimens. | sparsity | 1-of-N code |

# Using PSD to learn the features of an object recognition system



- Learning the filters of a ConvNet-like architecture with PSD

- 1. Train filters on images patches with PSD

- 2. Plug the filters into a ConvNet architecture

- 3. Train a supervised classifier on top

# "Modern" Object Recognition Architecture in Computer Vision



| Filter Bank | Non-Linearity | Spatial Pooling | Classifier |
|---|---|---|---|

| Oriented Edges | Sigmoid | Averaging |
|---|---|---|
| Gabor Wavelets | Rectification | Max pooling |
| Other Filters... | Vector Quant. | VQ+Histogram |
| | Contrast Norm. | Geometric Blurr |

- **Example:**
  - ▶ Edges + Rectification + Histograms + SVM [Dalal & Triggs 2005]
  - ▶ SIFT + classification

- **Fixed Features + "shallow" classifier**

*Yann LeCun*

| Filter Bank | Non-Linearity | feature Pooling | Filter Bank | Non-Linearity | feature Pooling | Classifier |

**Oriented Edges**   **WTA**   **Histogram (sum)**   **K-means**   **Pyramid Histogram (sum)**   **SVM with Histogram Intersection kernel**

**SIFT**

- 🔵 **Example:**
  - ▶ SIFT features with Spatial Pyramid Match Kernel SVM [Lazebnik et al. 2006]
- 🔵 **Fixed Features + unsupervised features + "shallow" classifier**

# Can't we get the same results with (deep) learning?



- Stacking multiple stages of feature extraction/pooling.

- Creates a hierarchy of features

- ConvNets and SIFT+PMK-SVM architectures are conceptually similar

- Can deep learning make a ConvNet match the performance of SIFT+PNK-SVM?

New York University

**Recognition Architecture**

# Procedure for a single-stage system

**1. Pre-process images**

- remove mean, high-pass filter, normalize contrast

**2. Train encoder-decoder on 9x9 image patches**

**3. use the filters in a recognition architecture**

- Apply the filters to the whole image
- Apply the tanh and D scaling
- Add more non-linearities (rectification, normalization)
- Add a spatial pooling layer

**4. Train a supervised classifier on top**

- Multinomial Logistic Regression or Pyramid Match Kernel SVM

# Using PSD Features for Recognition

**64 filters on 9x9 patches trained with PSD**

▶ with Linear-Sigmoid-Diagonal Encoder



weights :-0.2828 - 0.3043

New York University

# Feature Extraction

- **C** Convolution/sigmoid layer: filter bank? Learning, fixed Gabors?

