# An old Artificial Intelligence dream that comes true: Merging language and vision modalities
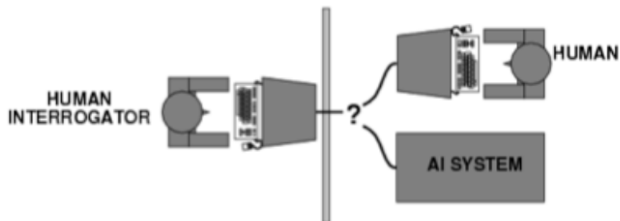
Raffaella Bernardi

University of Trento

November, 2017

# Credits

L. Fei Fei, Tamara Berg, Andrej Karpathy, Angeliki Lazaridou, Elia Bruni, Marco Baroni, Chris McCormick
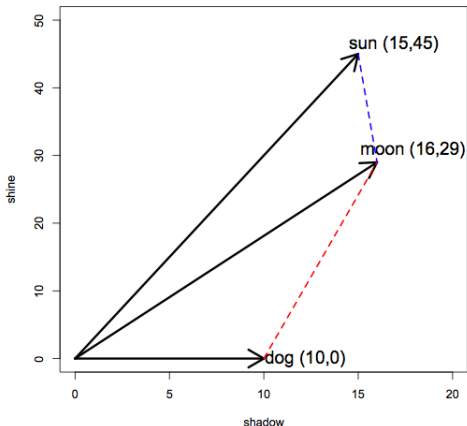
# An old AI dream

# Recall: From words to Meaning Representation

## Compute the word distribution

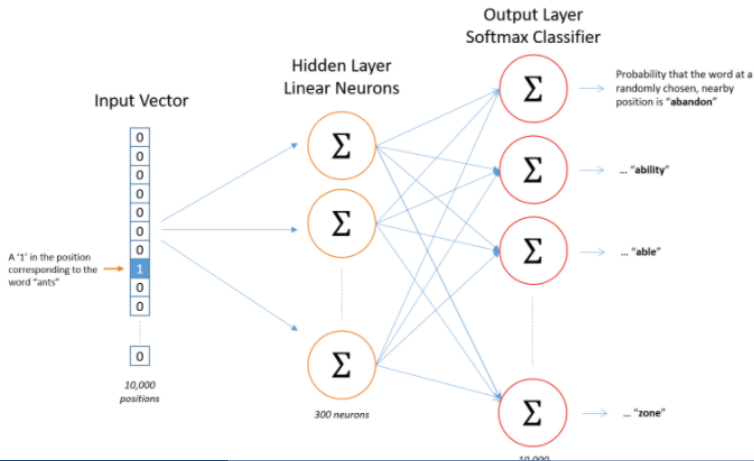Words can be represented by vectors harvested from a corpus of texts *counting* word co-occurences.

|      | shadow | shine |
|------|--------|-------|
| moon | 16     | 29    |
| sun  | 15     | 45    |
| dog  | 10     | 0     |

# From words to Meaning Representation
Predict the context: Word2Vec (Skip-Gram)

Instead counting words co-occurrences, the vector representing a word can be learned by *predicting* its nearby word. (Mikolov et al, 2013)
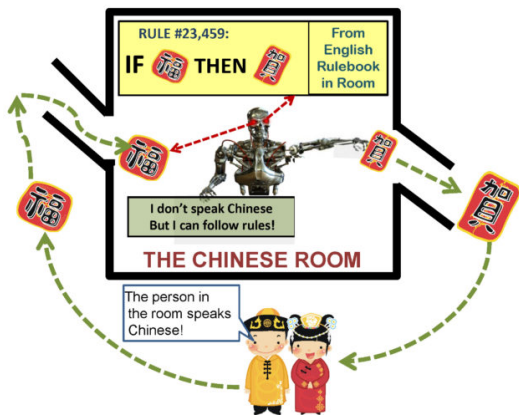
# Recall: Vector Representations
Successful

