

Deep Learning for Visual and Multimodal Recognition
Séminaire TIM 2017
Traitement de l'Information Multimodale

Nicolas Thome

Conservatoire National des Arts et Métiers (Cnam)

Équipe MSDMA - Département Informatique

Prenom.Nom@cnam.fr

<http://cedric.cnam.fr/~thomen/>

6 Juillet 2017

Context

Big Data: Images & videos everywhere



BBC: 2.4M videos

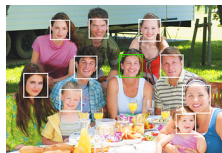


Social media,
e.g. Facebook: 1B each day



100M monitoring cameras

- Obvious need to access, search, or classify these data: **Visual Recognition**
- Huge number of applications: mobile visual search, robotics, autonomous driving, augmented reality, medical imaging etc
- Leading track in major ML/CV conferences during the last decade



Outline

- 1 Deep Learning & Model Architectures
- 2 Applications of Deep Learning for Visual Recognition
- 3 Open Issues & Perspectives in Artificial Intelligence

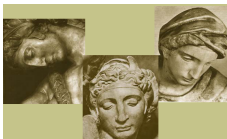
Recognition of low-level signals

Challenge: filling the semantic gap



What we perceive vs
 What a computer sees

141	129	240	222	206	202	186	218	211	208	216	212
142	208	218	110	497	91	88	162	216	208	208	201
143	242	122	58	94	82	132	77	100	100	100	112
135	217	119	212	242	216	247	139	91	209	208	211
138	108	181	221	218	218	186	114	74	208	218	214
132	217	181	114	77	188	89	96	82	201	208	181
131	132	182	184	184	179	159	113	91	212	216	131
162	198	201	184	218	181	129	81	178	202	241	140
135	208	220	128	172	126	91	63	124	149	241	242
137	226	247	143	55	55	10	94	155	248	247	251
134	237	240	181	55	51	112	144	112	241	241	251
140	148	181	128	148	109	130	95	47	168	139	181
180	167	38	162	94	73	114	16	17	7	51	137
13	81	81	148	148	203	179	43	27	17	12	8
17	16	12	163	156	235	189	12	16	19	16	14



- Illumination variations
- View-point variations
- Deformable objects
- intra-class variance
- etc

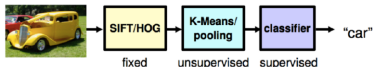
⇒ How to design "good" intermediate representation ?

Deep Learning (DL) & Recognition of low-level signals

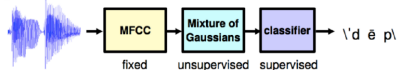
- Before DL: **handcrafted intermediate representations** for each task

- ⊖ Needs expertise in each field
- ⊖ Shallow: low-level features

VISION



SPEECH

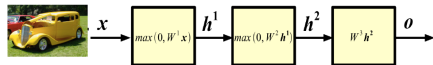


Credit: I. Kokkinos

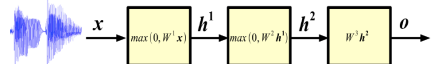
- Since DL: **Representation Learning**

- ⊕ Deep: hierarchy, gradually learning higher-level representations
- ⊕ Common learning methodology ⇒ field independent, no expertise

VISION



SPEECH

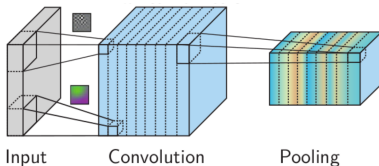


Credit: I. Kokkinos

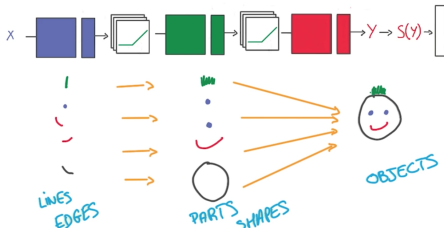
- All parameters trained with backpropagation with class labels

