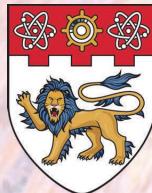


# *Convolutional Neural Networks on Graphs*

Xavier Bresson

School of Computer Science and Engineering  
Nanyang Technological University



**NANYANG  
TECHNOLOGICAL  
UNIVERSITY  
SINGAPORE**

**NATIONAL RESEARCH FOUNDATION**  
PRIME MINISTER'S OFFICE  
SINGAPORE

Joint work with M. Defferrard (EPFL), P. Vandergheynst (EPFL),  
Y. Seo (EPFL), M. Bronstein (USI), F. Monti (USI), T. Laurent (LMU)

Institute for Mathematical Sciences  
National University of Singapore  
June 1<sup>st</sup> 2017

Xavier Bresson

*Note: Background  
image is generated by  
an artificial neural  
network with random  
input, i.e. the AI is  
dreaming... 1*

# Outline

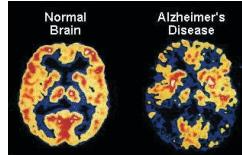
- Data Science and Deep Learning
  - *A brief history*
  - *What makes DL great*
- Convolutional Neural Networks
  - *Architecture*
- Convolutional Neural Networks on Graphs
  - *Non-Euclidean data*
  - *Graph convolution, graph downsampling and pooling*
  - *Numerical experiments*
- Natural Language Processing
- Recommender Systems
- Conclusion

# Outline

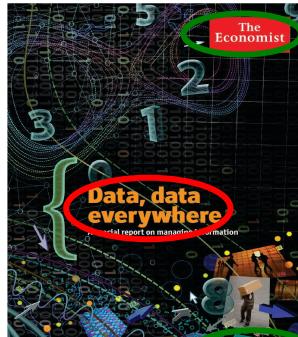
- Data Science and Deep Learning
  - *A brief history*
  - *What makes DL great*
- Convolutional Neural Networks
  - *Architecture*
- Convolutional Neural Networks on Graphs
  - *Non-Euclidean data*
  - *Graph convolution, graph downsampling and pooling*
  - *Numerical experiments*
- Natural Language Processing
- Recommender Systems
- Conclusion

# Data Science and Artificial Neural Networks

- Data Science is an *interdisciplinary* field that aims at *transforming raw data into meaningful knowledge* to provide smart decisions for real-world problems:



- In the news:



McKinsey Global Institute

Big data: The next frontier for innovation, competition, and productivity



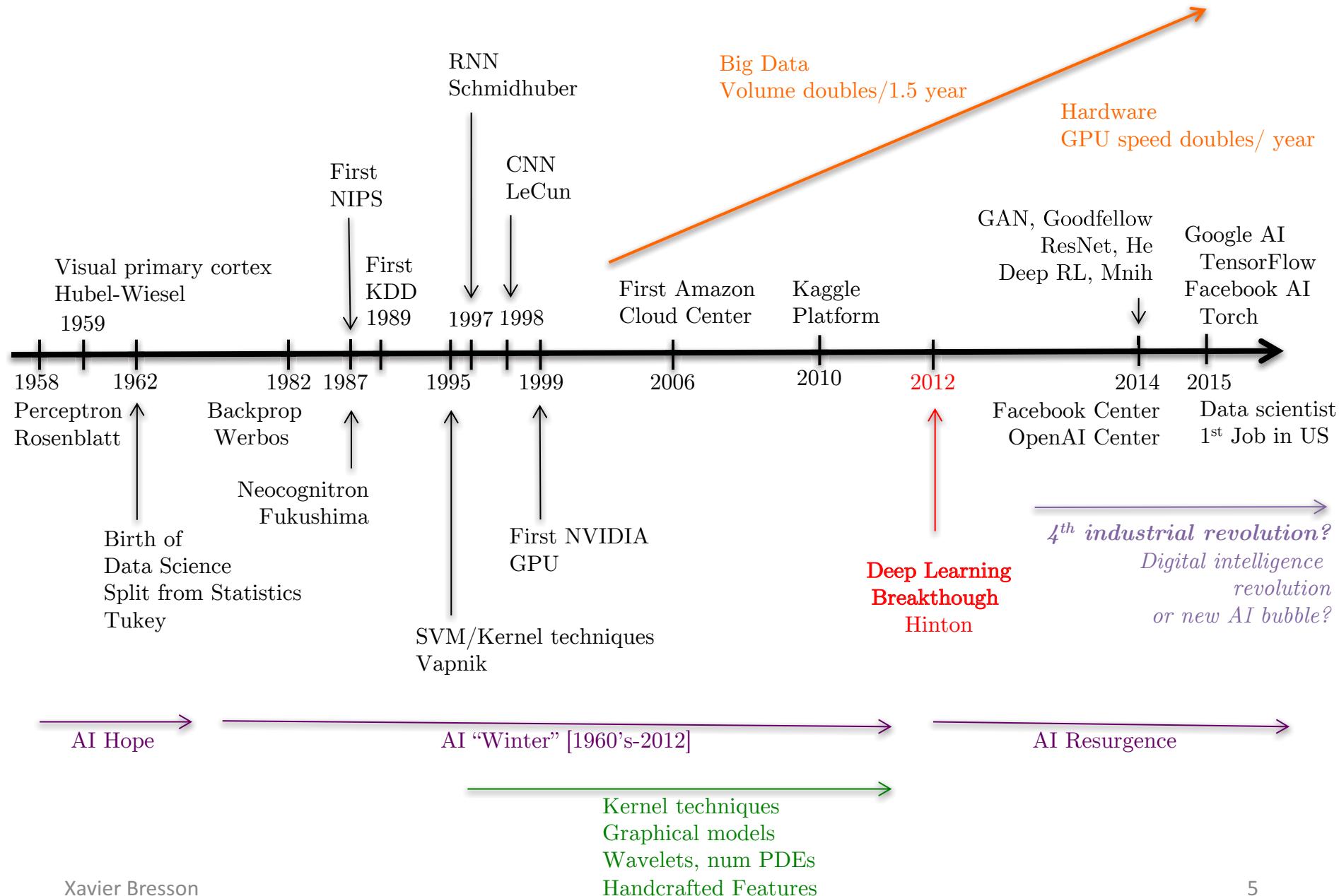
Data Scientist: The Sexiest Job of the 21st Century



South Korea trumpets \$860-million AI fund after AlphaGo 'shock'



# A Brief History of Artificial Neural Networks

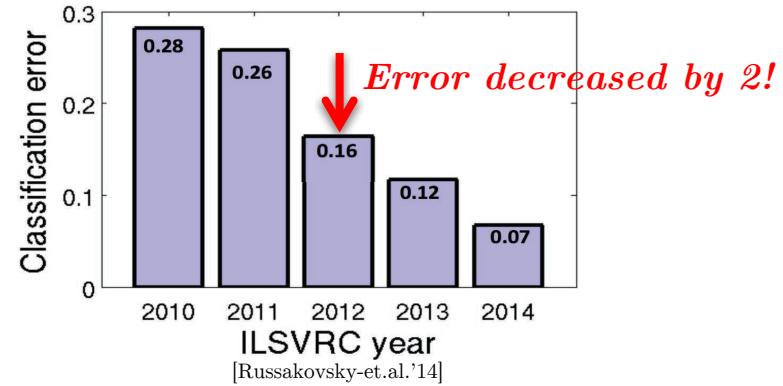


# 2012: Deep Learning Breakthrough

- ImageNet [Fei Fei et.al.'09]: International Image Classification Challenge  
1,000 object classes and 1,431,167 images

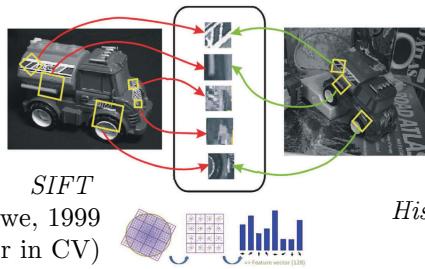


[L. Fei-Fei et.al.'09]

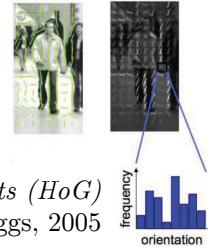


- Before 2012: Handcrafted filters

[Fei-Fei Li et.al.'15]

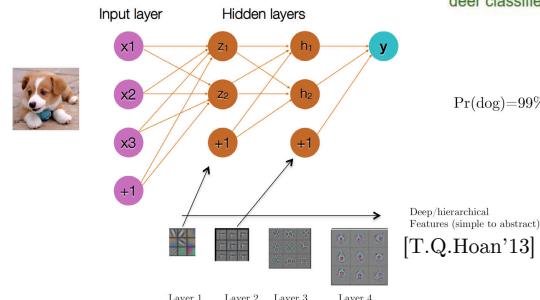
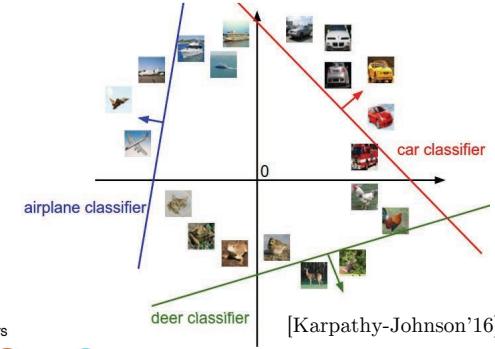


Histogram of Gradients (HoG)  
Dalal & Triggs, 2005

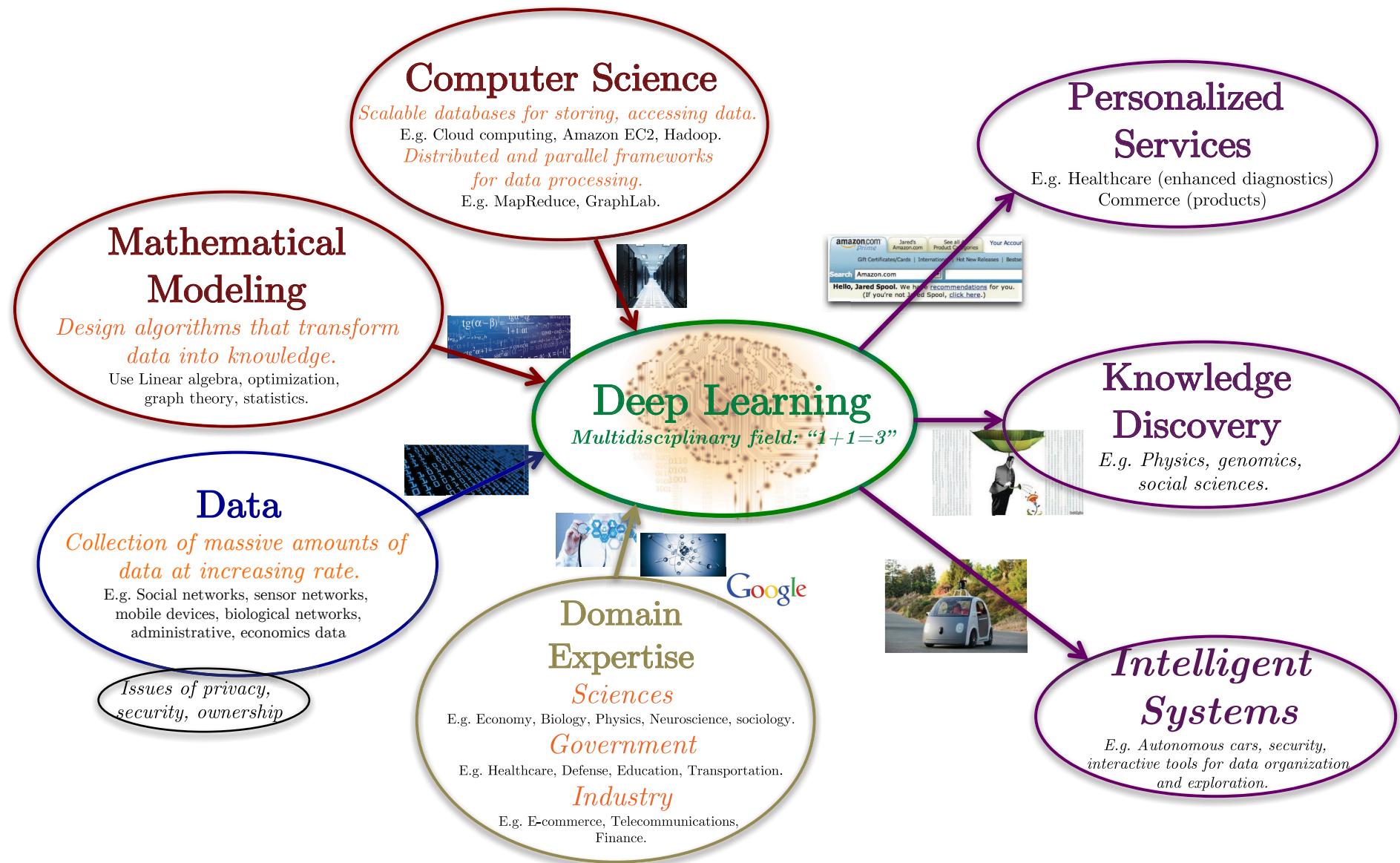


After 2012: Learn filters and classifier with neural networks (end-to-end system).