- Lexical meaning
    - Synonyms
    - Concept categorization (eg. car ISA vehicle)
    - Selectional preferences (e.g. eat chocolate vs. *eat sympathy)
    - relation classification (exam-anxiety CAUSE-EFFECT relation)
    - salient properties (car-wheels)
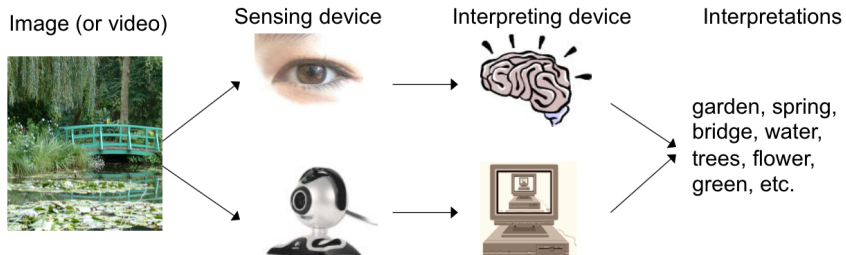- Compositionality: Phrase and Sentence
    - similairity
    - entailment

# Vector Representations
Grounding Problem



*Grounding* language representations into the world.
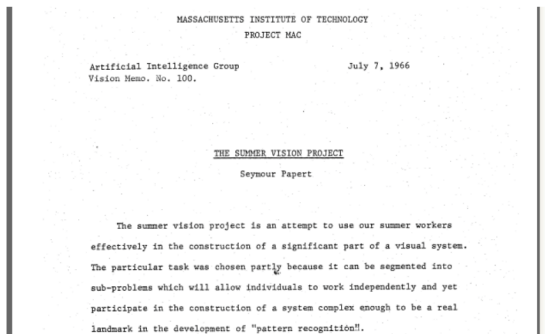Point to the *reference* of our mental representation.

# What is (Computer) Vision



Image (or video)    Sensing device    Interpreting device    Interpretations

garden, spring,
bridge, water,
trees, flower,
green, etc.

# How did it started?

## Origins of computer vision: an MIT undergraduate <u>summer project</u>



MASSACHUSETTS INSTITUTE OF TECHNOLOGY
PROJECT MAC

Artificial Intelligence Group              July 7, 1966
Vision Memo. No. 100.

THE SUMMER VISION PROJECT

Seymour Papert

The summer vision project is an attempt to use our summer workers
effectively in the construction of a significant part of a visual system.
The particular task was chosen partly because it can be segmented into
sub-problems which will allow individuals to work independently and yet
participate in the construction of a system complex enough to be a real
landmark in the development of "pattern recognition".

# From Pixels to Meaning Representation
## Gap

- To bridge the gap between pixels and "meaning"



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

What we see                    What a computer sees

Source: S. Narasimhan

# From Pixels to Meaning
Pixels

Raw images representation consists of pixels (a pixel is the minimum element of an image).

Pixels, identified by their physical coordinates, are stored as numbers encoding their color intensity. For instance, a black and white image is a 1-D representation of the pixel brightness);

A colored image is a 3-D arity of intensity values:

$$f(x, y) = \begin{bmatrix} red(x, y), \\ green(x, y), \\ blue(x, y) \end{bmatrix}$$

where color(x,y) is the intensity of that color (x) at position (y).

# How to represent an image: Keep all the pixels

# How to represent an image: Compute average pixel

# How to represent an image: Spatial grid of average pixel colors?



Photo by: marielito
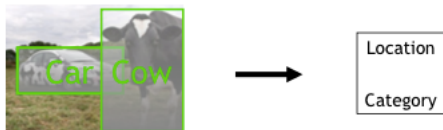
# From Pixels to Meaning Representation
## Challanges

# Applications: Traditional CV tasks
Objects

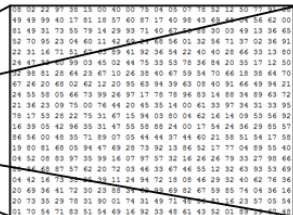**Image classification**: assigning a label to the image.