Convolutional Neural Networks (ConvNets)

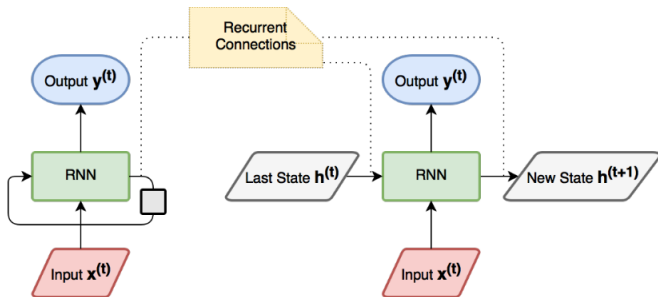
- Convolution on tensors, *i.e.* multidimensional arrays: T of size $W \times H \times D$
 - Convolution: $C[T] = T'$, T' tensor of size $W' \times H' \times K$
 - Each filter locally connected with shared weights (K number of filters)



- **An elementary block: Convolution + Non linearity (e.g. ReLU) + pooling**
- **Convolution:** structure (local processing), **Pooling:** invariance
- **Stacking several Blocks:** intuitive hierarchical information extraction



Recurrent Neural Networks (RNNs) for Sequence Modeling



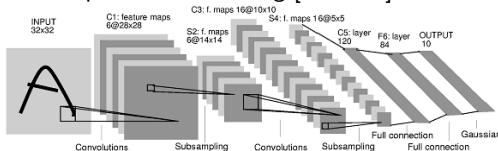
- **Sequences:** 1d/2d signals (e.g. audio, videos), molecules, text, etc
- Input vector $x^{(t)}$, e.g. word (text) or image representation (CNN)
- Input/Output $h^{(t)}$: vector representing model "short-term memory"
- Output vector $y^{(t)}$: task dependent
- All parameters trained with backpropagation through time

Outline

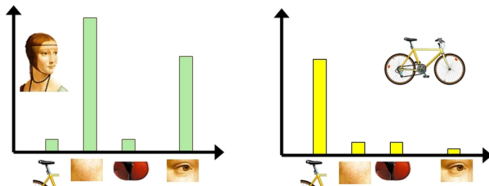
- 1 Deep Learning & Model Architectures
- 2 Applications of Deep Learning for Visual Recognition
- 3 Open Issues & Perspectives in Artificial Intelligence

Visual Recognition History: Trends and methods in the last four decades

- 80's: training Convolutional Neural Networks (CNN) with back-propagation \Rightarrow postal code reading [LBD⁺89]

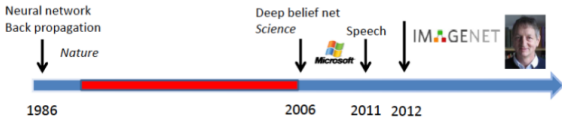


- 90's: golden age of kernel methods, NN = black box
- 2000's: BoW + SVM : state-of-the-art CV



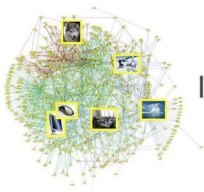
Visual Recognition History: Trends and methods in the last four decades

- Deep learning revival in 2012: CNN success of ConvNets in ImageNet [KSH12]



Rank	Name	Error rate	Description
1	U. Toronto	0.15315	Deep learning
2	U. Tokyo	0.26172	Hand-crafted
3	U. Oxford	0.26979	features and learning models.
4	Xerox/INRIA	0.27058	Bottleneck.