# Feature Extraction

- **C**    Convolution/sigmoid layer: filter bank? Learning, fixed Gabors?
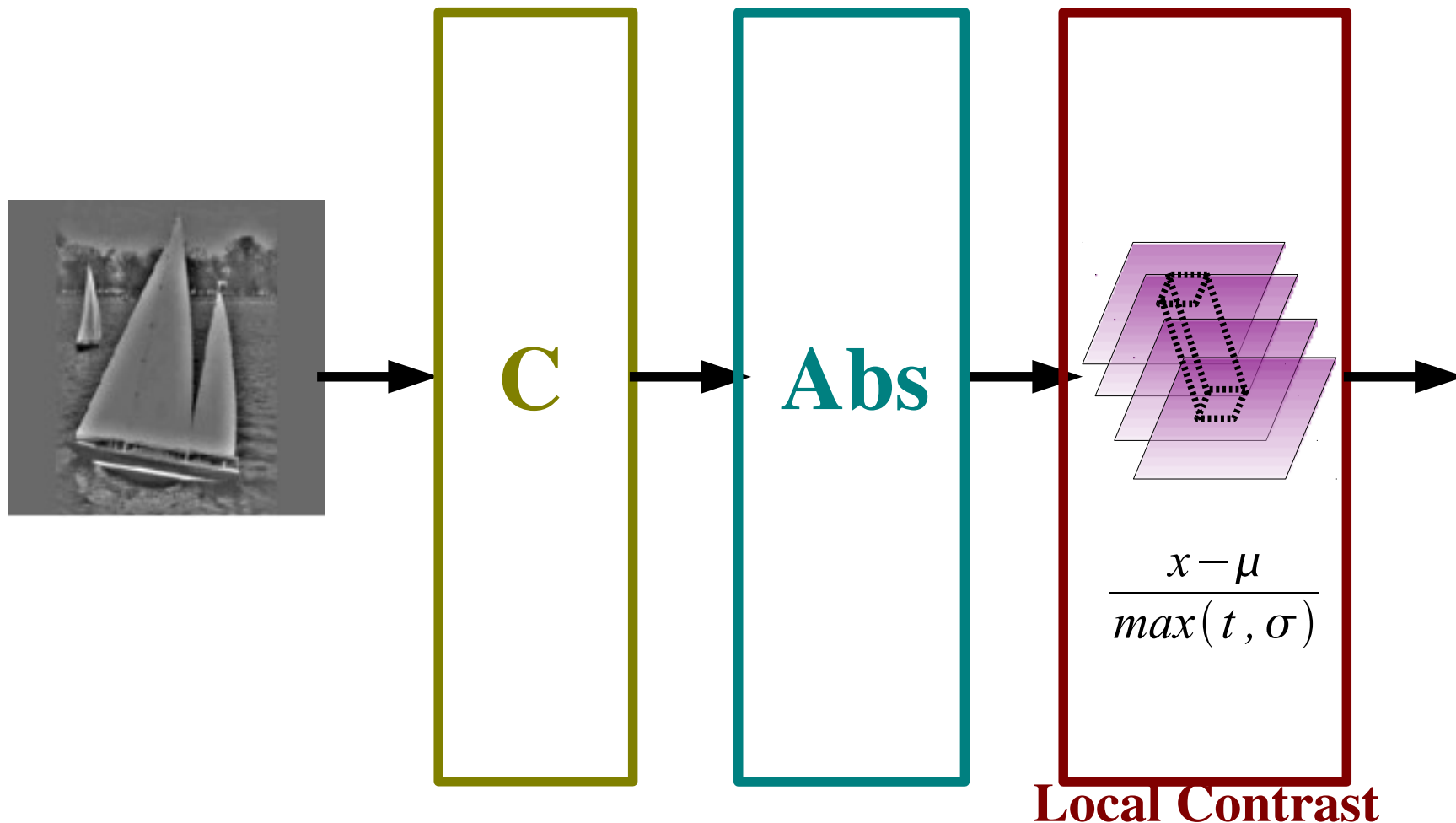