+ SVM classification



# What Makes Deep Learning Great? [long answer]



**Major challenges:** Multidisciplinary integration, large-scale databases, scalable computational infrastructures, design math algorithms for massive datasets, trade-off speed and accuracy for real-time decisions, interactive visualization tools.

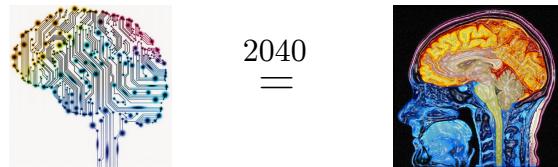
# What Makes Deep Learning Great? [short answer]

1. Data.. lots of data: Inherited from 3<sup>rd</sup> industrial revolution (digital revolution)



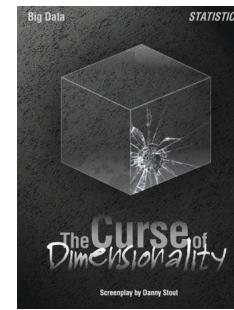
2. Computational power: Moore's law (CPUs, GPUs) and cloud computing

Note: Human brain  $10^{14}$  synapses/ $10^9$  Mflops =  
Computers in 2040 (10,000 times faster than today)



3. Neural networks can beat the curse of dimensionality  
(more next slide).

4. Neural network is a universal/flexible strategy to solve many challenging real-world problems:  
DL/AI is the new “electricity” [A. Ng, Baidu], it will change major industries in the coming decade.



# Curse of Dimensionality

- Data lie in high-dim spaces:

$$\dim(\text{Go}) = 19 \times 19 \approx 10^3$$

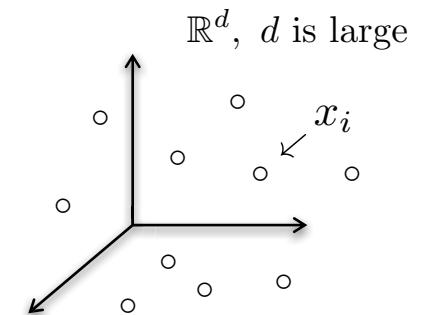
$$\dim(\text{Images}) = 512 \times 512 \approx 10^6$$

$$\dim(\text{Books}) = 80,000 \approx 10^5$$

- Consequence 1: Curse of dimensionality [Beyer'98]:

*In high dimensions, distances between data is meaningless  $\Rightarrow$  all data are close to each other.*

$$\lim_{d \rightarrow \infty} \mathbb{E} \left( \frac{d_{\max}^{\ell_2}(x_i, V - x_i) - d_{\min}^{\ell_2}(x_i, V - x_i)}{d_{\min}^{\ell_2}(x_i, V - x_i)} \right) \rightarrow 0$$



- Consequence 2: All data are intractable.

$$\#\text{Go} = 3^{19 \times 19} \approx 10^{172}$$

$$\#\text{Images} = 256^{512 \times 512} \approx 10^{630,000}$$

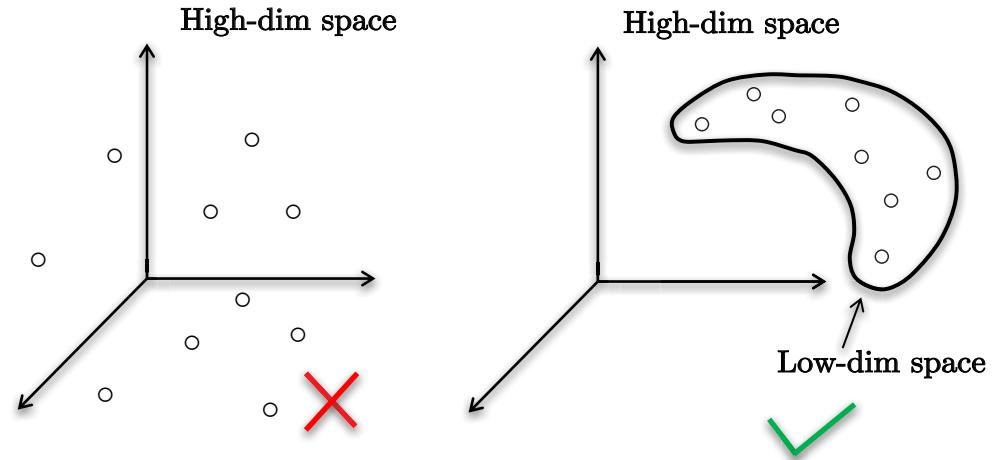
$$\#\text{Books} = 8,000^{80,000} \approx 10^{310,000}$$

$\Rightarrow$  Can we learn exponential functions  $f(x)$ ,  $\#x=10^{172}$ ? No, in general.

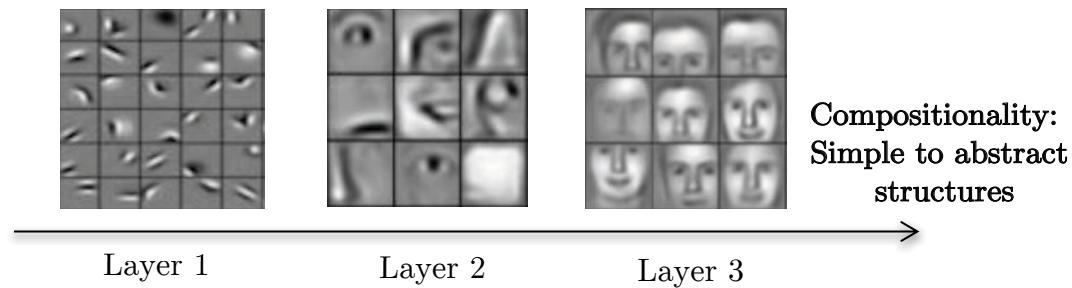
# Good News

- In our universe:
  1. Data have structures.
  2. These structures are compositional.

- **Blessing of structure:** Data are not random, they have structures, meaning data are concentrated in low-dim spaces

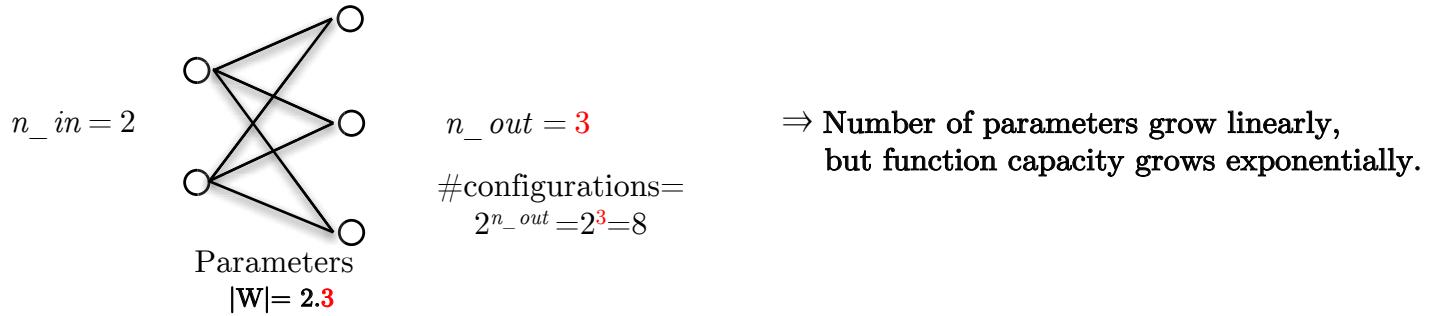


- **Compositionality:** Data are formed of simple structures that agglomerate to compose more complex and abstract structures and so on.



# Learning Exponential Functions

- Neural networks leverage structureness and compositionality by learning **parametric** exponential functions with **linear** number of parameters.



- Another good news: Neural networks can learn **a very large number of parameters** in a consistent way (w/ backpropagation), order of  $10^9 \Rightarrow$  huge capacity to learn (and to generalize).
- Another good news: **Theorem of Universal Approximation [Cybenko'89]** is today obsolete.

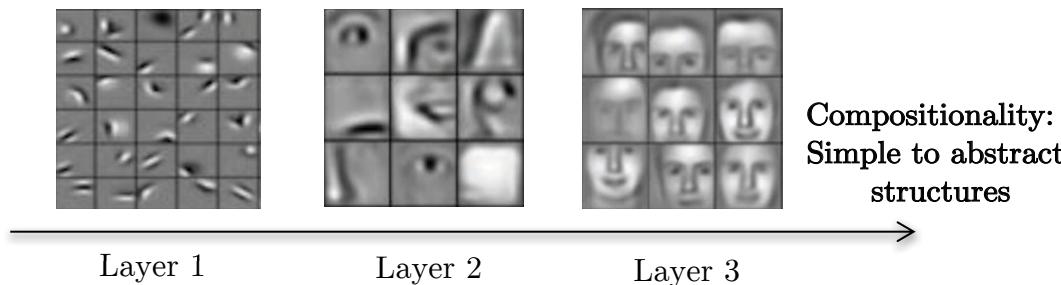
Theorem states that a neural network with a single hidden layer can approximate nearly every function, but exponential functions of interest require an exponential number of neurons with a single layer.

*Not true with multiple layers ⇒ deep property [Bengio-et.al'14].*

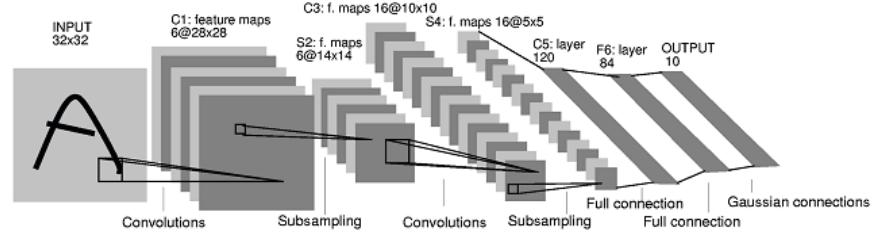
# Outline

- Data Science and Deep Learning
  - *A brief history*
  - *What makes DL great*
- Convolutional Neural Networks
  - *Architecture*
- Convolutional Neural Networks on Graphs
  - *Non-Euclidean data*
  - *Graph convolution, graph downsampling and pooling*
  - *Numerical experiments*
- Natural Language Processing
- Recommender Systems
- Conclusion

# How Compositionality is Captured in Images?

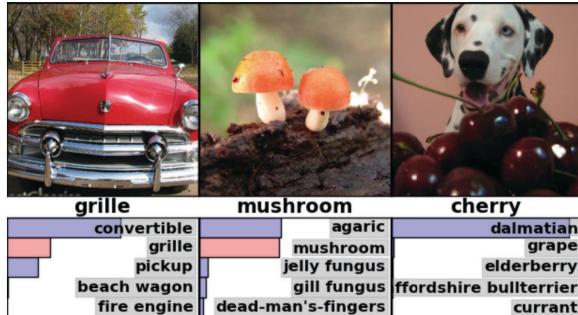


- Compositionality is captured by learning **local stationary structures** and compose them into **multiscale hierarchical patterns** (explained next).
- This class of deep learning is called **convolutional neural networks (CNNs)** [LeCun-Bottou-Bengio-Haffner'98].
- CNNs are **extremely efficient** at extracting meaningful statistical patterns in large-scale and high-dimensional image datasets  
⇒ CNN has become the Swiss knife of Computer Vision.



# CNNs are Great

- CNNs are highly successful for computer vision tasks:



*Image/object recognition*

[Krizhevsky-Sutskever-Hinton'12]



*Image inpainting*

[Pathak-Efros-etal'16]



*Image captioning*

[Karpathy-FeiFei'15]