**Object localization**: define the location and the category.



Similarly, scene recognition.
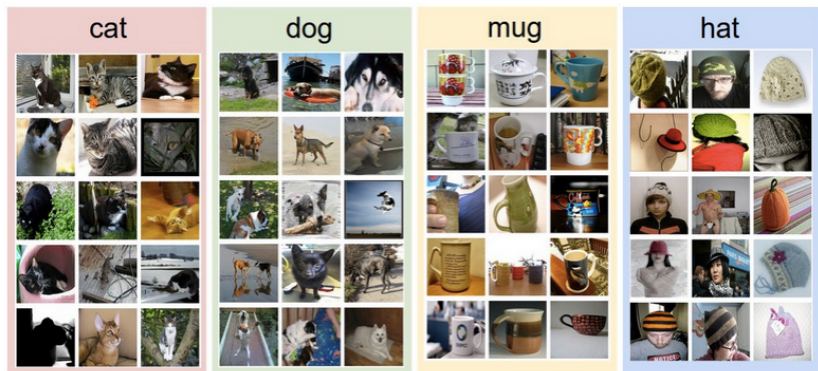
# Object Classification



What the computer sees

82% cat
15% dog
2% hat
1% mug

image classification

The task in Image Classification is to predict a single label (or a distribution over labels as shown here to indicates our confidence) for a given image. Images are 3-dimensional arrays of integers from 0 to 255, of size Width x Height x 3. The 3 is due to the three color channels Red, Green, Blue.

# Data Driven

**Data-driven approach**: it relies on first accumulating a training dataset of labeled images.



An example training set for four visual categories. In practice we may have thousands of categories and hundreds of thousands of images for each category.

## The image classification pipeline

- Input. Our input consists of a set of N images, each labeled with one of K different classes. We refer to this data as the *training set*.
- Learning. Our task is to use the training set to learn what every one of the classes looks like. We refer to this step as *training a classifier*, or learning a model.
- Evaluation. In the end, we evaluate the quality of the classifier by asking it to *predict labels* for a new set of images that it has never seen before (*test set*). We will then compare the true labels of these images to the ones predicted by the classifier. Intuitively, we're hoping that a lot of the predictions match up with the true answers (which we call the ground truth).

# Nearest Neighbor Classifier

The nearest neighbor (NN) classifier will

1. take a test image,
2. compare it to every single one of the training images, and
3. predict the label of the closest training image.

# Nearest Neighbor examples

In only about 3 out of 10 examples an image of the same class is retrieved, while in the other 7 examples this is not the case. For example, in the 8th row the nearest training image to the horse head is a red car, presumably due to the strong black background. As a result, this image of a horse would in this case be mislabeled as a car.

# Image distance

The difference (or the familiar cosine similarity)

| | test image | | | | | training image | | | | | pixel-wise absolute value differences | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 56 | 32 | 10 | 18 | | 10 | 20 | 24 | 17 | | | 46 | 12 | 14 | 1 | |
| 90 | 23 | 128 | 133 | − | 8 | 10 | 89 | 100 | = | | 82 | 13 | 39 | 33 | → 456 |
| 24 | 26 | 178 | 200 | | 12 | 16 | 178 | 170 | | | 12 | 10 | 0 | 30 | |
| 2 | 0 | 255 | 220 | | 4 | 32 | 233 | 112 | | | 2 | 32 | 22 | 108 | |

An example of using pixel-wise differences to compare two images with L1 distance (for one color channel in this example). Two images are subtracted elementwise and then all differences are added up to a single number. If two images are identical the result will be zero. But if the images are very different the result will be large.

# Evaluation

The CIFAR-10 training set of 50,000 images (5,000 images for every one of the labels), and we wish to label the remaining 10,000.
The NN Classifier based on raw pixel representation and the image distance measure above reaches 38.6 % accuracy (vs. upper bound: 94% – human).
State of the art classifier, Convolutional Neural Network reaches 95%.

# K-Nearest Neighbor Classifier

You may have noticed that it is strange to only use the label of the nearest image when we wish to make a prediction. Indeed, it is almost always the case that one can do better by using what's called a k-Nearest Neighbor Classifier. The idea is very simple: instead of finding the single closest image in the training set, we will find the top k closest images, and have them vote on the label of the test image.
Which is the best k?
K is an hyperparameter. There are others too.

# Validation dataset vs Test dataset

They can be trained/learned.

- Training data-set: to train the classifier.
- Trail data-set (or development dataset or validatation dataset): to tune the parameters.
- Test data-set. To test the classifier.