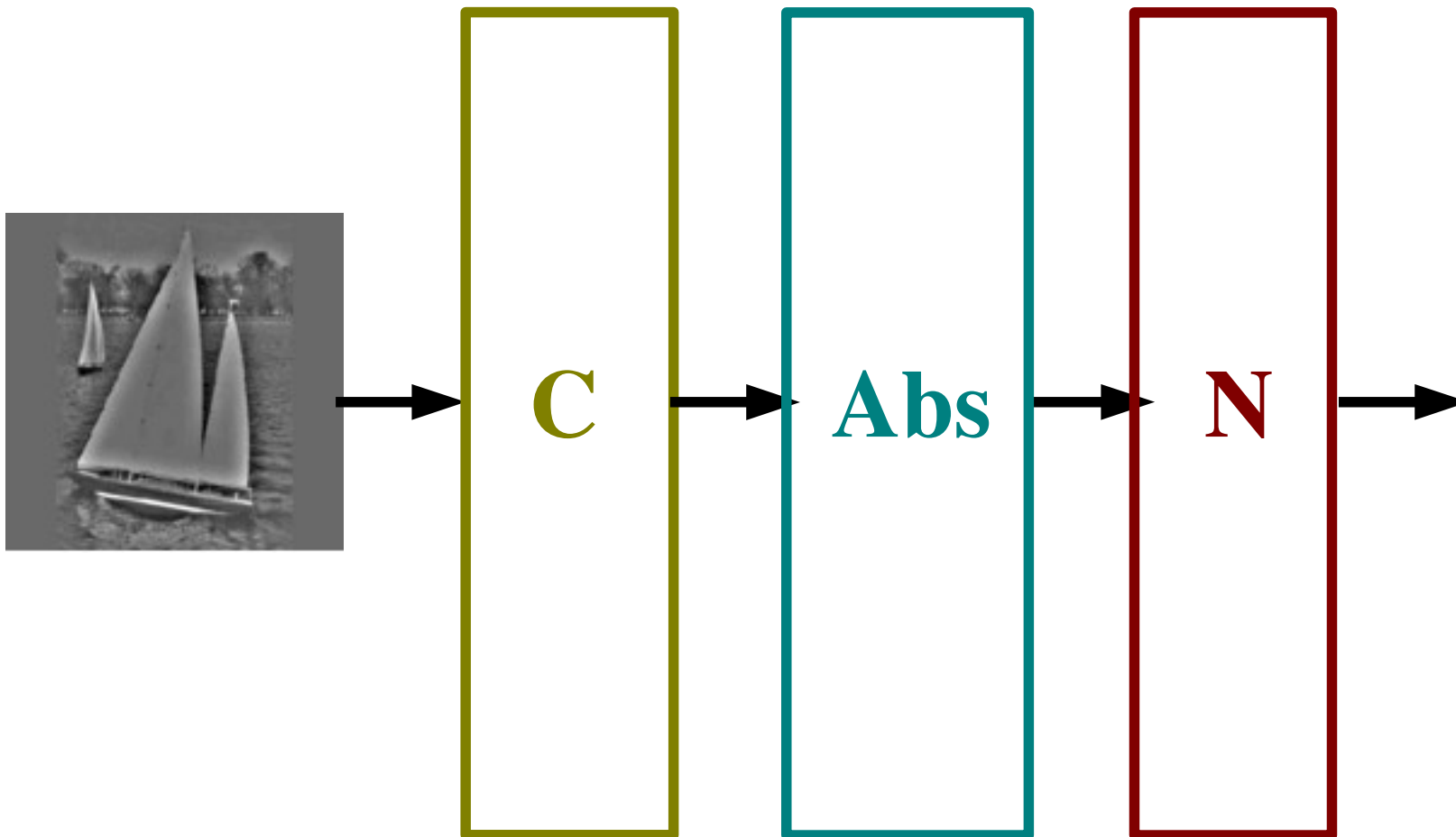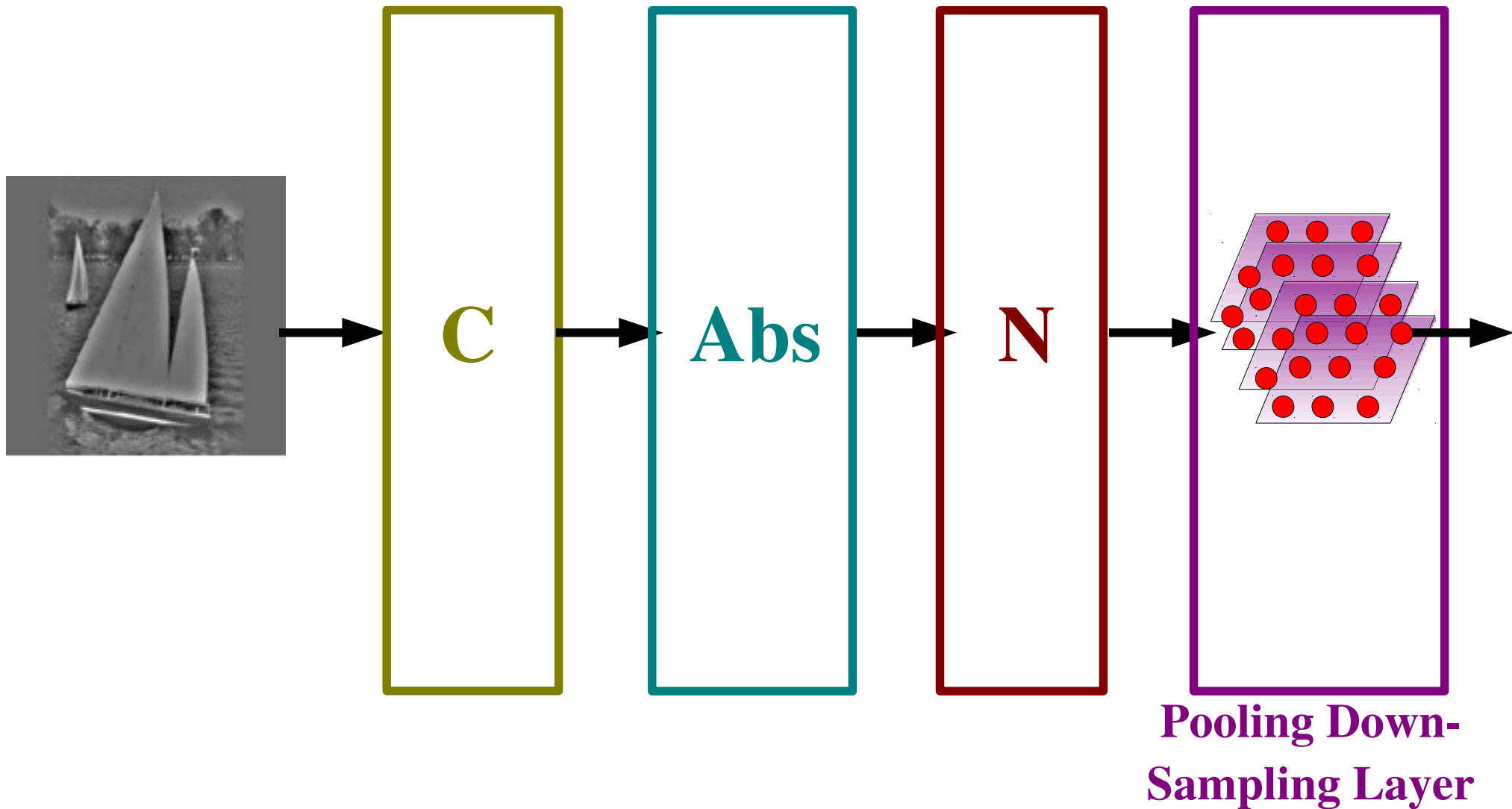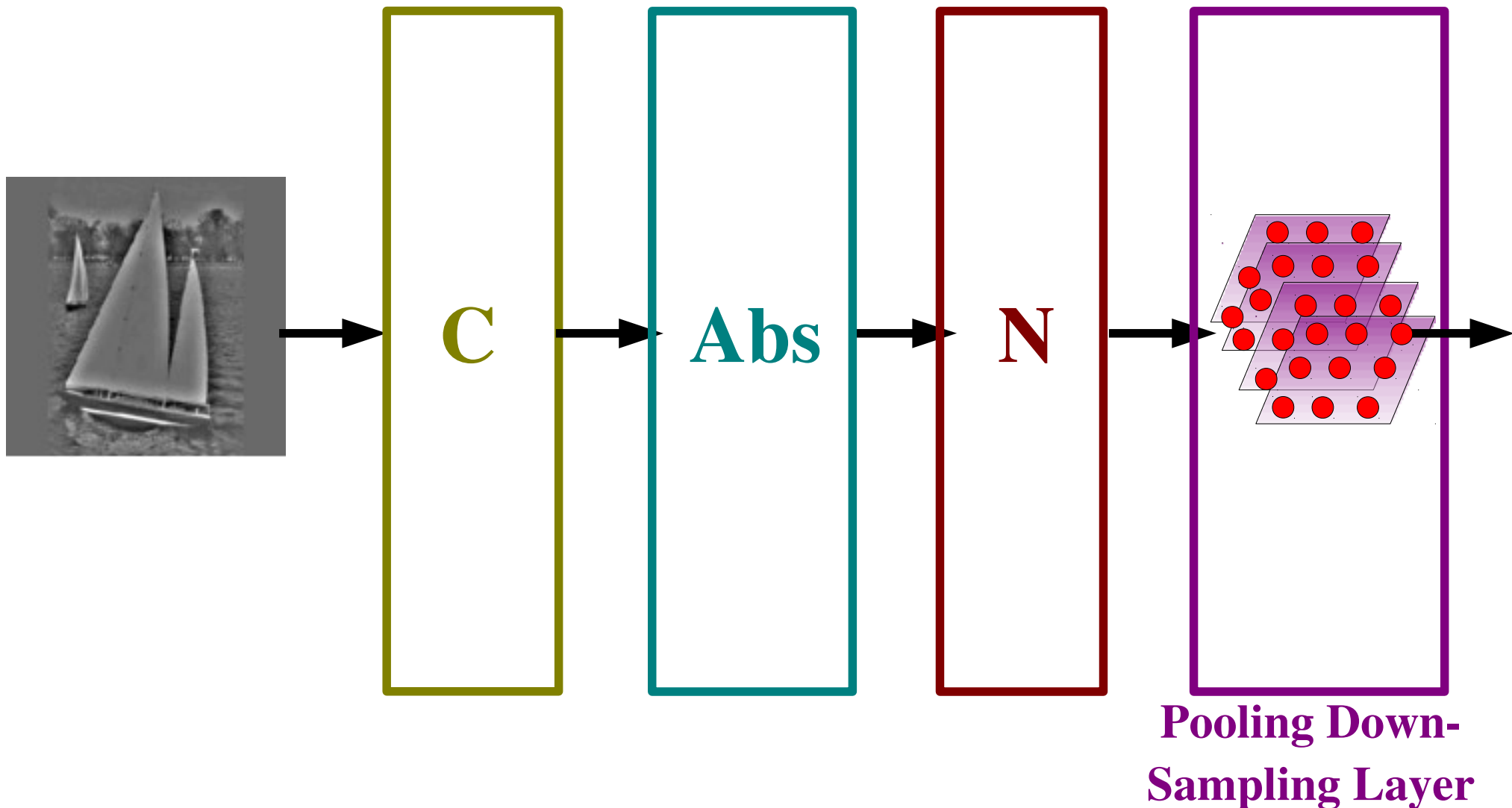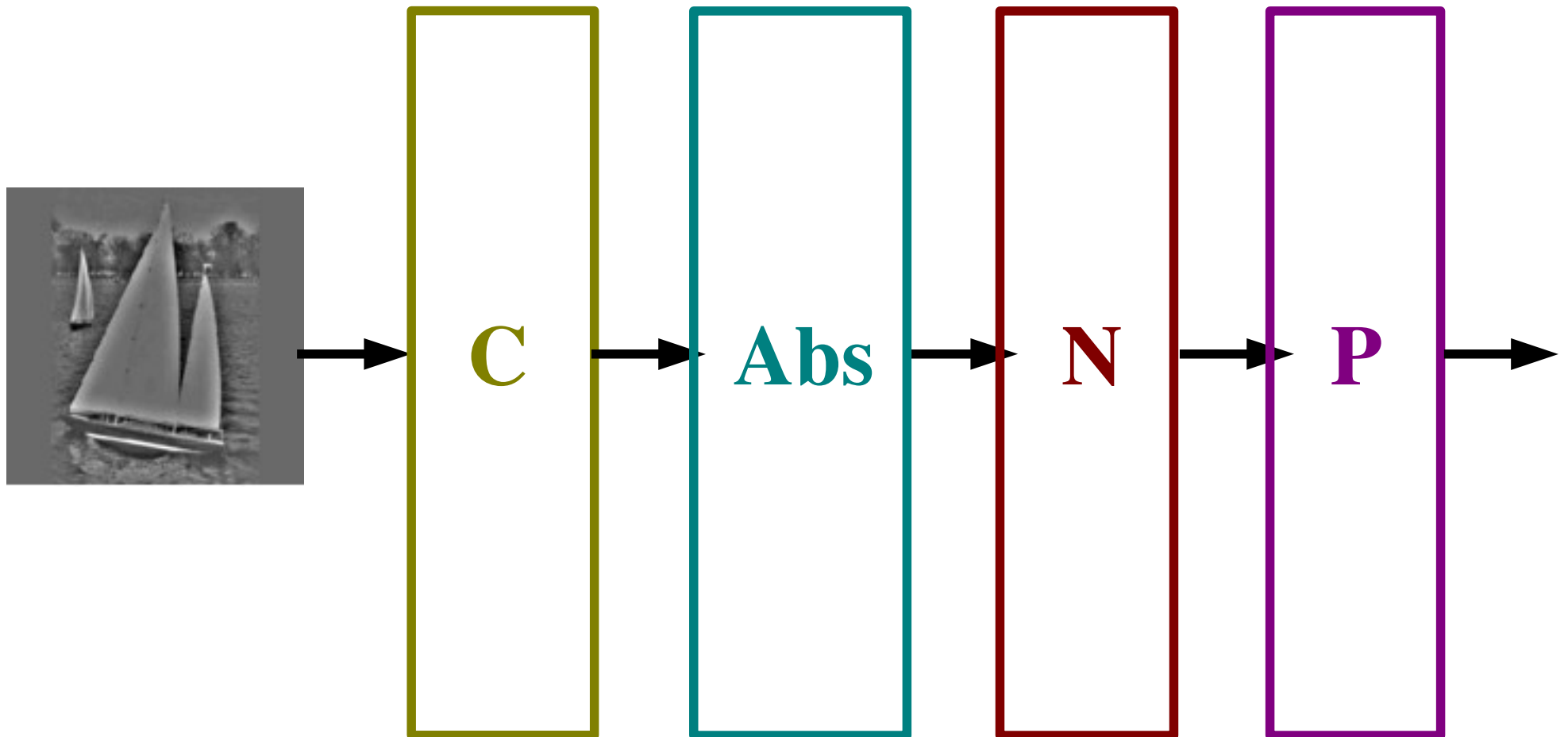