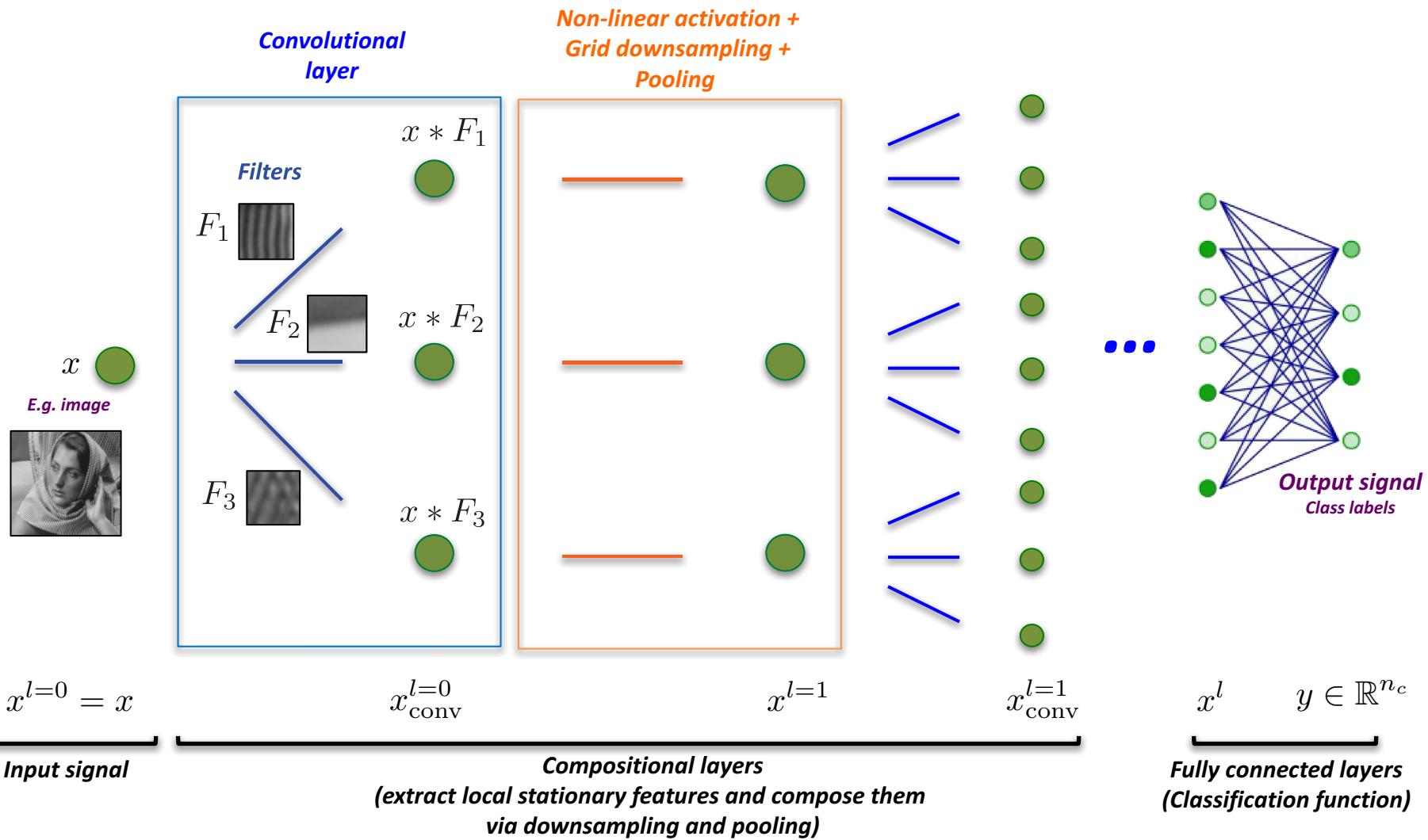
*Self-driving cars*

⇒ CNNs are best implemented by:

1. Google TensorFlow
2. Facebook Torch

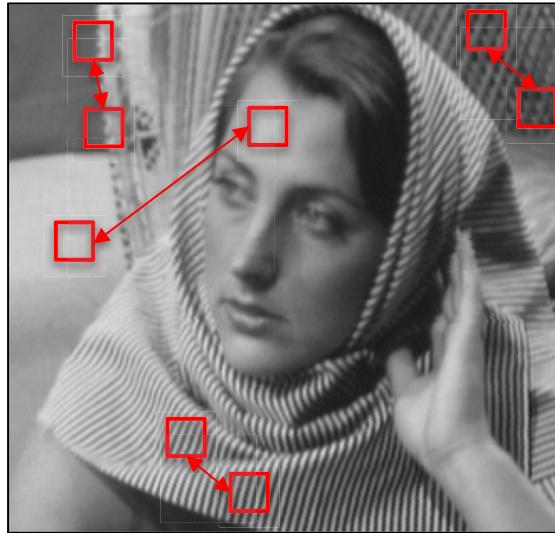


# Architecture of CNNs



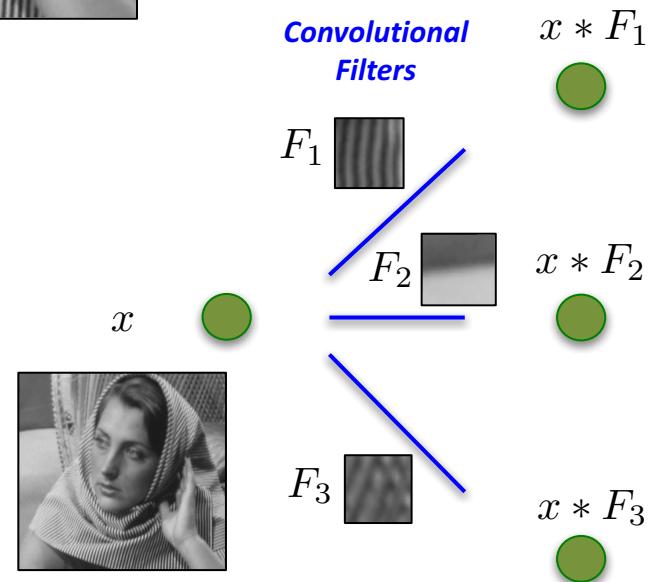
# Convolutional Layer

- **Assumption:** Images are *locally stationary*  $\Leftrightarrow$  similar local patches are shared across the data domain:



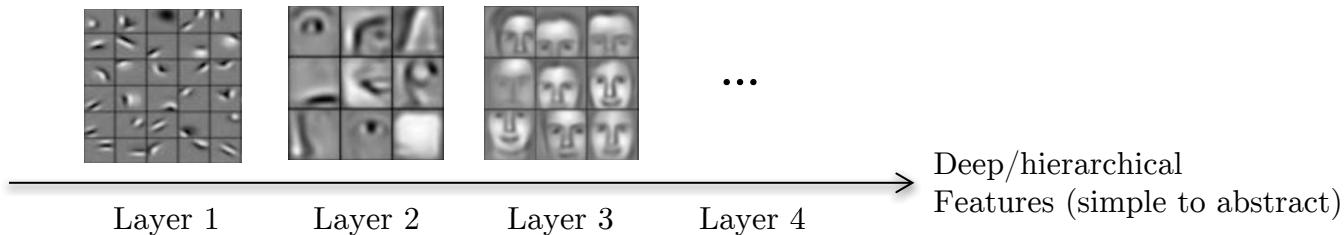
- How to extract local stationary patterns?

*Convolutional filters* (filters with compact support)



# Compositional Layers

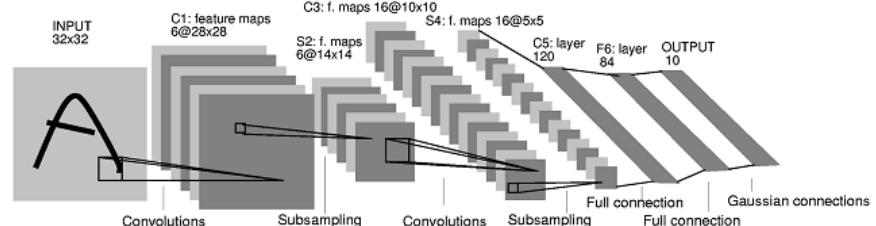
- **Assumption:** Images are *compositional*  $\Leftrightarrow$  local stationary patterns can be composed to form more abstract complex patterns:



- **How to extract these hierarchical patterns?**  
*Downsampling* of data domain (image grid) with *pooling* (max, average).



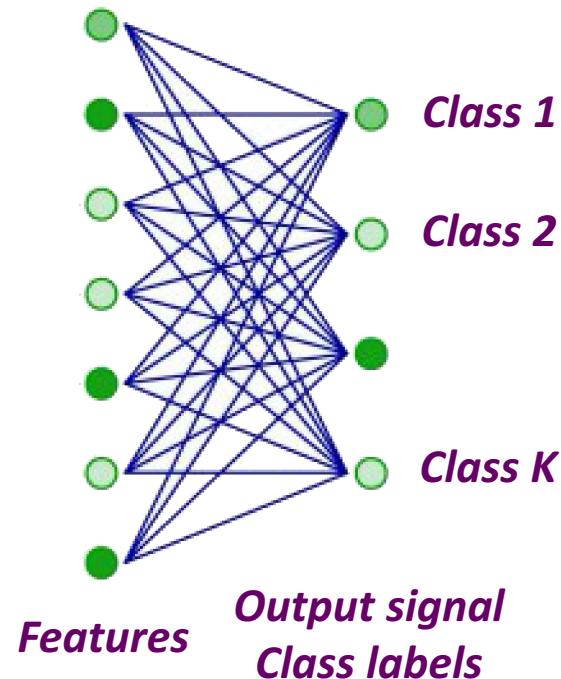
- **Other advantage:** Reduce computational complexity while increasing #filters.



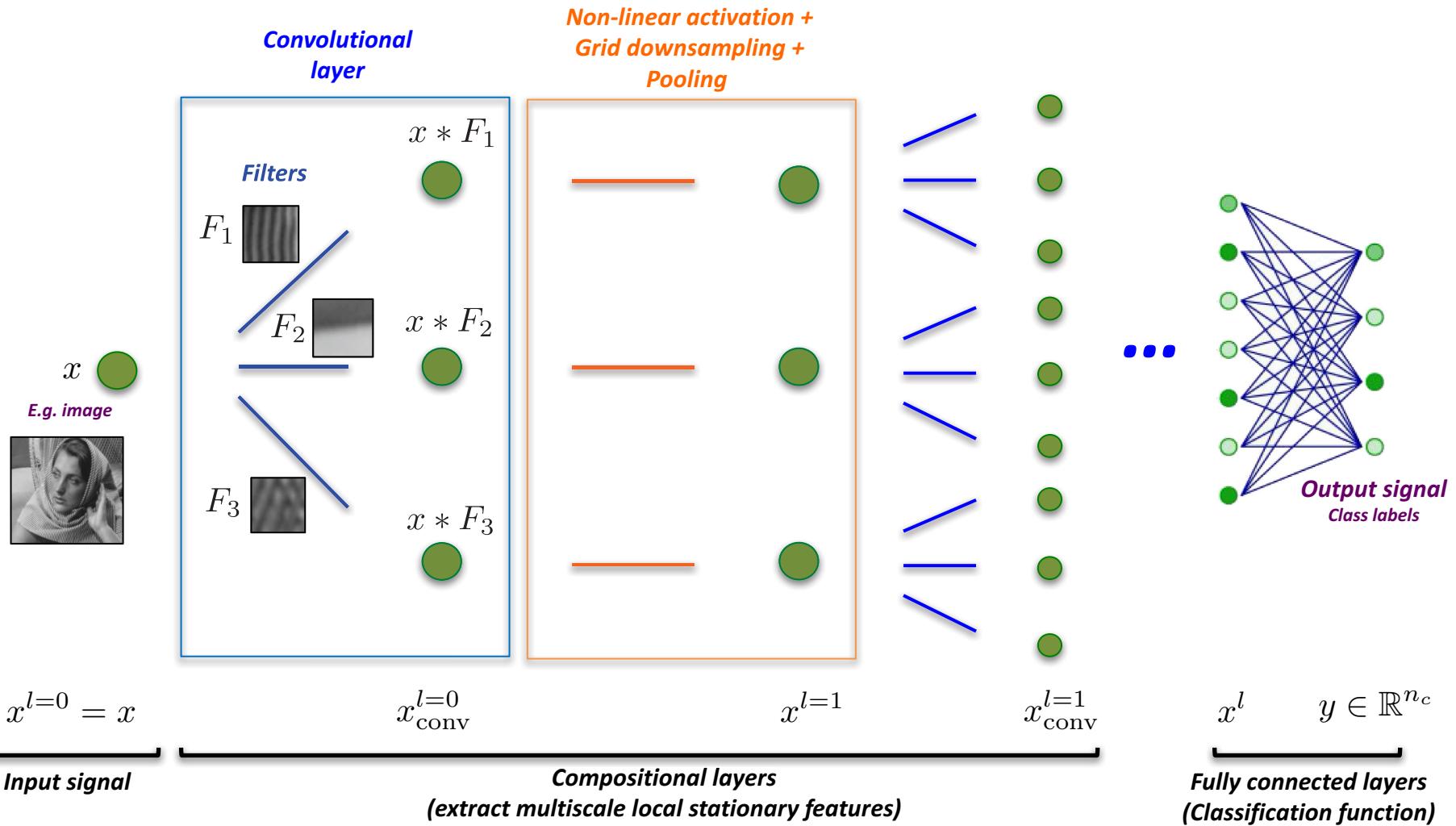
# Fully Connected Layer

- **Classifier:** After extracting multiscale locally stationary features, feed them to a classification function.
- How to design a (linear) classifier?  
*Fully connected neural networks.*

$$x_{\text{class}} = W x_{\text{feat}}$$



# Architecture of CNNs



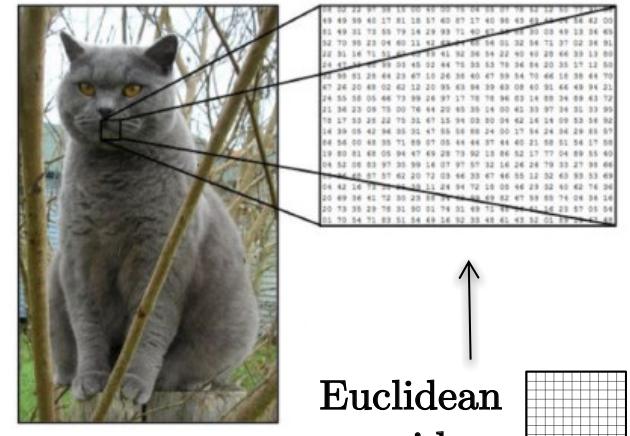
# Outline

- Data Science and Deep Learning
  - *A brief history*
  - *What makes DL great*
- Convolutional Neural Networks
  - *Architecture*
- **Convolutional Neural Networks on Graphs**
  - *Non-Euclidean data*
  - *Graph convolution, graph downsampling and pooling*
  - *Numerical experiments*
- Natural Language Processing
- Recommender Systems
- Conclusion

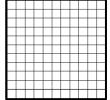
# CNNs are Great but...