# First problem: the Raw Pixel representation

- Using the NN classifier over the raw pixel representation images that are nearby each other are much more a function of the general color distribution of the images, or the type of background rather than their semantic identity.
- For example, a dog can be seen very near a frog since both happen to be on white background.
- Ideally we would like images of all of the 10 classes to form their own clusters, so that images of the same class are nearby to each other regardless of irrelevant characteristics and variations (such as the background).

To get this property we will have to go beyond raw pixels.

# From Pixels to Meaning

## Abstract Features

# Second problem: the classifier

NN Classifier: pro and contra.

- the classifier takes no time to train, since all that is required is to store and possibly index the training data.
- However, we pay that computational cost at test time, since classifying a test example requires a comparison to every single training example.
- This is backwards, since in practice we often care about the test time efficiency much more than the efficiency at training time.

In CV it's better to use other classifiers.
State of the art: Deep neural networks are very expensive to train, but once the training is finished it is very cheap to classify a new test example. This mode of operation is much more desirable in practice.

# First Revolution: Big dataset
ImageNet

Image database organized according to the WordNet hierarchy.
Stanford Vision Lab, Stanford University & Princeton University.

- Challenges: 2007-present
- AMT: 48,940 annotators from 167 countries
- 15M images
- 22K categories of objects

# From Pixels to Features
Two methods

- Bag of Visual words (BoVW) (Sivic and Zisserman, 2003)
- Convolutional neural network (CNN) (LeCun et al., 1998, Krizhevsky et ali. 2012)

# From Pixels to Features: BoVW
Pipeline

**Keypoints detectors** To locate interesting points/content, various kinds of low-level features detectors exists:

- edge detection: the lines we would draw – encode shape info
- corner detection

**Local description** The identified interesting points are then described: clustered into regions and transformed into *vectors representing the region*. Several local descriptors exist, e.g:

- SIFT: Scale-invariante feature transform (Lowe '99) – edge based features.
- Textons (Leung and Malik '01)
- HoG (Histograms of Oriented Gradients) (Dalal and Triggs '05)

The low-level features can capture eg. Color, Texture, Shape,

**Bag of Visual Words** The local descriptions are clustered to obtain the Visual Words that are used to obtain the vector representation of the image.

# From Pixels to Features: BoVW
## BoVW's pipeline

# Second Revolution: End-to-end systems
Convolutional Neural Networks

ImageNet Classification with Deep Convolutional Neural Networks Alex Krizhevsky, Ilya Sutskever and Geoffrey E. Hinton, 2012
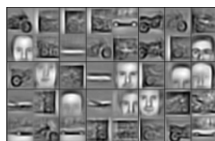
- 2012: Krizhevsky outperformed the other systems using CNN
- 2013: half of the systems used CNN
- 2014: All of the systems used CNN.
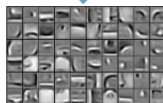
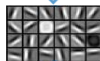# End-to-end Systems
Hierarchy of features

## Deep Learning

– Deep architectures can be representationally efficient.

– Natural progression from low level to high level structures.

– Can share the lower-level representations for multiple tasks.
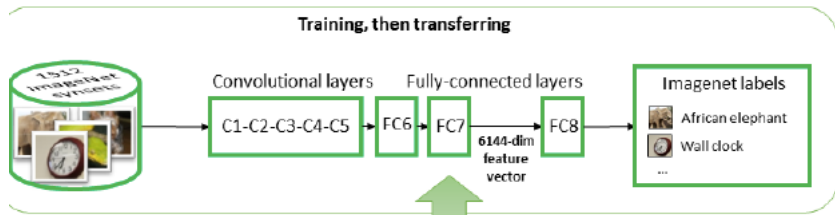


3rd layer
"Objects"

2nd layer
"Object parts"

1st layer
"edges"

Input

# End-to-end systems
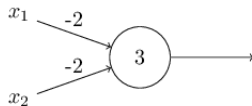CNN: off-the-shelf vector representation



- Train a CNN on a vision task (e.g. AlexNet on ImageNet)
- Do a forward pass given an image input
- Transfer one or more layers (e.g. FC7 or C5)

# Neural Networks
Example to compute a logical operator "Not And"

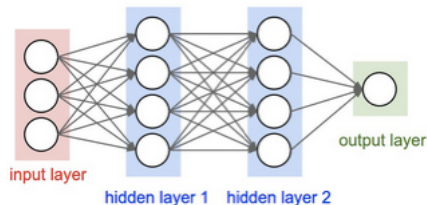| $x_1$ | $x_2$ | "Not and" |
|-------|-------|-----------|
| 0 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |



Input 00 produces output 1 (since: (-2)*0 + (-2) * 0 + 3 = 3 is positive) and similarly, 01 and 10; but the input 11 produces output 0 (since: (-2)*1 + (-2) * 1 + 3 = -1 is negative.)