**RECTIFICATION LAYER**

Pinto, Cox and DiCarlo, PloS 08

# Feature Extraction

- **C** Convolution/sigmoid layer: filter bank? Learning, fixed Gabors?
- **Abs** Rectification layer: needed?



**C**

OR

**RECTIFICATION**

**LAYER**

Pinto, Cox and DiCarlo, PloS 08

# Feature Extraction

- **C** Convolution/sigmoid layer: filter bank? Learning, fixed Gabors?
- **Abs** Rectification layer: needed?



Pinto, Cox and DiCarlo, PloS 08

# Feature Extraction

- **C** Convolution/sigmoid layer: filter bank? Learning, fixed Gabors?
- **Abs** Rectification layer: needed?



**C**

**Abs**

$$\frac{x-\mu}{max(t,\sigma)}$$

**Local Contrast**
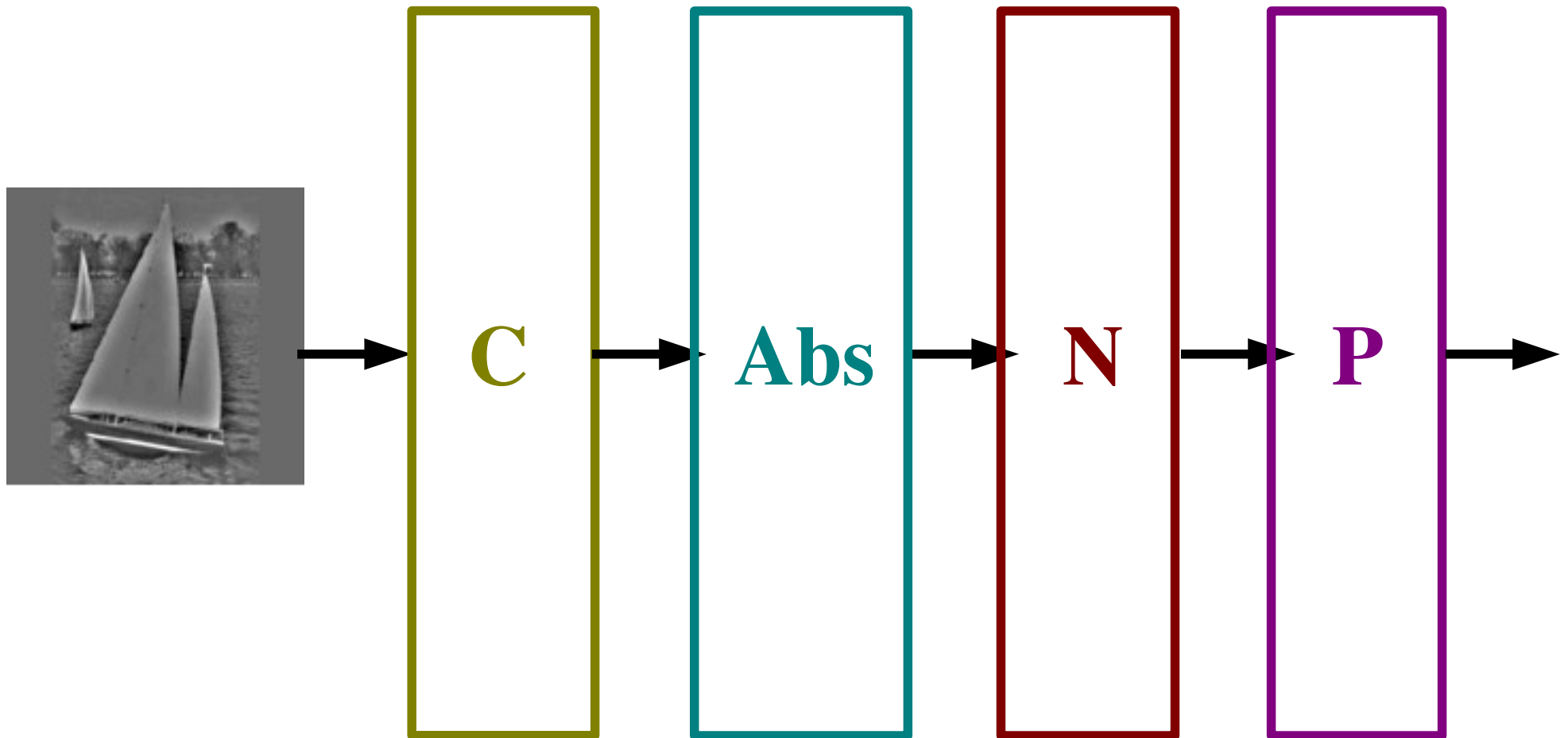
**Normalization Layer**

Pinto, Cox and DiCarlo, PloS 08

# Feature Extraction

- **C** Convolution/sigmoid layer: filter bank? Learning, fixed Gabors?
- **Abs** Rectification layer: needed?
- **N** Normalization layer: needed?



$$\frac{x-\mu}{max(t,\sigma)}$$

**Local Contrast**

**Normalization Layer**

Pinto, Cox and DiCarlo, PloS 08

# Feature Extraction

- **C** Convolution/sigmoid layer: filter bank? Learning, fixed Gabors?
- **Abs** Rectification layer: needed?
- **N** Normalization layer: needed?



Pinto, Cox and DiCarlo, PloS 08

# Feature Extraction

- **C** Convolution/sigmoid layer: filter bank? Learning, fixed Gabors?
- **Abs** Rectification layer: needed?
- **N** Normalization layer: needed?



**C** → **Abs** → **N** → **Pooling Down-Sampling Layer**

# Feature Extraction

- **C** Convolution/sigmoid layer: filter bank? Learning, fixed Gabors?
- **Abs** Rectification layer: needed?
- **N** Normalization layer: needed?
- **P** Pooling down-sampling layer: average or max?



**Pooling Down-Sampling Layer**

# Feature Extraction

- **C** Convolution/sigmoid layer: filter bank? Learning, fixed Gabors?
- **Abs** Rectification layer: needed?
- **N** Normalization layer: needed?
- **P** Pooling down-sampling layer: average or max?

# Feature Extraction

- **C** Convolution/sigmoid layer: filter bank? Learning, fixed Gabors?
- **Abs** Rectification layer: needed?
- **N** Normalization layer: needed?
- **P** Pooling down-sampling layer: average or max?



**THIS IS ONE STAGE OF FEATURE EXTRACTION**

# Training Protocol

- **Training**

  - Logistic Regression on Random Features: $R$

  - Logistic Regression on PSD features: $U$

  - Refinement of whole net from random with backprop: $R^+$

  - Refinement of whole net starting from PSD filters: $U^+$

- **Classifier**

  - Multinomial Logistic Regression or Pyramid Match Kernel SVM



Feature Extraction → Classification → *BOAT*

# Using PSD Features for Recognition

| $[64.\mathrm{F}_{\mathbf{CSG}}^{\mathbf{9\times9}} - \mathbf{R/N/P}^{\mathbf{5\times5}}]$ - log_reg | | | | |
|---|---|---|---|---|
| $\mathbf{R/N/P}$ | $\mathbf{R_{abs}-N-P_A}$ | $\mathbf{R_{abs}-P_A}$ | $\mathbf{N-P_M}$ | $\mathbf{N-P_A}$ / $\mathbf{P_A}$ |
| $\mathrm{U}^+$ | 54.2% | 50.0% | 44.3% | 18.5% / 14.5% |
| $\mathrm{R}^+$ | 54.8% | 47.0% | 38.0% | 16.3% / 14.3% |
| U | 52.2% | $43.3(\pm1.6)\%$ | 44.0% | 17.2% / 13.4% |
| R | 53.3% | 31.7% | 32.1% | 15.3% / $12.1(\pm2.2)\%$ |
| $[64.\mathrm{F}_{\mathbf{CSG}}^{\mathbf{9\times9}} - \mathbf{R/N/P}^{\mathbf{5\times5}}]$ - PMK | | | | |
| U | 65.0% | | | |
| $[96.\mathrm{F}_{\mathbf{CSG}}^{\mathbf{9\times9}} - \mathbf{R/N/P}^{\mathbf{5\times5}}]$ - PCA - lin_svm | | | | |
| U | 58.0% | | | |
| 96.Gabors - PCA - lin_svm (Pinto and DiCarlo 2006) | | | | |
| Gabors | 59.0% | | | |
| SIFT - PMK (Lazebnik et al. CVPR 2006) | | | | |
| Gabors | 64.6% | | | |

# Using PSD Features for Recognition

- **Rectification makes a huge difference:**
  - 14.5% -> 50.0%, without normalization
  - 44.3% -> 54.2% with normalization

- **Normalization makes a difference:**
  - 50.0 → 54.2

- **Unsupervised pretraining makes small difference**

- **PSD works just as well as SIFT**

- **Random filters work as well as anything!**
  - If rectification/normalization is present

- **PMK_SVM classifier works a lot better than multinomial log_reg on low-level features**
  - 52.2% → 65.0%

# Comparing Optimal Codes Predicted Codes on Caltech 101

**Approximated Sparse Features Predicted by PSD give better recognition results than Optimal Sparse Features computed with Feature Sign!**

▶ PSD features are more stable.



Feature Sign (FS) is an optimization methods for computing sparse codes [Lee...Ng 2006]

New York University

# PSD Features are more stable

- **Approximated Sparse Features Predicted by PSD give better recognition results than Optimal Sparse Features computed with Feature Sign!**

- **Because PSD features are more stable. Feature obtained through sparse optimization can change a lot with small changes of the input.**



How many features change sign in patches from successive video frames (a,b), versus patches from random frame pairs (c)

# PSD features are much cheaper to compute

🔵 **Computing PSD features is hundreds of times cheaper than Feature Sign.**

# How Many 9x9 PSD features do we need?

- Accuracy increases slowly past 64 filters.

# Training a Multi-Stage Hubel-Wiesel Architecture with PSD

| | Filter Bank | Non-Linearity | Spatial Pooling | Filter Bank | Non-Linearity | Spatial Pooling | Classifier |

- 1. Train stage-1 filters with PSD on patches from natural images

- 2. Compute stage-1 features on training set

- 3. Train state-2 filters with PSD on stage-1 feature patches

- 4. Compute stage-2 features on training set

- 5. Train linear classifier on stage-2 features

- 6. Refine entire network with supervised gradient descent

- What are the effects of the non-linearities and unsupervised pretraining?