- CNNs are designed for data lying on *Euclidean grids*:
  - (1) Convolution on Euclidean grids (FFT)
  - (2) Downsampling on Euclidean grids
  - (3) Pooling on Euclidean grids

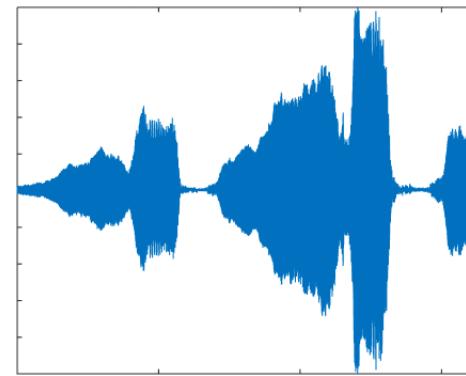
⇒ *Everything mathematically well defined and computationally fast!*



Euclidean grid



Images (2D, 3D)  
videos (2+1D)

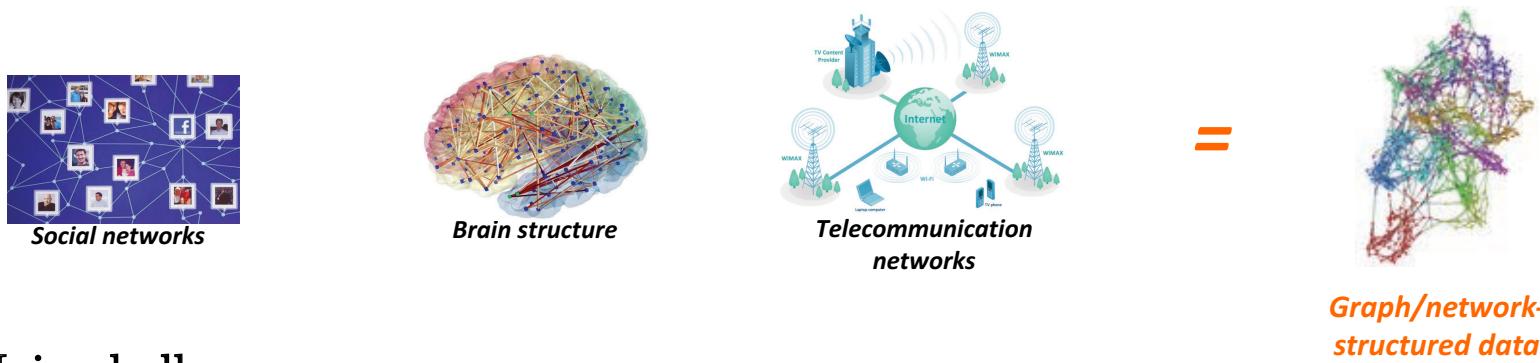


Sound (1D)

- But not all data lie on Euclidean grids.

# Non-Euclidean Data

- Examples of non-Euclidean data:
  - (1) *Social networks* (Facebook, Twitter)
  - (2) *Biological networks* (gene, brain connectivity)
  - (3) *Communication networks* (road, wireless, internet)
  - (4) *Word semantic networks* (NLP, translation, Q&A)



- Main challenges:
  1. *How to define compositionality on graphs (i.e. convolution, downsampling and pooling on graphs)?*
  2. *And how to make them numerically fast? (as fast as standard CNNs)*
- Current solution: Map graph-structured data to regular/Euclidean grids with e.g. kernel methods and apply standard CNNs. But handcrafting the mapping is (to my opinion) against CNN principle.

# CNNs for Graph-Structured Data

[Defferrard-B-Vandergheynst'16]

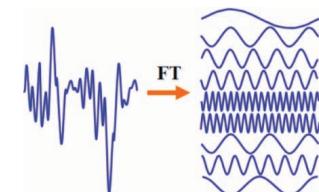
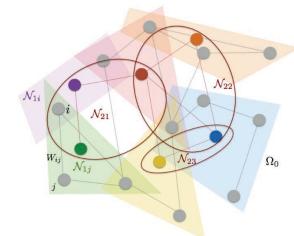
- Our contribution: Generalizing CNNs to any graph-structured data with *same linear time complexity* as standard CNNs.



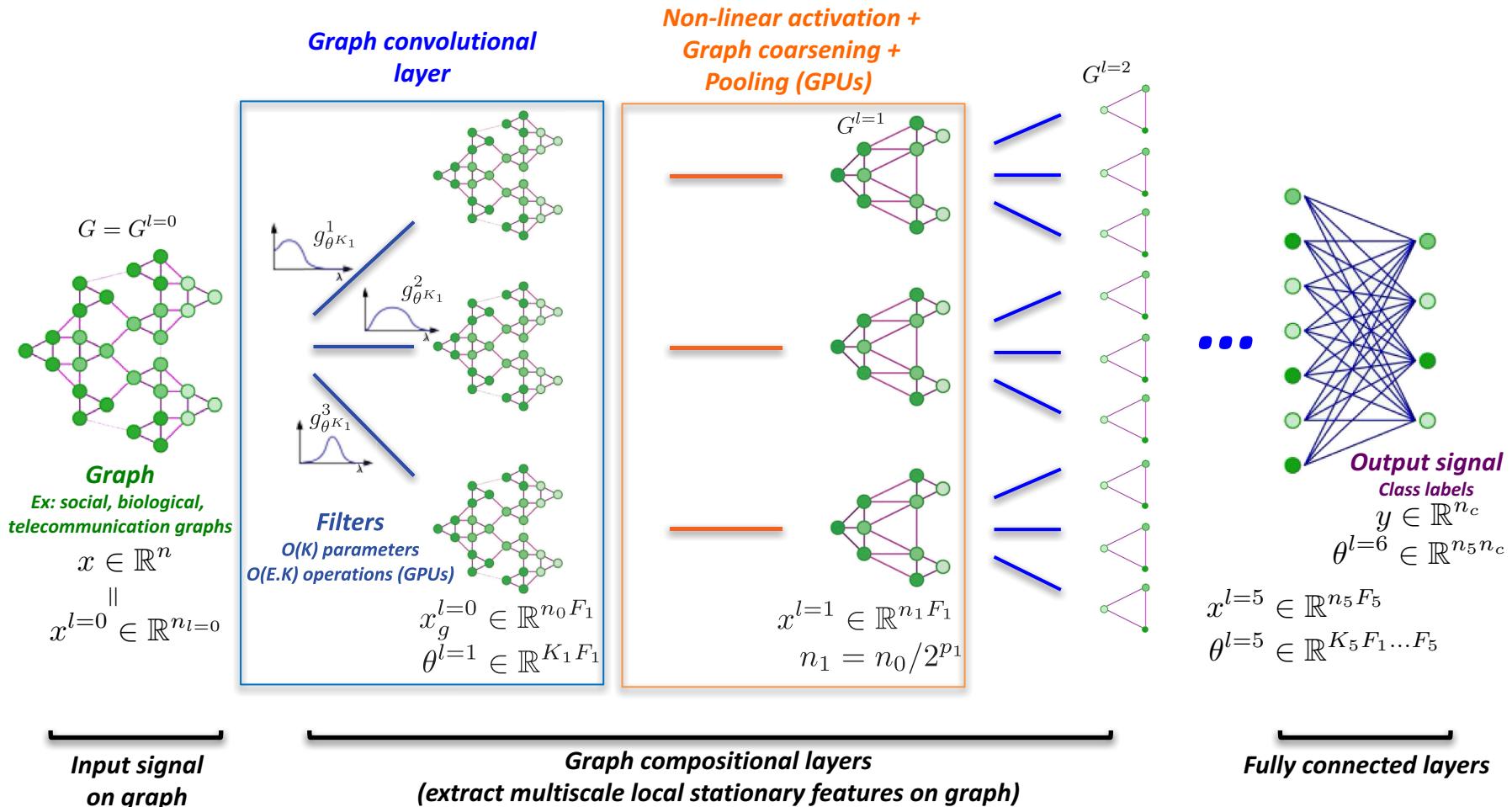
- What mathematical tools?
  1. *Graph spectral theory*,
  2. *Graph clustering*.
- This research was published at NIPS'16.  
Code/data are available on GitHub: [https://github.com/mdeff/cnn\\_graph](https://github.com/mdeff/cnn_graph)

# Related Works

- Categories of graph CNNs:
  - (1) *Spatial approach*
  - (2) *Spectral (Fourier) approach*
- Spatial approach:
  - *Local reception fields* [Coates-Ng'11, Gregor-LeCun'10]:  
Find compact groups of similar features, but no defined convolution.
  - *Locally Connected Networks* [Bruna-Zaremba-Szlam-LeCun'13]:  
Exploit multiresolution structure of graphs, but no defined convolution.
  - *Haar Scattering on Graphs* [Chen-Cheng-Mallat'14]:  
Use scattering transform on graphs to invariant signal descriptors.
  - *ShapeNet* [Bronstein-et.al.'15'16]:  
Generalization of CNNs to 3D-meshes. Convolution well-defined in these smooth low-dimensional non-Euclidean spaces. Handle multiple graphs.  
Obtained state-of-the-art results for 3D shape recognition.
- Spectral approach:
  - *Deep Spectral Networks* [Henaff-Bruna-LeCun'15]:  
Related to this work. We will compare later.

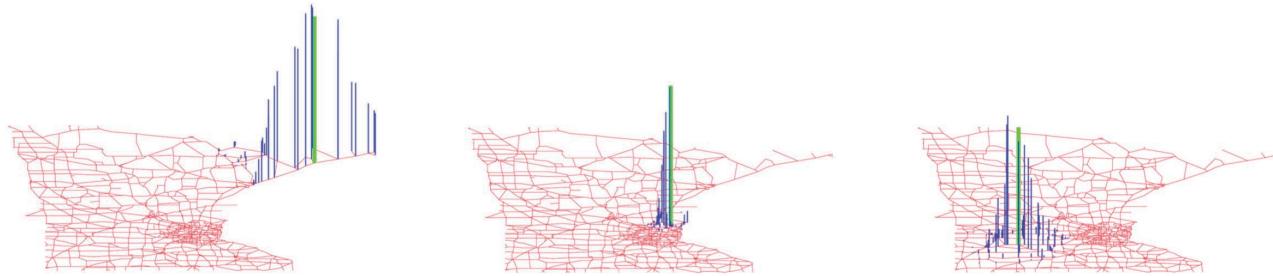


# Architecture of Graph CNNs

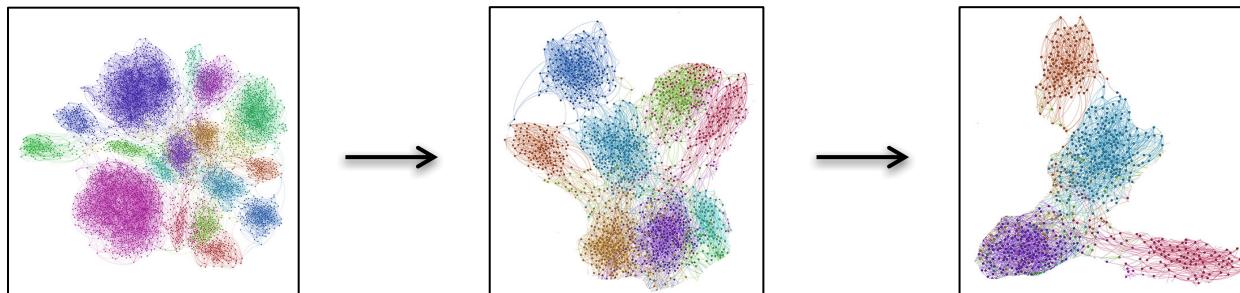


# Redefining Standard Layers

- Step 1. Convolutional layer: How to perform convolution on graphs?  
⇒ *Graph spectral theory*

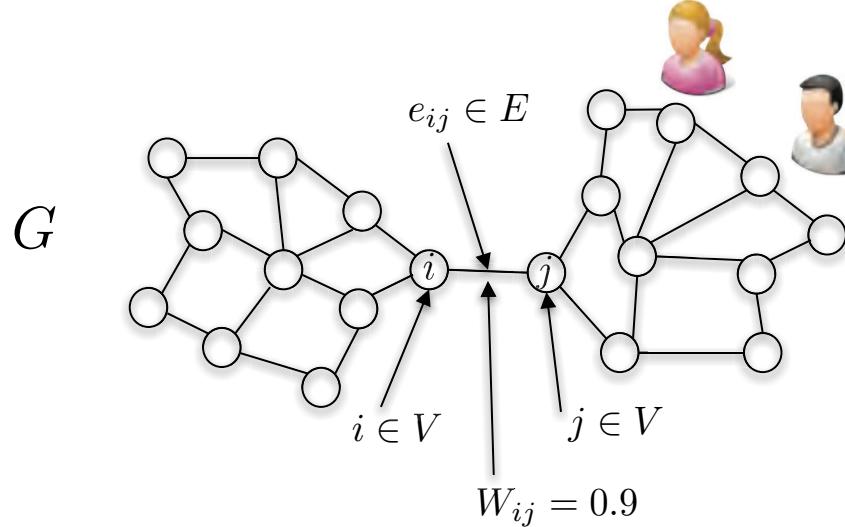


- Step 2. Compositional layer: How to “downsample” and pool on graphs?  
⇒ *Graph clustering*  
⇒ *Graph pooling*



# Step 1. Graph Convolutional Layer

- **Graphs:**  $G = (V, E, W)$ , with  $V$  set of vertices,  $E$  set of edges,  $W$  similarity matrix, and  $|V| = n$ .