- Two main practical reasons:
 - 1 Huge number of labeled images (10^6 images)
 - 2 GPU implementation for training

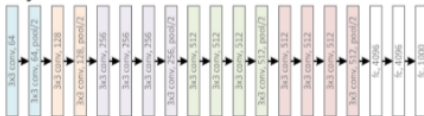


IMAGENET



Deep Learning since 2012: Larger & larger networks

VGG, 16/19 layers, 2014



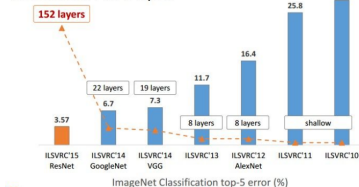
GoogleNet, 22 layers, 2014



ResNet, 152 layers, 2015



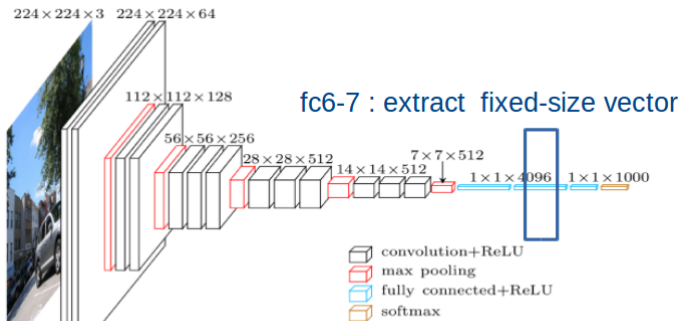
Revolution of Depth



Deep Learning since 2012

Transferring Representations learned from ImageNet

- Deep ConvNets require large-scale annotated datasets
- **BUT:** Extract layer \Rightarrow fixed-size vector: "Deep Features" (DF)



- Now state-of-the-art for any visual recognition task [ARS⁺16]
- Ex: Domain adaptation for image classif:
 - DF very robust to data variations, e.g. medical / astronomy images

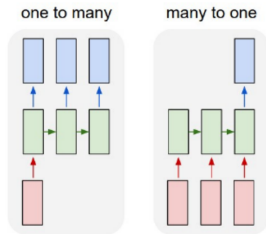
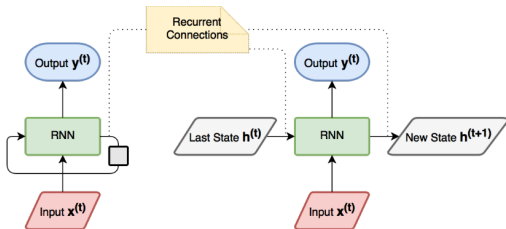
Outline

- 1 Deep Learning & Model Architectures
- 2 Applications of Deep Learning for Visual Recognition
- 3 Open Issues & Perspectives in Artificial Intelligence

Ongoing Issues in Deep Learning

New Tasks in Artificial Intelligence

- Intersection of vision and language research:
 - Improvements in Vision Understanding with ConvNets
 - Language (text) Modeling with RNNs



- One to Many: Image captioning
- Many to One: Visual Question Answering



Many to One - Visual Question Answering (VQA)

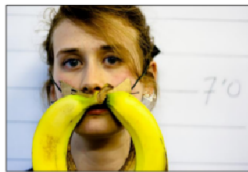
- Goal: build a system that can answer questions about images



How many slices of pizza are there?
Is this a vegetarian pizza?



Does it appear to be rainy?
Does this person have 20/20 vision?



What color are her eyes?
What is the mustache made of?

- Very complex task, that requires :
 - Precise image and text models
 - High level interaction modeling
 - Full scene understanding
 - Reasoning (e.g. spatial ...)



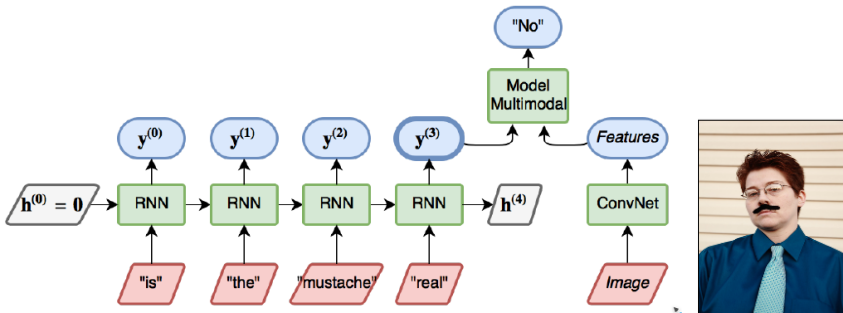
What color is the fire hydrant
on the right ? **yellow**