New York University

# Multistage Hubel-Wiesel Architecture on Caltech-101



Y (luminance)

U

V

CONVOLUTIONS (9x9)

MAX/SUBSAMPLING (4x4)

CONVOLUTIONS (9x9)

MAX/SUBSAMPLING (5x5)

INPUT 3@140x140  32@132x132  32@33x33  64@25x25  64@5x5

Yann LeCun

New York University

# Multistage Hubel-Wiesel Architecture

🔵 **Image Preprocessing:**

▶ High-pass filter, local contrast normalization (divisive)

🔵 **First Stage:**

▶ Filters: 64  9x9 kernels producing 64 feature maps
▶ Pooling: 10x10 averaging with 5x5 subsampling

🔵 **Second Stage:**

▶ Filters: 4096  9x9 kernels producing 256 feature maps
▶ Pooling: 6x6 averaging with 3x3 subsampling
▶ Features: 256 feature maps of size 4x4 (4096 features)

🔵 **Classifier Stage:**

▶ Multinomial logistic regression

🔵 **Number of parameters:**

▶ Roughly 750,000

New York University

| Single Stage System: $[64.\mathrm{F}^{9\times9}_{\mathrm{CSG}} - \mathrm{R/N/P}^{5\times5}]$ - log_reg | | | | | |
|---|---|---|---|---|---|
| R/N/P | $\mathrm{R_{abs}} - \mathrm{N} - \mathrm{P_A}$ | $\mathrm{R_{abs}} - \mathrm{P_A}$ | $\mathrm{N} - \mathrm{P_M}$ | $\mathrm{N} - \mathrm{P_A}$ | $\mathrm{P_A}$ |
| $\mathrm{U}^+$ | 54.2% | 50.0% | 44.3% | 18.5% | 14.5% |
| $\mathrm{R}^+$ | 54.8% | 47.0% | 38.0% | 16.3% | 14.3% |
| U | 52.2% | 43.3%($\pm$1.6) | 44.0% | 17.2% | 13.4% |
| R | 53.3% | 31.7% | 32.1% | 15.3% | 12.1%($\pm$2.2) |
| G | 52.3% | | | | |

| Two Stage System: $[64.\mathrm{F}^{9\times9}_{\mathrm{CSG}} - \mathrm{R/N/P}^{5\times5}] - [256.\mathrm{F}^{9\times9}_{\mathrm{CSG}} - \mathrm{R/N/P}^{4\times4}]$ - log_reg | | | | | |
|---|---|---|---|---|---|
| R/N/P | $\mathrm{R_{abs}} - \mathrm{N} - \mathrm{P_A}$ | $\mathrm{R_{abs}} - \mathrm{P_A}$ | $\mathrm{N} - \mathrm{P_M}$ | $\mathrm{N} - \mathrm{P_A}$ | $\mathrm{P_A}$ |
| $\mathrm{U}^+\mathrm{U}^+$ | 65.5% | 60.5% | 61.0% | 34.0% | 32.0% |
| $\mathrm{R}^+\mathrm{R}^+$ | 64.7% | 59.5% | 60.0% | 31.0% | 29.7% |
| UU | 63.7% | 46.7% | 56.0% | 23.1% | 9.1% |
| RR | 62.9% | 33.7%($\pm$1.5) | 37.6%($\pm$1.9) | 19.6% | 8.8% |
| GT | 55.8% | $\leftarrow$ like HMAX model | | | |

| Single Stage: $[64.\mathrm{F}^{9\times9}_{\mathrm{CSG}} - \mathrm{R/N/P}^{5\times5}]$ - PMK-SVM | |
|---|---|
| U | 64.0% |

| Two Stages: $[64.\mathrm{F}^{9\times9}_{\mathrm{CSG}} - \mathrm{R/N/P}^{5\times5}] - [256.\mathrm{F}^{9\times9}_{\mathrm{CSG}} - \mathrm{R/N}]$ - PMK-SVM | |
|---|---|
| UU | 52.8% |

*Yann LeCun*

# Two-Stage Result Analysis

- **Second Stage + logistic regression = PMK_SVM**

- **Unsupervised pre-training doesn't help much :-(**

- <span style="color:red">**Random filters work amazingly well with normalization**</span>

- **Supervised global refirnement helps a bit**

- **The best system is really cheap**

- **Either use rectification and average pooling or no rectification and max pooling.**

# Multistage Hubel-Wiesel Architecture: Filters

**After PSD**

**After supervised refinement**

**Stage 1**

**Stage2**

# Demo: real-time learning of visual categories



Input
high-pass filtered
contrast-normalized
83x83 (raw: 91x91)

**STAGE 1**
Filter Bank + Tanh + Gain
Abs + Contrast Norm + Pooling + Downsampling

64 features 75x75
64 filters
9x9 kernels

64 features 14x14
5x5 subsampling
10x10 pooling

**STAGE 2**
Filter Bank + Tanh + Gain
Abs + Contrast Norm + Pooling + Downsampling

256 features 6x6
4096 filters
9x9 kernels

256 features 1x1
4x4 subsampling
6x6 pooling

**CLASSIFIER**
Parzen Windows Classifier

*Yann LeCun*

New York University

# MNIST dataset

- 10 classes and up to 60,000 training samples per class

# Architecture

**MNIST dataset**

- $U^+U^+$: 0.53% error (this is a record on the undistorted MNIST!)

- Comparison: $RR$ versus $UU$ and $R^+R^+$



Classification error on the MNIST dataset

# Why Random Filters Work?

# The Competition: SIFT + Sparse-Coding + PMK-SVM

**Replacing K-means with Sparse Coding**

▶ [Yang 2008] [Boureau, Bach, Ponce, LeCun 2010]

| | Method | Caltech 15 | Caltech 30 | Scenes |
|---|---|---|---|---|
| Boiman et al. [1] | Nearest neighbor + spatial correspondence | $65.00 \pm 1.14$ | 70.40 | - |
| Jain et al. [8] | Fast image search for learned metrics | 61.00 | 69.60 | - |
| Lazebnik et al. [12] | Spatial Pyramid + hard quantization + kernel SVM | 56.40 | $64.40 \pm 0.80$ | $81.40 \pm 0.50$ |
| van Gemert et al. [24] | Spatial Pyramid + soft quantization + kernel SVM | – | $64.14 \pm 1.18$ | $76.67 \pm 0.39$ |
| Yang et al. [26] | SP + sparse codes + max pooling + linear | **$67.00 \pm 0.45$** | **$73.2 \pm 0.54$** | $80.28 \pm 0.93$ |
| Zhang et al. [27] | $k$NN-SVM | $59.10 \pm 0.60$ | $66.20 \pm 0.50$ | - |
| Zhou et al. [29] | SP + Gaussian mixture | – | – | **$84.1 \pm 0.5$** |
| Baseline: | SP + hard quantization + avg pool + kernel SVM | $56.74 \pm 1.31$ | $64.19 \pm 0.94$ | $80.89 \pm 0.21$ |
| Unsupervised coding | SP + soft quantization + avg pool + kernel SVM | $59.12 \pm 1.51$ | $66.42 \pm 1.26$ | $81.52 \pm 0.54$ |
| $1 \times 1$ features | SP + soft quantization + max pool + kernel SVM | $63.61 \pm 0.88$ | – | $83.41 \pm 0.57$ |
| 8 pixel grid resolution | SP + sparse codes + avg pool + kernel SVM | $62.85 \pm 1.22$ | $70.27 \pm 1.29$ | $83.15 \pm 0.35$ |
| | SP + sparse codes + max pool + kernel SVM | $64.62 \pm 0.94$ | **$71.81 \pm 0.96$** | **$84.25 \pm 0.35$** |
| | SP + sparse codes + max pool + linear | **$64.71 \pm 1.05$** | $71.52 \pm 1.13$ | $83.78 \pm 0.53$ |
| Macrofeatures + | SP + sparse codes + max pool + kernel SVM | **$69.03 \pm 1.17$** | **$75.72 \pm 1.06$** | **$84.60 \pm 0.38$** |
| Finer grid resolution | SP + sparse codes + max pool + linear | $68.78 \pm 1.09$ | $75.14 \pm 0.86$ | $84.41 \pm 0.26$ |

# Small NORB dataset

- 5 classes and up to 24,300 training samples per class

# NORB Generic Object Recognition Dataset

- **50** toys belonging to 5 categories: **animal, human figure, airplane, truck, car**
- **10** instance per category: 5 instances used for training, 5 instances for testing
- **Raw dataset: 972** stereo pair of each object instance. **48,600** image pairs total.