- **Graph Laplacian** (*core operator to spectral graph theory [1]*):  
2<sup>nd</sup> order derivative operator on graphs

$$L = D - W$$

with  $D = \text{diag}\left(\sum_j W_{ij}\right)$

[1] Chung, 1997

# Graph Fourier Transform

- Fourier transform on graphs [2]:

Eigen-decomposition of  $L$ :

$$L = U \Lambda U^T$$

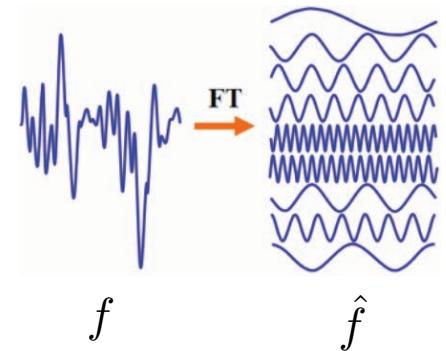
*Fourier modes*                    *Graph frequencies*

The Graph Fourier Transform of  $f \in \mathbb{R}^n$  is

$$\mathcal{F}_G f = \hat{f} = U^T f$$

The inverse GFT is defined as:

$$\mathcal{F}_G^{-1} \hat{f} = U \hat{f} = UU^T f = f,$$



[2] Hammond, Vandergheynst, Gribonval, 2011

# Convolution on Graphs

- Convolution on graphs is easy in the Fourier domain:

$$f *_G g = \mathcal{F}_G^{-1}(\mathcal{F}_G f \odot \mathcal{F}_G g)$$

Note (for later): It is also convenient to see:

$$f *_G g = \hat{g}(L)f,$$

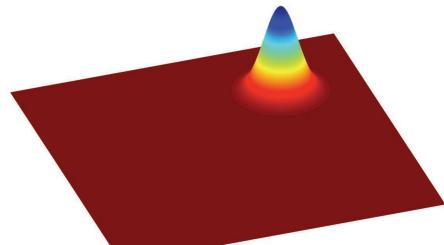
- [Bonus] Translation on graphs: A signal  $f$  defined on the graph can be translated to any vertex  $i$  using the graph convolution operator [3]:

$$T_i f := f *_G \delta_i,$$

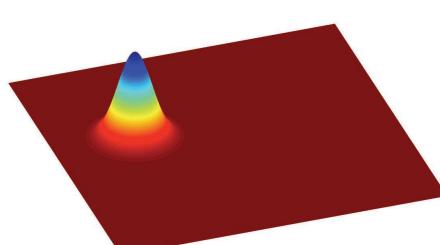
[3] Shuman, Ricaud, Vandergheynst, 2016

# Translation: Euclidean vs Non-Euclidean Domains

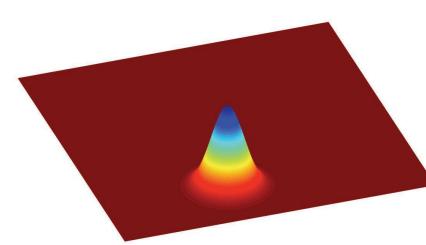
- Translated signals in the Euclidean domain:



(a)  $T_s f$

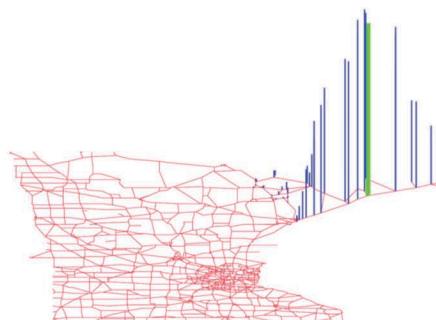


(b)  $T_{s'} f$

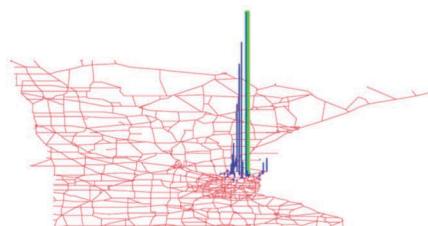


(c)  $T_{s''} f$

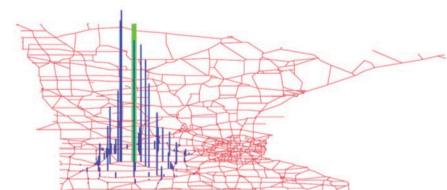
- Translated signals in a non-Euclidean domain:



(d)  $T_i f$



(e)  $T_{i'} f$

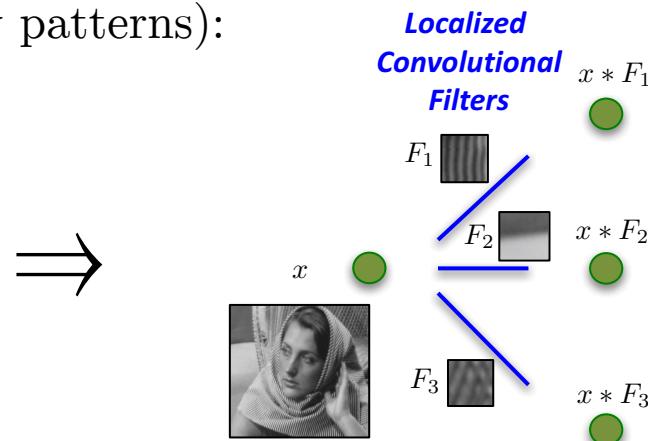
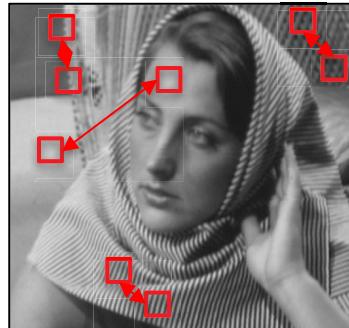


(f)  $T_{i''} f$

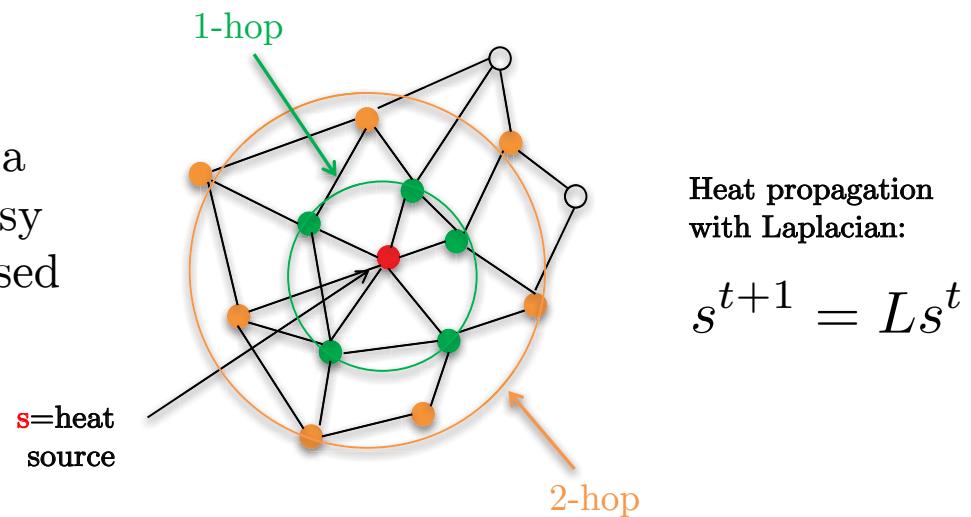
[Shuman-Ricaud-Vandergheynst'16]

# Localized Filters on Graphs

- As in standard CNNs, we need to define **localized filters on graphs** (to find local stationary patterns):



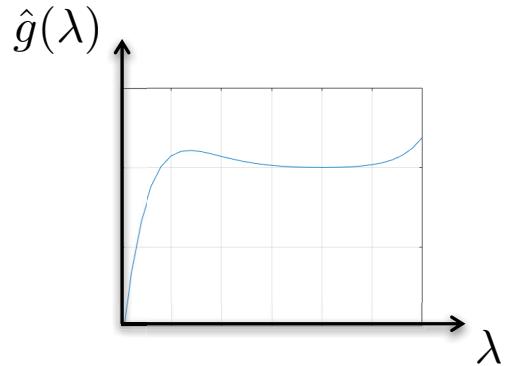
- Key observation:** Graph Laplacian operation increases the support of a Dirac function by 1-hop  $\Rightarrow$  it is easy to control the size of Laplacian-based filters!



# Localized Filters on Graphs

- Spectral kernels as Laplacian polynomial [2]:

$$\hat{g}(\lambda_l) = \hat{p}_K(\lambda_l) := \sum_{k=0}^K a_k \lambda_l^k$$

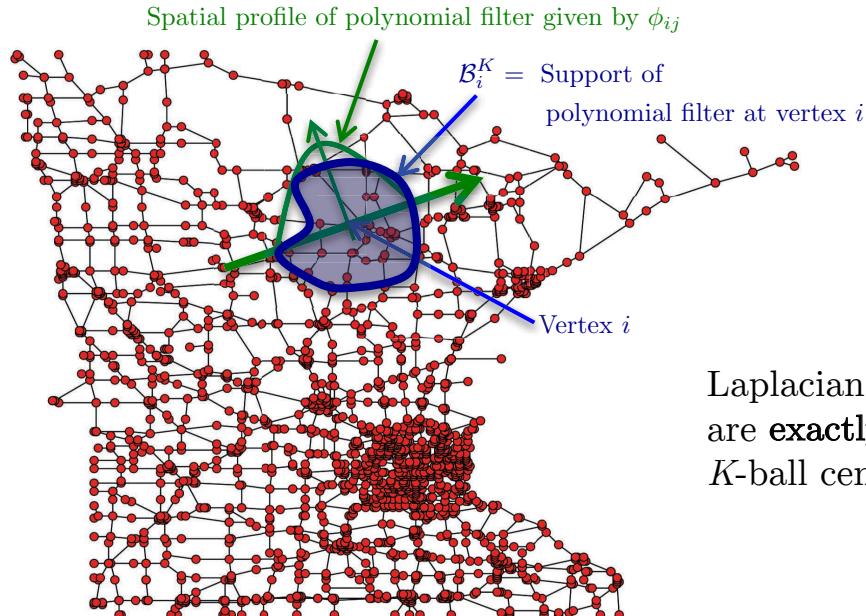


**Corollary 1.** Consider the function  $\phi_{ij}$  such that

$$\phi_{ij} = (T_i p_K)(j) = (p_K *_G \delta_i)(j) = (\hat{p}_K(L)\delta_i)(j) = \left( \sum_{k=0}^K a_k L^k \delta_i \right)(j).$$

Then

$$\phi_{ij} = 0 \quad \text{if} \quad d_G(i, j) > K.$$

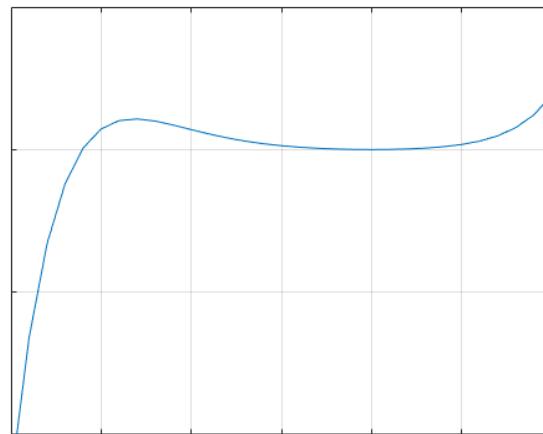


# Localized Filters on Graphs

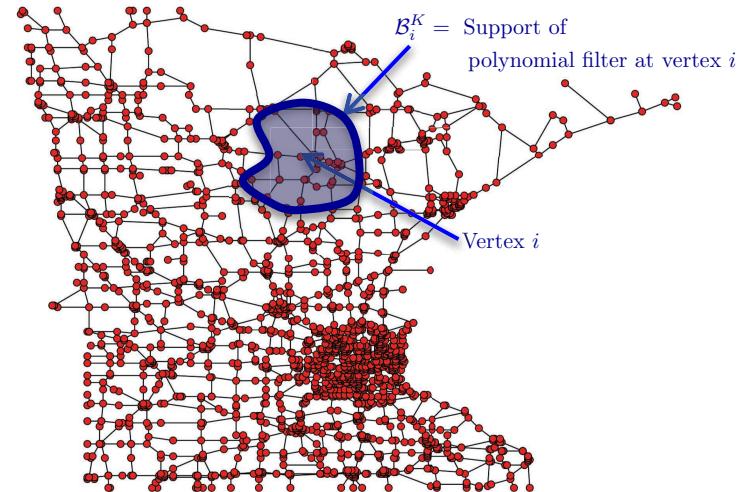
**Corollary 2.** *Localized filters on graphs are defined according to the principle:*

$$\text{Frequency smoothness} \Leftrightarrow \text{Spatial graph localization}$$

This is obviously the Heisenberg's uncertainty principle extended to the graph setting. Recent papers have studied the uncertainty principle on graphs.