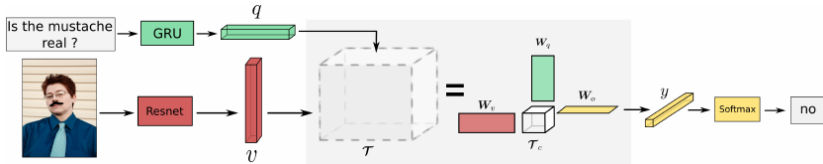
What color is the fire hydrant
on the left ? **green**

Many to One - Visual Question Answering (VQA)

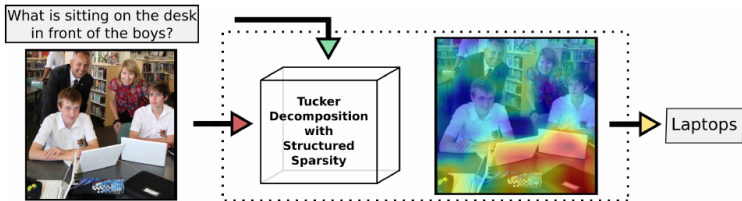
- Input: question & image
- Output: answer



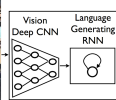
Multimodal Tucker Fusion for Visual Question Answering (MUTAN)



- State-of-the-art mono-modal representations:
 - Visual representation: ResNet-152
 - Question representation: pre-trained GRU (Gated Recurrent Units)
- How to perform multi-modal fusion ?
 - State-of-the-art: bilinear models [FPY⁺16, KOK⁺17] \Rightarrow accurate interactions
 - BUT full bilinear models intractable: factorization based on Tucker decomposition [BCCT17]

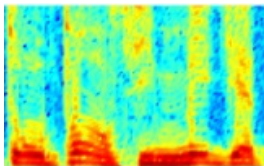


Conclusion

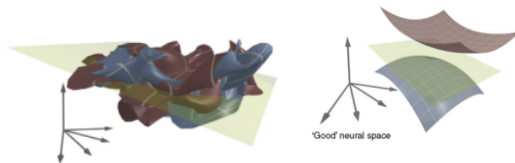


A group of people shopping at an outdoor market.

There are many vegetables at the fruit stand.



- Deep Learning: huge impact in terms of experimental results
- Still a long way to go toward real AI: reasoning, memory, predictive models, common knowledge, etc
- Formal understanding still limited:
 - Optimization: non-convex problem
 - Model: ability to untangle manifold
 - Robustness to over-fitting & generalization



References I



Hossein Azizpour, Ali Sharif Razavian, Josephine Sullivan, Atsuto Maki, and Stefan Carlsson, *Factors of transferability for a generic convnet representation*, IEEE Trans. Pattern Anal. Mach. Intell. 38 (2016), no. 9, 1790–1802.



Hedi Ben-younes, Rémi Cadène, Matthieu Cord, and Nicolas Thome, *MUTAN: multimodal tucker fusion for visual question answering*, CoRR abs/1705.06676 (2017).



Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach, *Multimodal compact bilinear pooling for visual question answering and visual grounding*, arXiv:1606.01847 (2016).



Jin-Hwa Kim, Kyoung-Woon On, Jeonghee Kim, Jung-Woo Ha, and Byoung-Tak Zhang, *Hadamard Product for Low-rank Bilinear Pooling*, 5th International Conference on Learning Representations, 2017.



Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, *Imagenet classification with deep convolutional neural networks*, Advances in neural information processing systems, 2012, pp. 1097–1105.



Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel, *Backpropagation applied to handwritten zip code recognition*, Neural computation 1 (1989), no. 4, 541–551.