- **For each instance:**
- **18 azimuths**
  - 0 to 350 degrees every 20 degrees
- **9 elevations**
  - 30 to 70 degrees from horizontal every 5 degrees
- **6 illuminations**
  - on/off combinations of 4 lights
- **2 cameras (stereo)**
  - 7.5 cm apart
  - 40 cm from the object



**Training instances**          **Test instances**

Yann LeCun

# Small NORB dataset

- **Two-stage system: error rate versus number of labeled training samples**

# Learning Complex Cells with Invariance Properties

[Kavukcuoglu et al. CVPR 2008]

- **Unsupervised PSD ignores the spatial pooling step.**

- **Could we devise a similar method that learns the pooling layer as well?**

- **Idea [Hyvarinen & Hoyer 2001]: group sparsity on pools of features**
  - Minimum number of pools must be non-zero
  - Number of features that are on within a pool doesn't matter
  - Polls tend to regroup similar features

$$\|Y^i - \tilde{Y}\|^2 \qquad W_d Z \qquad \lambda \sum_j .$$

INPUT $\quad Y \qquad Z \quad \sqrt{\left(\sum_{(k \in P_j)} Z_k^2\right)}$

FEATURES

$$g_e(W_e, Y^i) \qquad \|Z - \tilde{Z}\|^2$$

# Learning the filters and the pools

**Using an idea from Hyvarinen: topographic square pooling (subspace ICA)**

- ▶ 1. Apply filters on a patch (with suitable non-linearity)
- ▶ 2. Arrange filter outputs on a 2D plane
- ▶ 3. square filter outputs
- ▶ 4. minimize sqrt of sum of blocks of squared filter outputs

$$\text{Overall Sparsity term:} \quad \sum_{i=1}^{K} \sqrt{v_i^2}$$

$$v_1^2 = \sum_{j \in P_1} (w_j z_j)^2 \qquad v_K^2 = \sum_{j \in P_K} (w_j z_j)^2$$

$P_1$ $P_2$ ......

Overlapping Neighborhoods $P_i$

Gaussian Window $w_j$

Map of features

$z$

$P_K$

Units in the code Z       Define pools and enforce sparsity across pools

# Learning the filters and the pools

- **The filters arrange themselves spontaneously so that similar filters enter the same pool.**

- **The pooling units can be seen as complex cells**

- **They are invariant to local transformations of the input**
  - ▶ For some it's translations, for others rotations, or other transformations.

# Pinwheels?

New York University

# Invariance Properties Compared to SIFT

- **Measure distance between feature vectors (128 dimensions) of 16x16 patches from natural images**
  - ▶ Left: normalized distance as a function of translation
  - ▶ Right: normalized distance as a function of translation when one patch is rotated 25 degrees.

- **Topographic PSD features are more invariant than SIFT**

# Learning Invariant Features

**Recognition Architecture**

▶ ->HPF/LCN->filters->tanh->sqr->pooling->sqrt->Classifier

▶ Block pooling plays the same role as rectification

▶ A/B Comparison with SIFT (128x34x34 descriptors)
▶ 32x16 topographic map with 16x16 filters
▶ Pooling performed over 6x6 with 2x2 subsampling
▶ 128 dimensional feature vector per 16x16 patch
▶ Feature vector computed every 4x4 pixels (128x34x34 feature    maps)
▶ Resulting feature maps are spatially smoothed

| Method | Av. Accuracy/Class (%) |
|---|---|
| local norm$_{5\times5}$ + boxcar$_{5\times5}$ + PCA$_{3060}$ + linear SVM | |
| IPSD (24x24) | 50.9 |
| SIFT (24x24) (non rot. inv.) | 51.2 |
| SIFT (24x24) (rot. inv.) | 45.2 |
| Serre et al. features [25] | 47.1 |
| local norm$_{9\times9}$ + Spatial Pyramid Match Kernel SVM | |
| SIFT [11] | 64.6 |
| IPSD (34x34) | 59.6 |
| IPSD (56x56) | 62.6 |
| IPSD (120x120) | 65.5 |

Yann

# Recognition Accuracy on Tiny Images & MNIST

▶ A/B Comparison with SIFT (128x5x5 descriptors)
▶ 32x16 topographic map with 16x16 filters.

| Performance on Tiny Images Dataset | |
| --- | --- |
| Method | Accuracy (%) |
| IPSD (5x5) | 54 |
| SIFT (5x5) (non rot. inv.) | 53 |

| Performance on MNIST Dataset | |
| --- | --- |
| Method | Error Rate (%) |
| IPSD (5x5) | 1.0 |
| SIFT (5x5) (non rot. inv.) | 1.5 |

New York University

# Learning fields of Convolutional Filters

# Convolutional Training

🌐 **Problem:**

▶ With patch-level training, the learning algorithm must reconstruct the entire patch with a single feature vector

▶ But when the filters are used convolutionally, neighboring feature vectors will be highly redundant



weights :-0.2828 - 0.3043

# Convolutional Training

🔵 **Problem with patch-based training: high correlation between outputs of filters from overlapping receptive fields.**

# Convolutional Training

🔵 **Filters and Basis Functions obtained with 1, 2, 4, 8, 16, 32, and 64 filters.**



*Yann LeCun*

New York University

# Learning fields of Simple Cells and Complex Cells

[Gregor and LeCun, 2010]

# Training Simple Cells with Local Receptive Fields over Large Input Images

🔵 **Training on 115x115 images. Kernels are 15x15**

🔵 **Training on 115x115 images. Kernels are 15x15**

K Obermayer and GG Blasdel, Journal of Neuroscience, Vol 13, 4114-4129 (**Monkey**)

119x119 Image Input

100x100 Code

20x20 Receptive field size

sigma=5

Michael C. Crair, et. al. The Journal of Neurophysiology Vol. 77 No. 6 June 1997, pp. 3381-3385 (**Cat**)

# Same Method, withTraining at the Image Level (vs patch)

- Color indicates orientation (by fitting Gabors)

New York University

# Deep Learning for Mobile Robot Vision

*Yann LeCun*

# DARPA/LAGR: Learning Applied to Ground Robotics

- Getting a robot to drive autonomously in unknown terrain solely from vision (camera input).

- Our team (NYU/Net-Scale Technologies Inc.) was one of 8 participants funded by DARPA

- All teams received identical robots and can only modify the software (not the hardware)

- The robot is given the GPS coordinates of a goal, and must drive to the goal as fast as possible. The terrain is unknown in advance. The robot is run 3 times through the same course.

- Long-Range Obstacle Detection with on-line, self-trained ConvNet

- Uses temporal consistency!

New York University

Obstacles overlaid with camera image



Camera image                    Detected obstacles (red)

# Navigating to a goal is hard...



stereo perspective



human perspective

## especially in a snowstorm.

New York University

# Self-Supervised Learning

- Stereo vision tells us what nearby obstacles look like

- Use the labels (obstacle/traversible) produced by stereo vision to train a monocular neural network

- Self-supervised "near to far" learning

# Long Range Vision: Distance Normalization



**Pre-processing** (125 ms)

- Ground plane estimation
- Horizon leveling
- Conversion to YUV + local contrast normalization
- Scale invariant pyramid of distance-normalized image "bands"



112.3m to INF, scale: 1.0

50.7m to INF, scale: 1.4

24.2m to INF, scale: 1.9

13.8m to 86.8m, scale: 2.6

9.0m to 34.5m, scale: 3.5

5.8m to 17.6m, scale: 5.0

4.1m to 11.3m, scale: 6.7

net>SCALE
Technologies, Inc.

NEW YORK UNIVERSITY

# Convolutional Net Architecture

- Operates on 12x25 YUV windows from the pyramid

**Logistic regression 100 features -> 5 classes**

**100 features per**

**3x12x25 input window**

100x1x1 input window

**Convolutions with 6x5 kernels**

20x6x5 input window

**Pooling/subsampling with 1x4 kernels**

20x6x20 input window

**Convolutions with 7x6 kernels**

**YUV image band**

**20-36 pixels tall,**

**36-500 pixels wide**

3x12x25 input window

net SCALE
Technologies, Inc.

NEW YORK UNIVERSITY

# Convolutional Net Architecture

100@25x121

**CONVOLUTIONS (6x5)**

20@30x125

**MAX SUBSAMPLING (1x4)**

20@30x484

**CONVOLUTIONS (7x6)**

3@36x484

YUV input

net SCALE
Technologies, Inc.

NEW YORK UNIVERSITY

# Long Range Vision: 5 categories

**Online Learning** (52 ms)

- Label windows using stereo information – 5 classes



| | | | | |
|---|---|---|---|---|
| **super-ground** | **ground** | **footline** | **obstacle** | **super-obstacle** |

# Trainable Feature Extraction

- "Deep belief net" approach to unsupervised feature learning

- Two stages are trained in sequence
  - each stage has a layer of convolutional filters and a layer of horizontal feature pooling.
  - Naturally shift invariant in the horizontal direction

- Filters of the convolutional net are trained so that the input can be reconstructed from the features
  - 20 filters at the first stage (layers 1 and 2)
  - 300 filters at the second stage (layers 3 and 4)

- Scale invariance comes from pyramid.
  - for near-to-far generalization

# Long Range Vision Results



Input image

Stereo Labels

Classifier Output

Input image

Stereo Labels

Classifier Output

# Long Range Vision Results



Input image

Stereo Labels

Classifier Output

Input image

Stereo Labels

Classifier Output

net>SCALE
Technologies, Inc.