Frequency/spectral domain



Spatial/graph domain

# How fast is Convolution on Graphs?

- **Graph convolution:**

$$y = x *_G p_K = \hat{p}_K(L)x = \sum_{k=0}^K a_k L^k x$$

- Complexity is linear:  $O(n)$

Denote  $X_k = LX_{k-1}$  with  $X_0 = x$  and rewrite  $y = \sum_k a_k X_k$ .

Observe that all  $\{X_k\}$  are generated by multiplication of the *sparse* L matrix and a vector  $\Rightarrow$  Complexity is  $O(EK) = O(n)$  for *sparse* graphs.

- **GPU parallel implementation:** Linear algebra operations can be done in parallel, allowing a fast GPU implementation of graph convolution.

# Convolution on Graphs with Chebyshev

- Monomial graph convolution is actually not optimal for learning:

$$y = x *_G p_K = \hat{p}_K(L)x = \sum_{k=0}^K a_k L^k x$$

Monomial basis  $\{1, x, x^2, \dots, x^K\}$  do **not** form an **orthogonal** basis, which prevents learning good spectral filters.

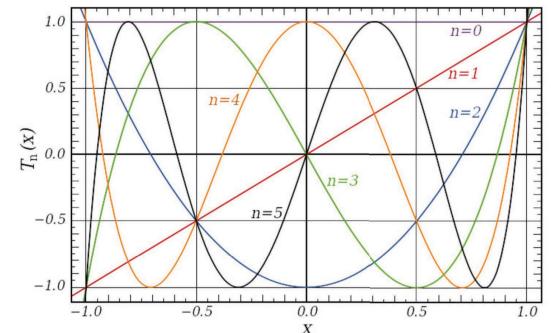
- Graph convolution with (orthogonal) Chebyshev polynomials :

$$y = x *_G q_K = \sum_{k=0}^K \theta_k T_k(L)x,$$

- Complexity is also linear:  $O(n)$

Denote  $X_k = T_k(L)X_{k-1} = 2LX_{k-1} - X_{k-2}$  and rewrite  
 $y = \sum_k \theta_k X_k$ .

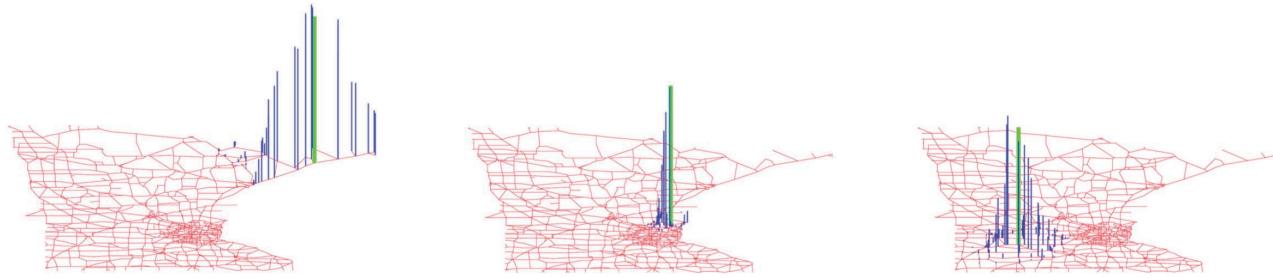
Complexity is  $O(n)$  for *sparse* graphs.



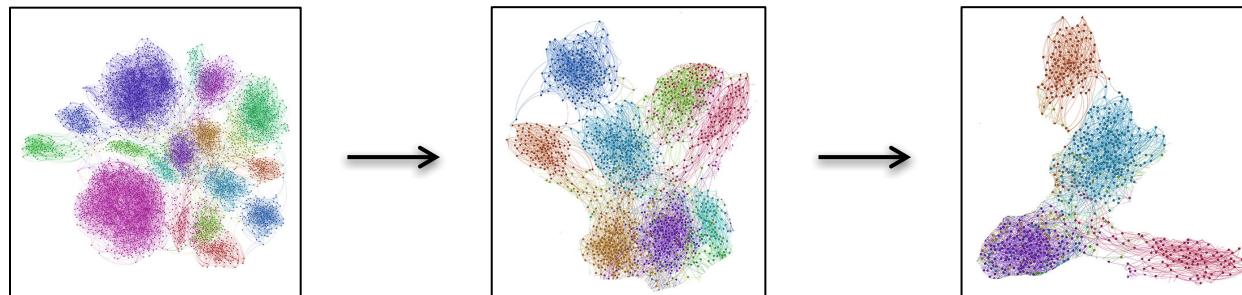
First six Chebyshev functions  
 $\{T_0, T_1, T_2, T_3, T_4, T_5\}$ .

# Redefining Standard Layers

- Step 1. Convolutional layer: How to perform convolution on graphs?  
⇒ *Graph spectral theory*



- Step 2. Compositional layer: How to “downsample” and pool on graphs?  
⇒ *Graph clustering*  
⇒ *Graph pooling*



## Step 2. Graph “Downsampling”

- Graph downsampling is graph coarsening, also graph clustering/partitioning.  
This step is essential to compose lower-level features to more complex patterns.

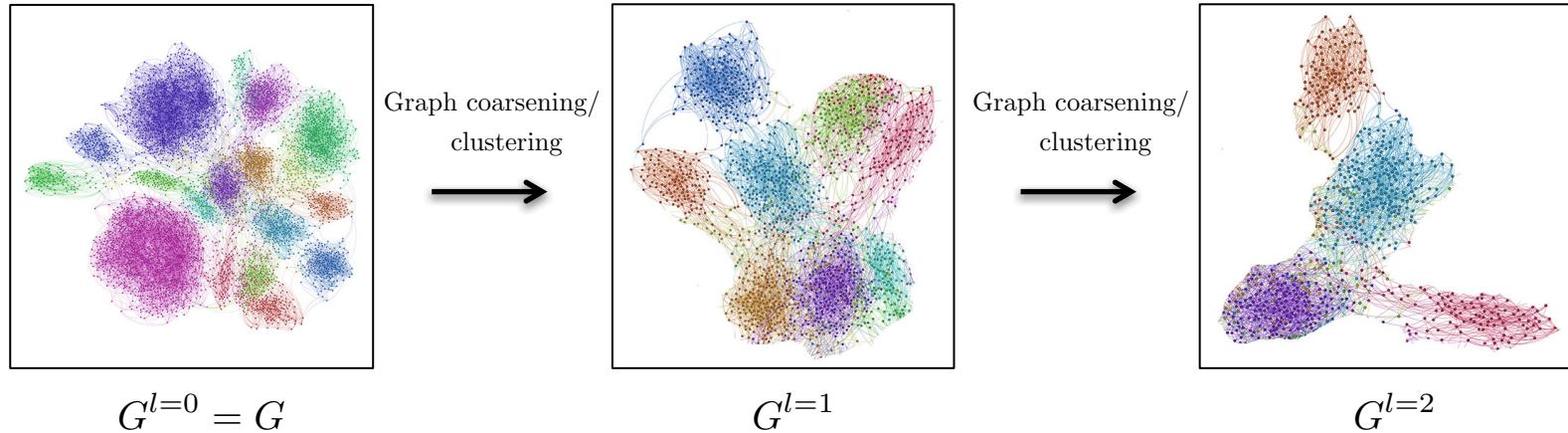


Fig. Illustration of graph coarsening.

- Graph partitioning is NP-hard  $\Rightarrow$  Approximation needed.

# Graph Partitioning

- **Balanced Cuts [4]:** Two powerful measures of graph clustering are the *Normalized Cut* and *Normalized Association* defined as:

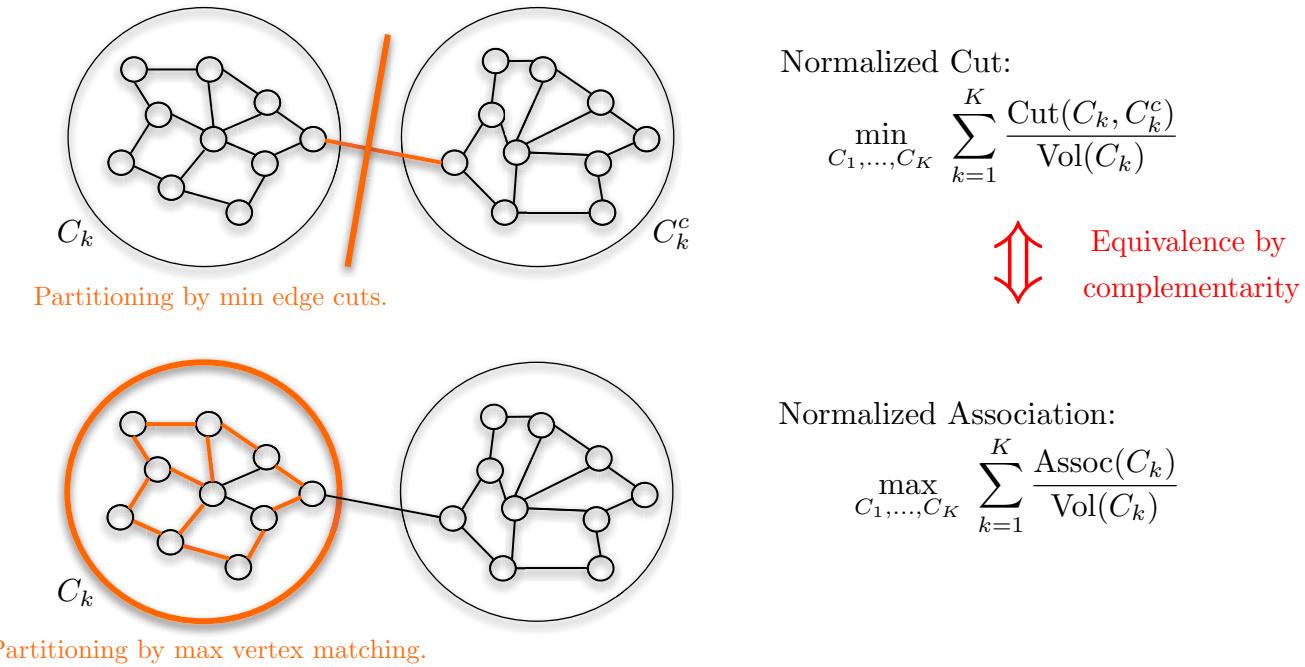


Figure 5: Equivalence between NCut and NAssoc partitioning.

where  $\text{Cut}(A, B) := \sum_{i \in A, j \in B} W_{ij}$ ,  $\text{Assoc}(A) := \sum_{i \in A, i \in B} W_{ij}$ ,  $\text{Vol}(A) := \sum_{i \in A, j \in B} d_i$ , and  $d_i := \sum_{j \in V} W_{ij}$  is the degree of vertex  $i$ .

[4] Shi, Malik, 2000

# Graclus Graph Partitioning

- **Graclus [5]:** It is a greedy (fast) technique that computes clusters that locally maximize the Normalized Association.

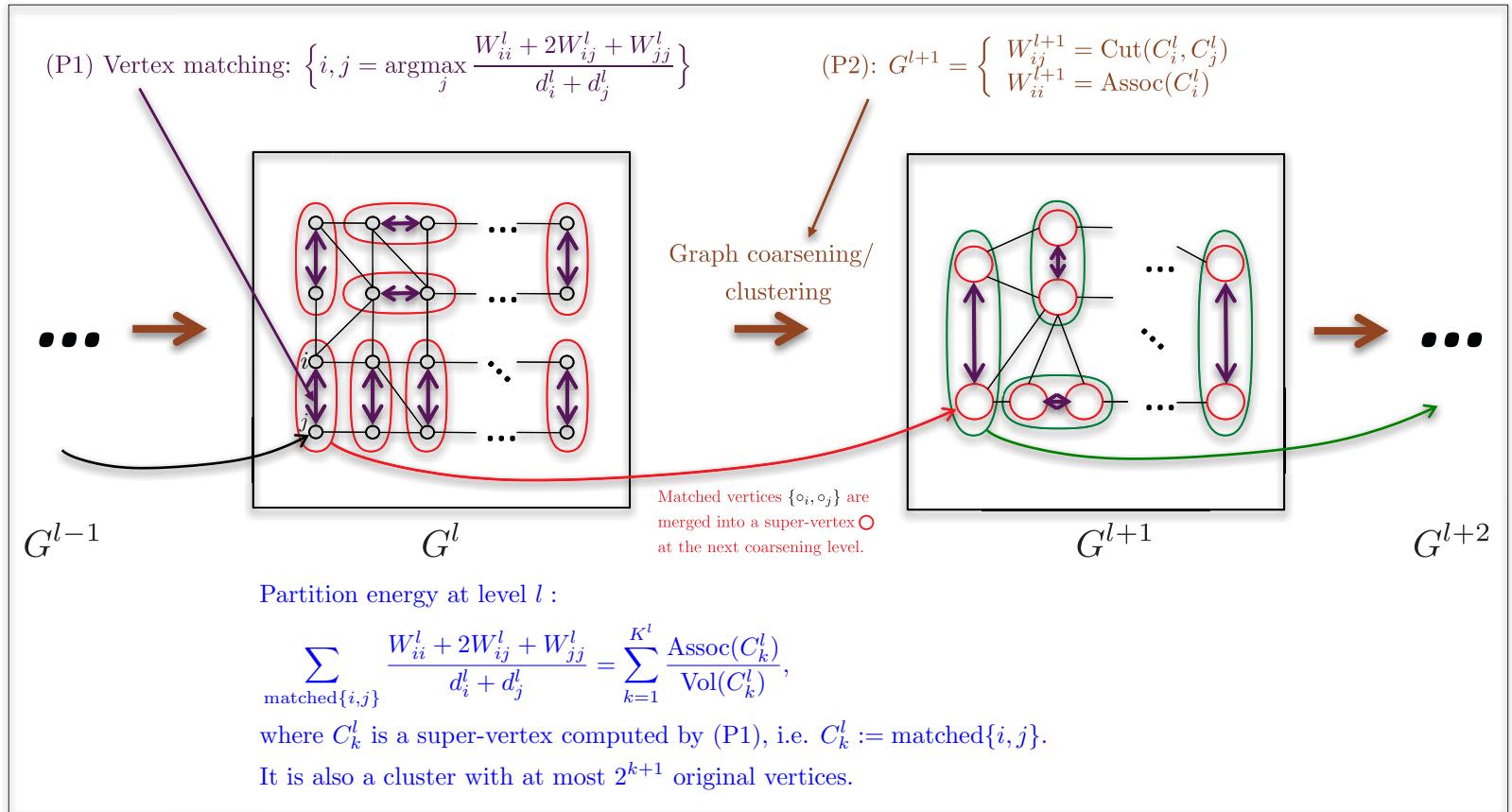


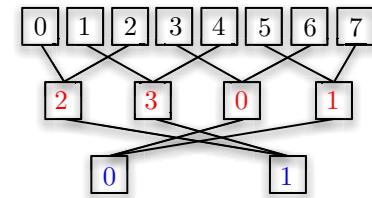
Figure 6: Graph coarsening with Graclus. Graclus proceeds by two successive steps: (P1) Vertex matching, and (P2) Graph coarsening. These two steps provide a local solution to the Normalized Association clustering problem at each coarsening level  $l$ .