# Neural Networks
Neural Networks

It's a composition of functions (neurons) that goes from an n-dimensional vector to class scores.



Each neuron receives some inputs, performs a dot product and optionally follows it with a non-linearity. On the last (fully-connected) layer, they have a loss function (e.g., Softmax).

# Neural Networks
Recurrent NN: intuitions

Traditional neural networks cannot use the information "about previous inputs" to inform later ones.

- *Recurrent neural networks* (RNNs) address this issue: They are networks with loops in them, allowing information to persist. They work well with short dependencies.
- *Long Short Term Memory* (LSTM) are a special kind of RNN, capable of learning long-term dependencies.
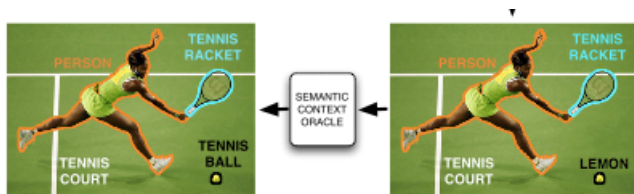
# Language and Vision

Language and Visual Space can be combined!

# Applications: Traditional CV tasks
Corpora as KB source: Object recognition

A. Rabinovich, A. Vedaldi, C. Galleguillos, E. Wiewiora, S. Belongie (ICCV 2007) Objects in Context.
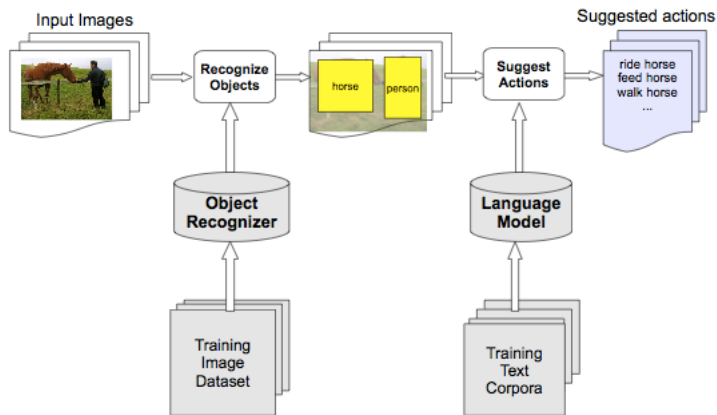


Not a Lemon, it's more probable a Tennis Ball. Info come from a KB (word similarity list, exctracted from internet – Google Sets).
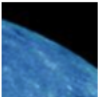
# Applications: Traditional CV tasks

Corpora as KB source: Action recognition

Thu Le Dieu, Jasper Uijlings and R. Bernardi (2010, 2011)

# Applications: Traditional NLP tasks

E. Bruni, G.B. Tran and M. Baroni (GEMS 2011, ACL 2012, Journal of AI 2014), E. Bruni, G. Boleda, M. Baroni and N. Tran (ACL 2012)

|        | planet | night |  |  |
|--------|--------|-------|--------|--------|
| moon   | 10     | 22    | 22     | 0      |
| sun    | 14     | 10    | 15     | 0      |
| dog    | 0      | 4     | 0      | 20     |

# Applications: Traditional NLP tasks

**Task 1** Predicting human **semantic relatedness** judgments

Improved!

**Task 2** **Concept categorization**, i.e. grouping words into classes based on their semantic relatedness (*car* ISA *vehicle*; *banana* ISA *fruit*)

Improved!

**Task 3** Find **typical color** of concrete objects (**cardboard is brown**, **tomato is red** )

Improved!

**Task 4** Distinguish **literal vs. non-literal** usages of color adjectives (**blue uniform** vs. **blue note**)

Improved!