NEW YORK UNIVERSITY

# Long Range Vision Results



Input image | Stereo Labels | Classifier Output

Input image | Stereo Labels | Classifier Output

net>SCALE
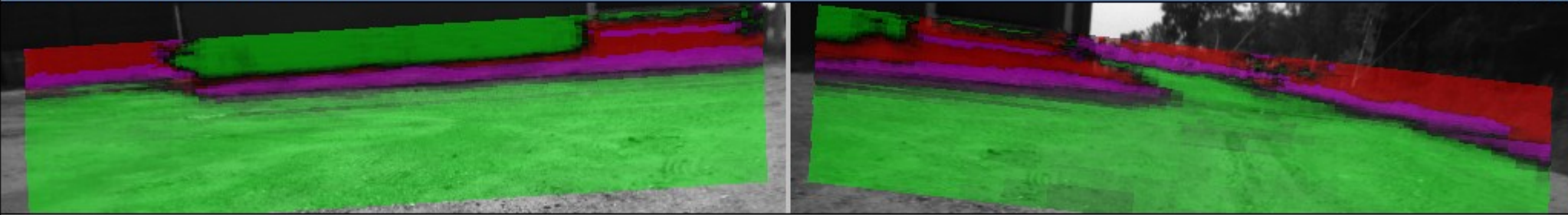Technologies, Inc.

NEW YORK UNIVERSITY

**Vehicle Map (Hyperbolic Polar map)**

Legend
- ○ Goal
- Path Planning
- Trajectories
- Traversible
- Uncertain
- Quasi-Lethal
- Lethal
- Bumper/Stuck
- Unseen

200m
100m
50m
25m
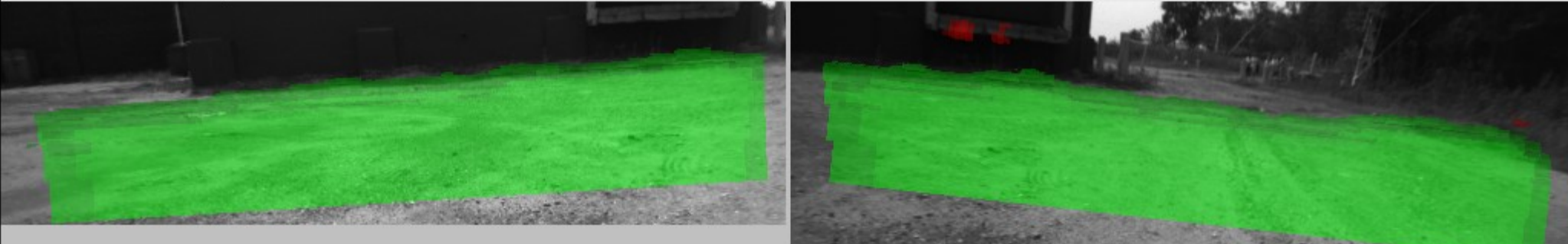15m
10m
5m
-5m
-10m
-15m
-25m
-50m
-100m
-200m

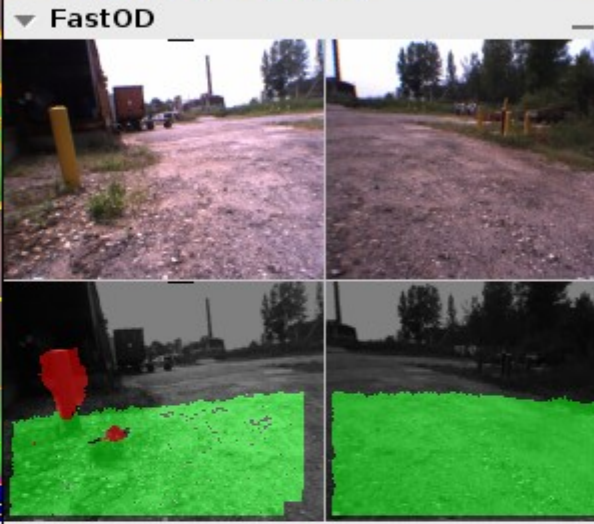Cost Map
(FastOD
& FarOD)

RGB Map
(FarOD Only)

**FastOD**

**FarOD Neural Network Labels**

**FarOD Stereo: Input labels to Neural Network**

netSCALE Technologies, Inc.

NEW YORK UNIVERSITY

# Vehicle Map (Hyperbolic Polar map)

Legend

- ○ Goal
- □ Path Planning
- ‖ Trajectories
- ■ Traversible
- ■ Uncertain
- ■ Quasi-Lethal
- ■ Lethal
- ■ Bumper/Stuck
- ■ Unseen

200m
100m
50m

25m

15m

10m

5m

-5m

-10m

-15m

-25m

-50m
-100m
-200m

Cost Map
(FastOD
& FarOD)

Goal

RGB Map
(FarOD Only)

Goal

FastOD

FarOD Neural Network Labels

FarOD Stereo: Input labels to Neural Network

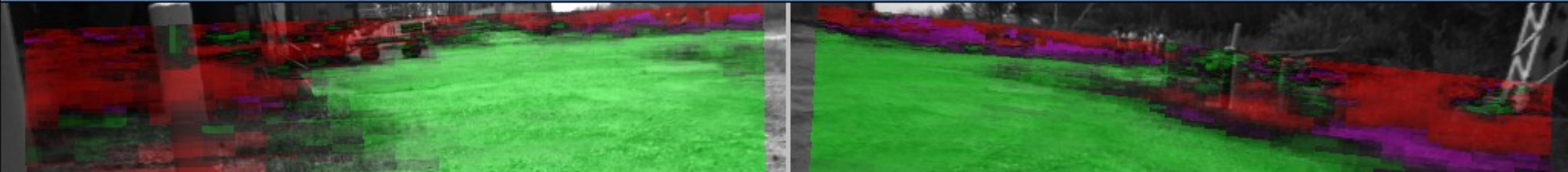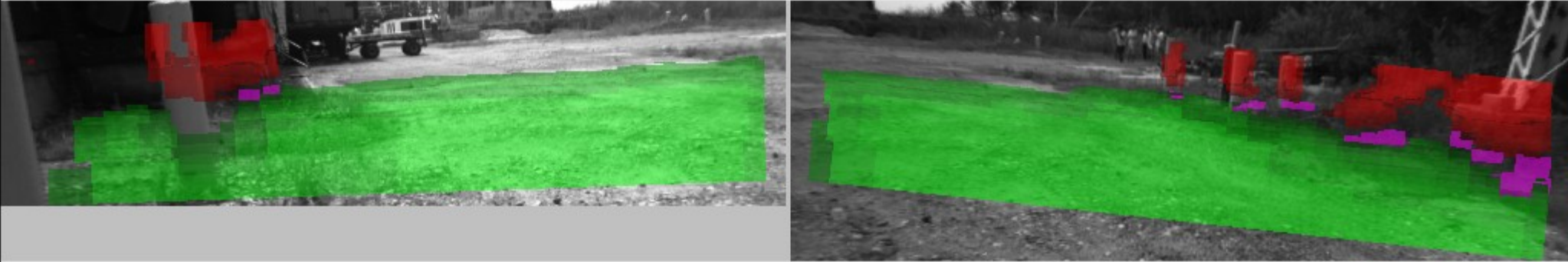NEW YORK UNIVERSITY

**Vehicle Map (Hyperbolic Polar map)**

Legend
- ○ Goal
- ☐ Path Planning
- Trajectories
- Traversible
- Uncertain
- Quasi-Lethal
- Lethal
- Bumper/Stuck
- Unseen

Cost Map
(FastOD
& FarOD)

RGB Map
(FarOD Only)

200m 100m 50m
25m
15m
10m
5m
-5m
-10m
-15m
-25m
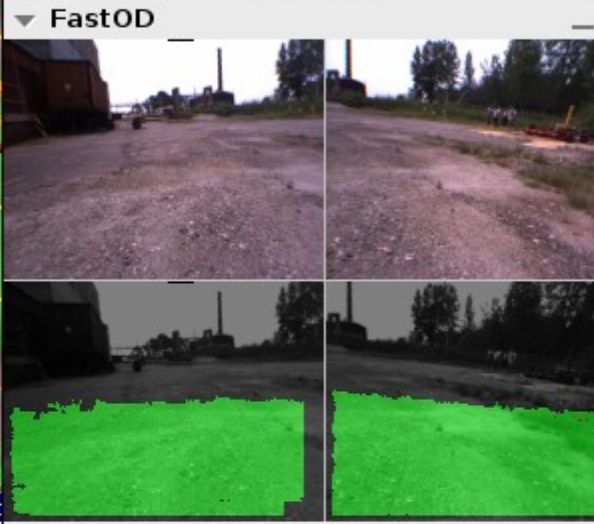-50m -100m -200m

Goal

**FastOD**

**FarOD Neural Network Labels**

**FarOD Stereo: Input labels to Neural Network**

NEW YORK UNIVERSITY

Technologies, Inc.