[5] Dhillon, Guan, Kulis, 2007

# Fast Graph Pooling

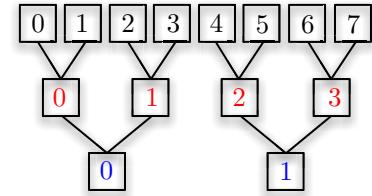
- **Graph pooling** such as **max** pooling or **average** pooling transfer feature value to the next layer.
- **Unstructured pooling is (highly) inefficient:** The graph and all its coarsened versions require to store a table with all matched vertices.  
⇒ Memory consuming and not (easily) parallelizable.
- Structured pooling can be as efficient as a 1D grid pooling (GPUs): Use a ***binary tree arrangement*** of the nodes such that adjacent nodes are hierarchically merged at the next coarser level.

Graph pooling:



Unstructured arrangement of vertices

Graph pooling:



Binary tree arrangement of vertices

# Reindexing Process

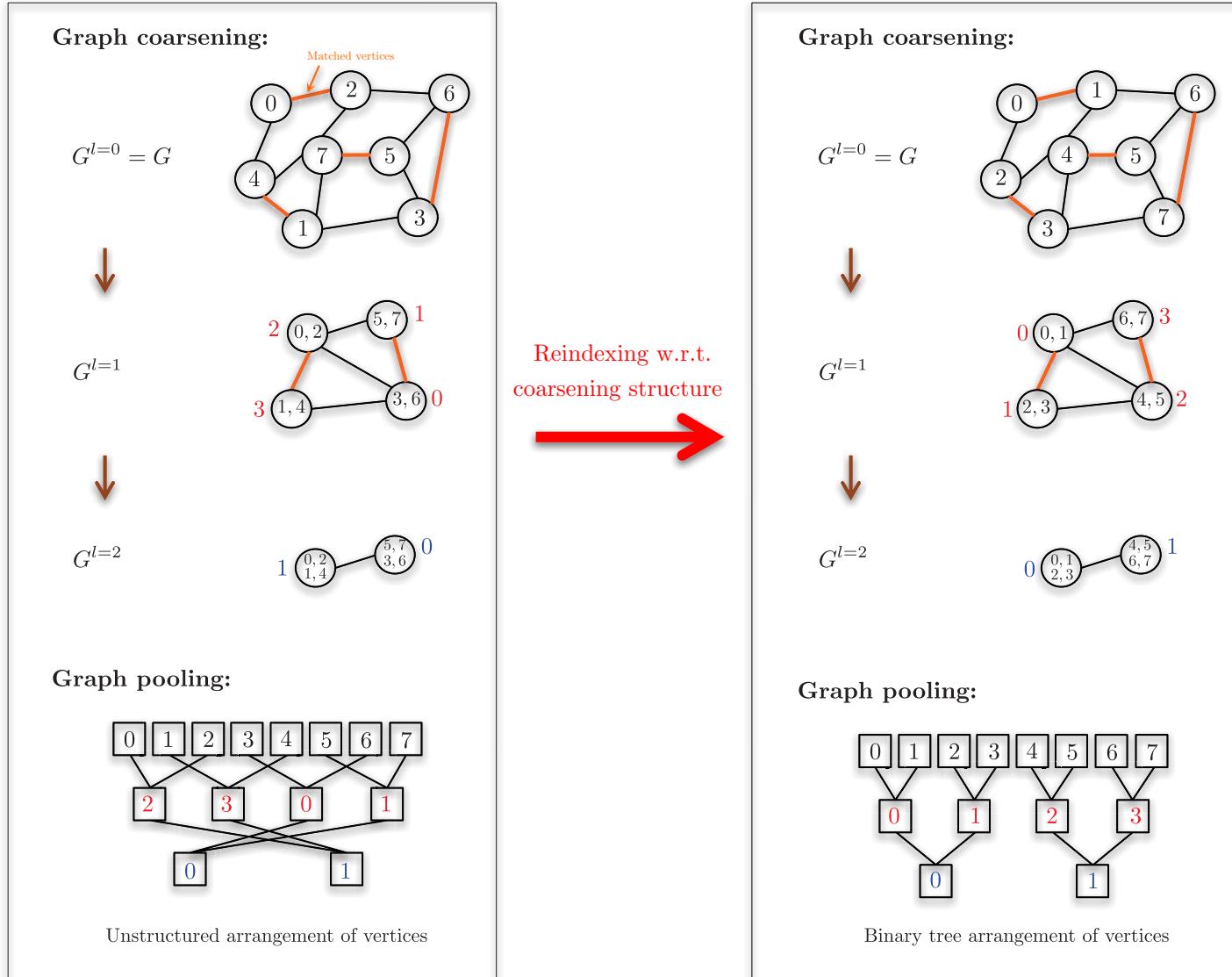
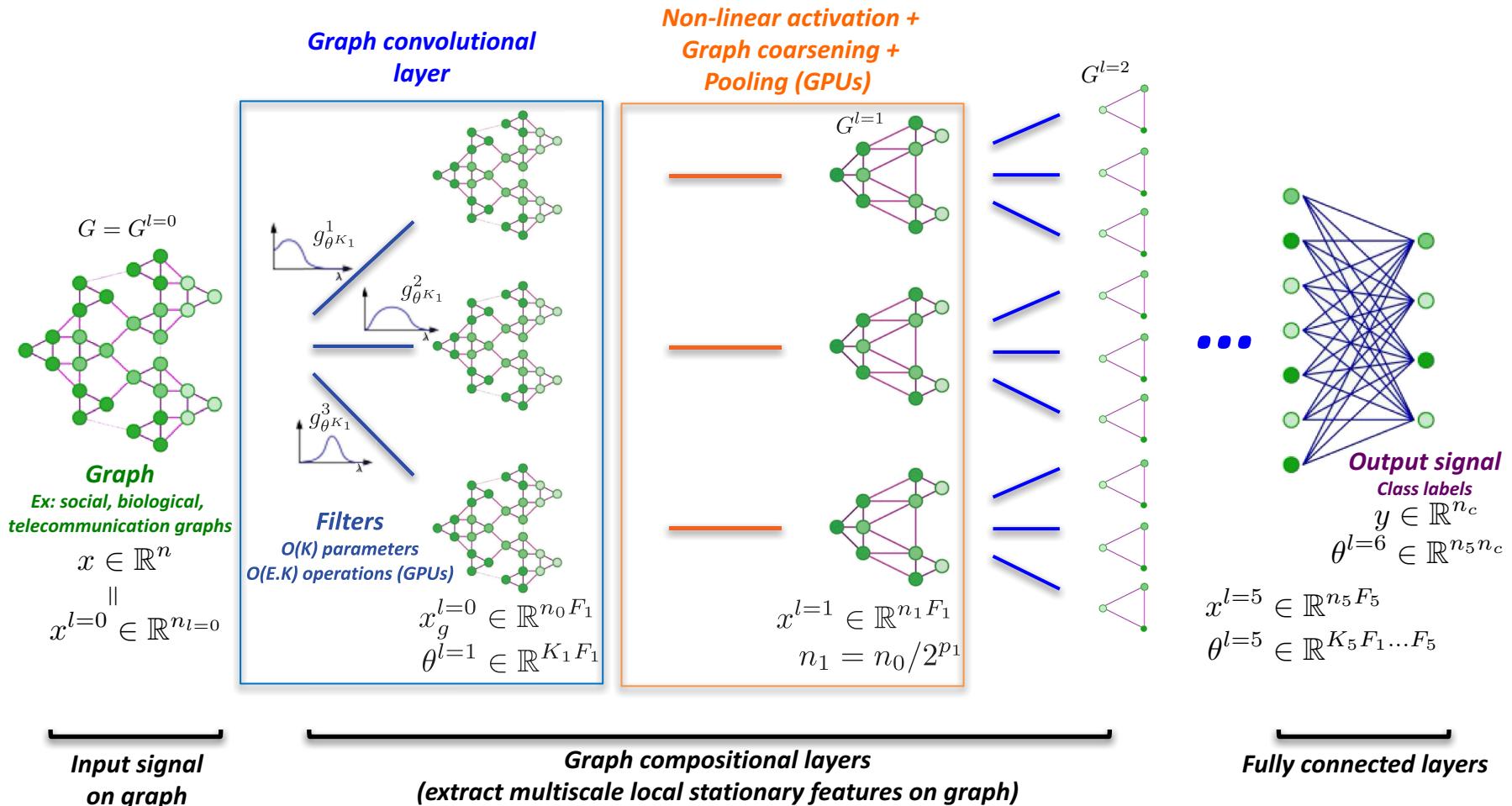


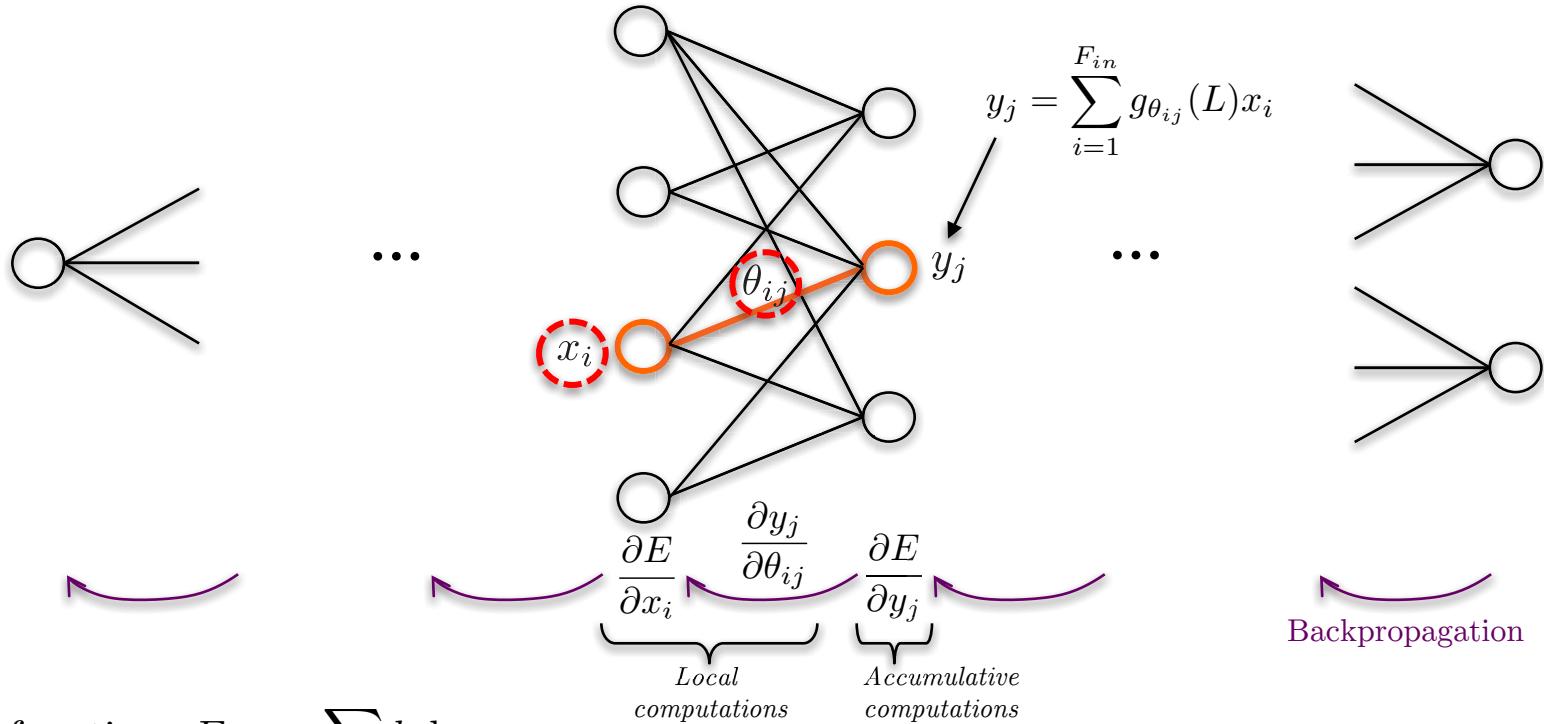
Figure 7: Fast graph pooling using graph coarsening structure. The binary tree arrangement of vertices allows a very efficient pooling on graphs, as fast as a regular 1D Euclidean grid pooling.