# Vehicle Map (Hyperbolic Polar map)

**Legend**
- ○ Goal
- ■ Path Planning
- Trajectories
- Traversible
- Uncertain
- Quasi-Lethal
- Lethal
- Bumper/Stuck
- Unseen

200m
100m
50m
25m

15m

10m

5m

-5m

-10m

-15m

-25m

-50m
-100m
-200m

Cost Map (FastOD & FarOD)   Goal
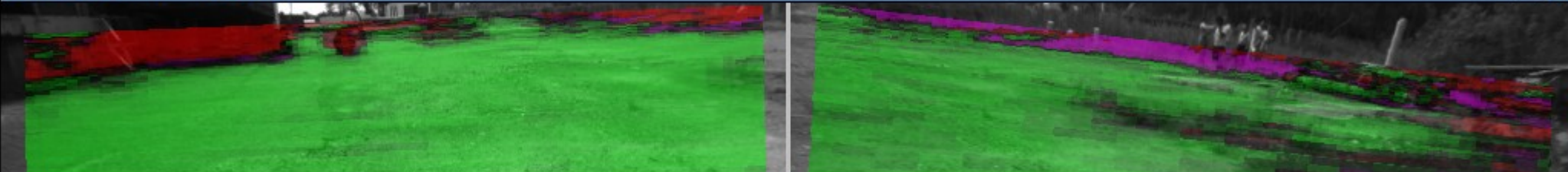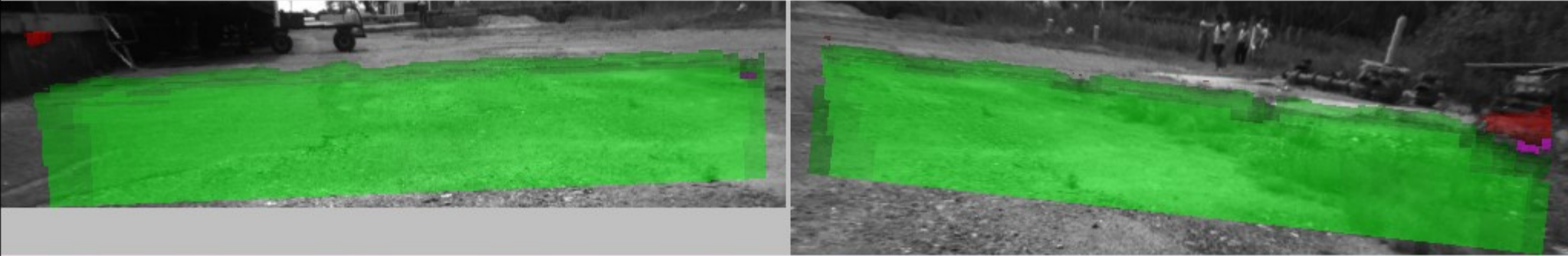
RGB Map (FarOD Only)   Goal

## ▾ FastOD

## ▾ FarOD Neural Network Labels

## ▾ FarOD Stereo: Input labels to Neural Network

NEW YORK UNIVERSITY

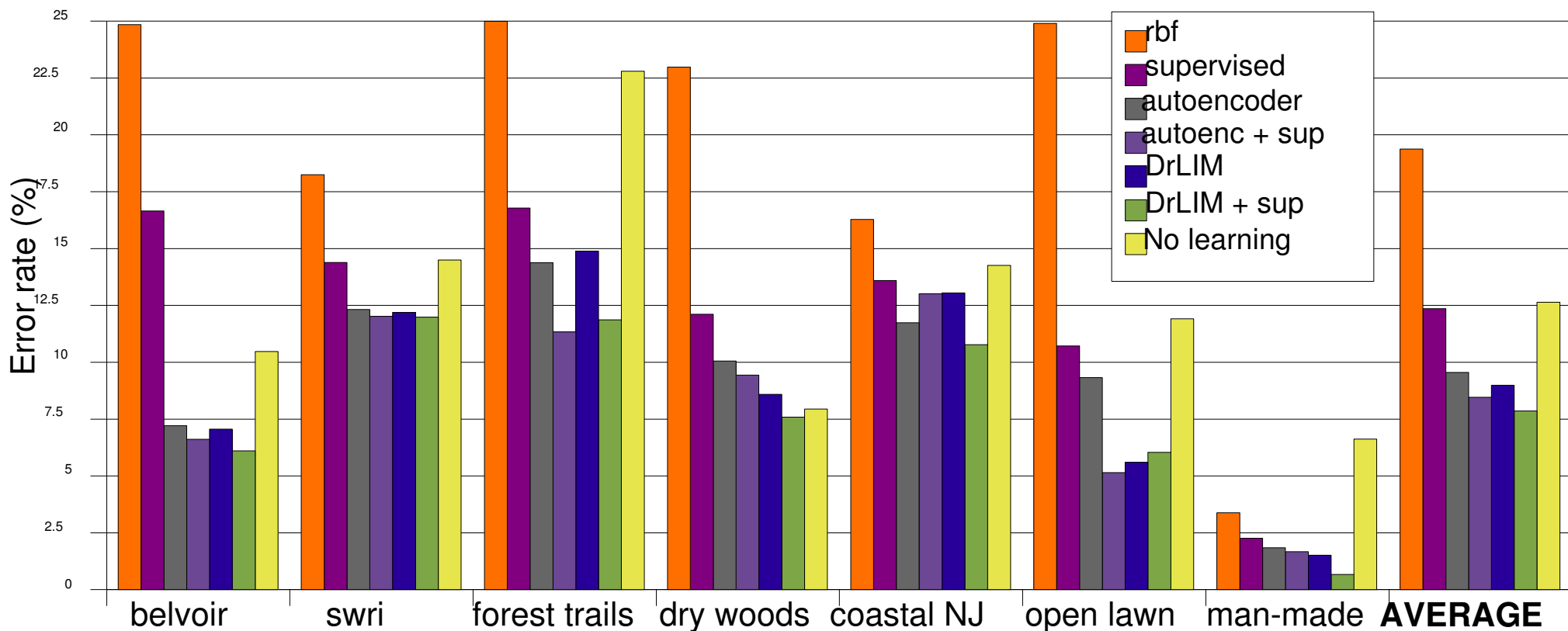# Feature Learning for traversability prediction (LAGR)

**Comparing**

- purely supervised
- stacked, invariant auto-encoders
- DrLIM invariant learning

**Testing on hand-labeled groundtruth frames – binary labels**

## Comparison of Feature Extractors on Groundtruth Data

# Collaborators

- **Current PhD students:**
  - Y-Lan Boureau, Koray Kavukcuoglu, Pierre Sermanet

- **Former PhD students:**
  - Raia Hadsell, Fu-Jie Huang, Marc'Aurelio Ranzato

- **Postdocs and Research Scientists**
  - Clément Farabet, Karol Gregor, Marco Scoffier

- **Senior Collaborators**
  - Rob Fergus (NYU): invariant feature learning
  - Eugenio Culurciello (Yale): FPGA/ASIC design
  - Yoshua Bengio (U. Montreal): deep learning
  - Leon Bottou (NEC Labs): handwriting recognition
  - Jean Ponce (ENS/INRIA), Francis Bach (ENS/INRIA): sparse coding.

New York University

**The End**

New York University