# Architecture of Graph CNNs



# Optimization

- Backpropagation [6] = Chain rule applied to the neurons at each layer.



$$\text{Loss function: } E = - \sum_{s \in S} l_s \log y_s$$

$$\text{Gradient descents: } \begin{cases} \theta_{ij}^{n+1} = \theta_{ij}^n - \tau \cdot \frac{\partial E}{\partial \theta_{ij}} \\ x_i^{n+1} = x_i^n - \tau \cdot \frac{\partial E}{\partial x_i} \end{cases}$$

[6] Werbos 1982 and Rumelhart, Hinton, Williams, 1985

Xavier Bresson

Local gradients:

$$\frac{\partial E}{\partial \theta_{ij}} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial \theta_{ij}} = \sum_{s \in S} [X_{0,s}, \dots, X_{K-1,s}]^T \frac{\partial E}{\partial y_{j,s}}$$

$$\frac{\partial E}{\partial x_i} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial x_i} = \sum_{j=1}^{F_{out}} g_{\theta_{ij}}(L) \frac{\partial E}{\partial y_j}$$

# Numerical Experiments

- **Platforms and hardware:** All experiments are carried out with:
  1. *TensorFlow (Google AI software) [7]*
  2. *GPU NVIDIA K40*
- **Reproducibility:** Code/data are on GitHub.
- **Goals of the experiments:**
  1. *Check linear complexity*
  2. *Compare with [LeCun-et al. '15]*

[7] Abadi-et.al. 2015

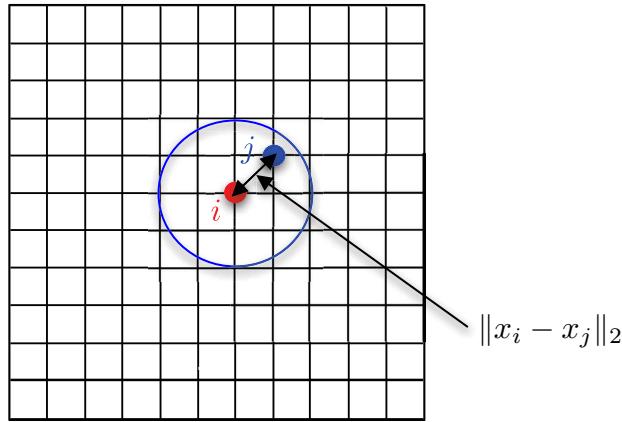
# Graph CNNs for Image Classification

- MNIST is the most popular dataset in deep learning [8].

It has 70,000 images represented on a 2D grid of size 28x28 ( $\text{dim}=28^2=784$ ) of handwritten digits, from 0 to 9.

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

- Graph: A k-NN graph ( $k=8$ ) of the Euclidean grid:



$$W_{ij} = e^{-\|x_i - x_j\|_2^2 / \sigma}$$

[8] LeCun, Bottou, Bengio, 1998

# Performances

- **Results:** Classification accuracies

Algorithm	Accuracy
Linear SVM	91.76
Softmax	92.36
CNNs [LeNet5]	99.33
graph CNNs: CN32-P4-CN64-P4-FC512-softmax	99.18

- **CPU vs. GPU:** A GPU implementation of graph CNNs is **8x faster** than a CPU implementation, same order of speed up as standard CNNs.

Architecture	Time	Speedup
CNNs with CPU	210ms	-
CNNs with GPU NVIDIA K40c	31ms	6.77x
graph CNNs with CPU	1600ms	-
graph CNNs with GPU NVIDIA K40c	200ms	8.00x

# Numerical Linear Complexity

- Advantages over [Henaff-Bruna-LeCun'15]:

(1) Computational complexity  $O(n)$  vs.  $O(n^2)$ :

$$y = x *_G h_K = \hat{h}_K(L)x = U(\hat{h}_K(\Lambda)(U^T x)),$$

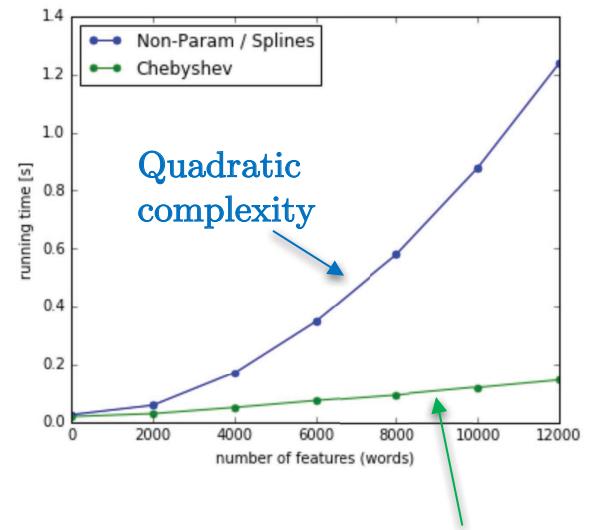
where  $\hat{h}_K$  are spectral filters based on  $K^{th}$ -order splines,  $\Lambda, U$  are the eigenvalue, eigenvector matrices of the graph Laplacian  $L$ , and  $U^T, U$  are full  $O(n^2)$  matrices.

(2) no EVD  $O(n^3)$  vs. [HBL15]

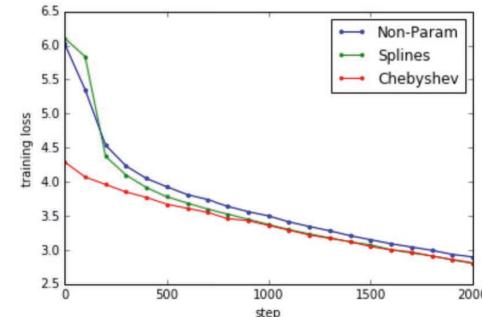
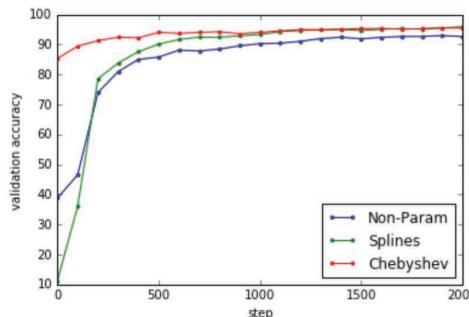
(3) Accuracy:

Architecture	Accuracy
CN10 [HBL15]	97.26
CN10 [Ours]	97.48
CN32-P4-CN64-P4-FC512 [HBL15]	97.75
CN32-P4-CN64-P4-FC512 [Ours]	99.18

(4) Learn faster:



Quadratic complexity  
Linear complexity



# Graph CNNs as a Tool for Data Analysis?

- Standard CNN is the Swiss knife of Computer Vision.  
⇒ Graph CNNs may have the same potential for any (un)structured data.
- PI' works:
  1. Recurrent neural networks on graphs for NLP [Seo-Defferrard-Vandergheynst-B, Structured Sequence Modeling with Graph Convolutional Recurrent Networks'16]
  2. Recommender systems [Monti-Bronstein-B, Geometric Matrix Completion with Recurrent Multi-Graph Neural Networks'17] – state-of-the-art
- Recent works that build on graph CNNs:
  1. Brain analysis [Distance Metric Learning using Graph Convolutional Networks: Application to Functional Brain Networks'17]
  2. Temporal time series analysis [Dynamic Graph Convolutional Networks'17]
  3. Traffic Control [Deep Deterministic Policy Gradient for Urban Traffic Light Control'17]
  4. Chemistry [Neural Message Passing for Quantum Chemistry'17]

# Outline

- Data Science and Deep Learning
  - *A brief history*
  - *What makes DL great*
- Convolutional Neural Networks
  - *Architecture*
- Convolutional Neural Networks on Graphs
  - *Non-Euclidean data*
  - *Graph convolution, graph downsampling and pooling*
  - *Numerical experiments*
- **Natural Language Processing**
- Recommender Systems
- Conclusion

# Natural Language Processing

- NLP problem:

*Input:* A sequence of words

Yesterday I went to the beach and I saw a ...

*Output:* Probability of the next word →

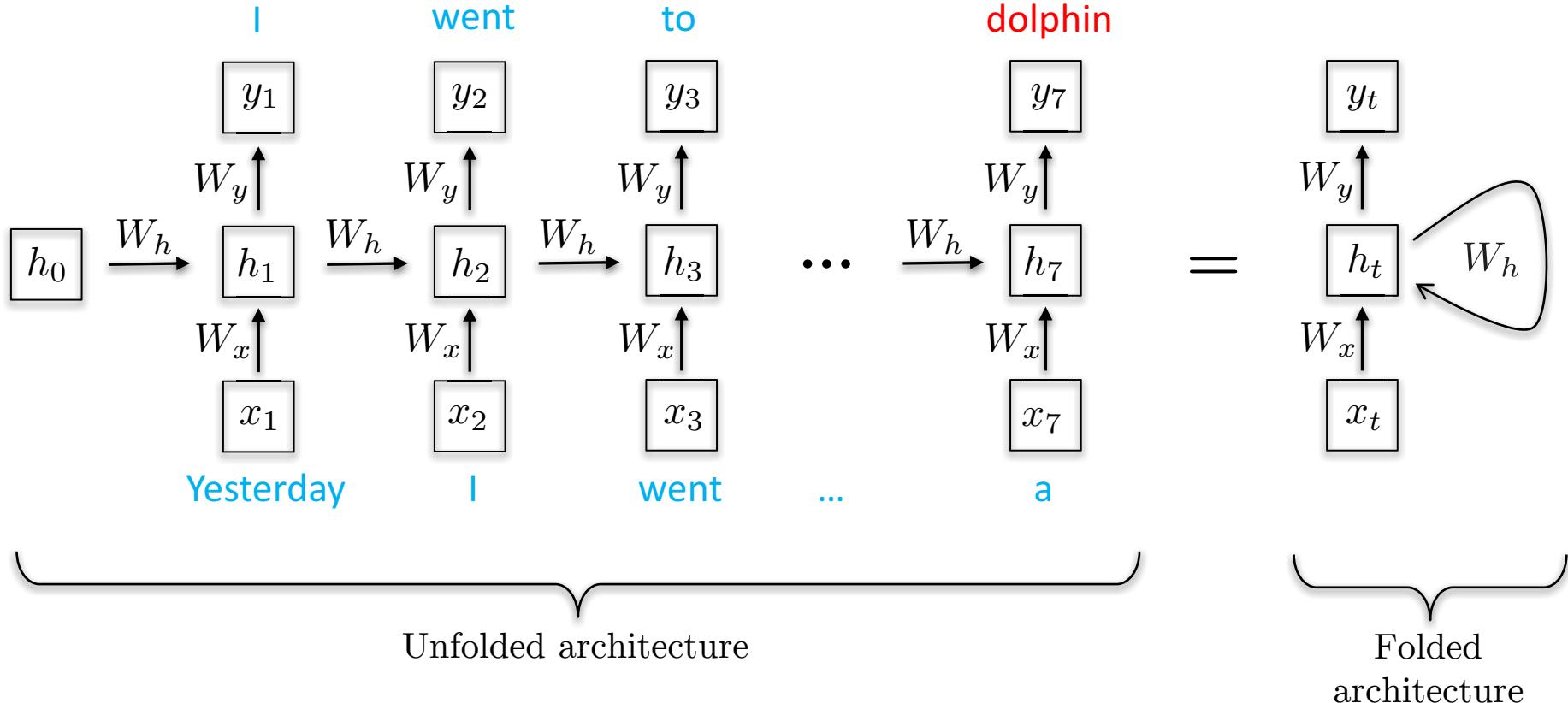
next word	probability
dolphin	11.5%
boat	8.2%
:	
elephant	0.1%
tiger	0.1%
:	
cooking	0.0%
with	0.0%

- State-of-the-art solution since 2015: Learn to predict next word with RNN.
- Google's Neural Machine Translation system (Sept 2016):

“Today we announce the Google Neural Machine Translation system (GNMT), which utilizes state-of-the-art training techniques to achieve the largest improvements to date for machine translation quality.”

# RNN for Natural Language Processing

- Vanilla RNN architecture for NLP:

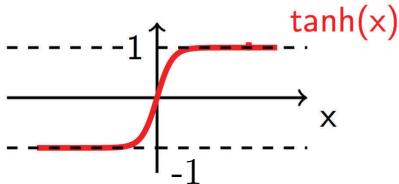
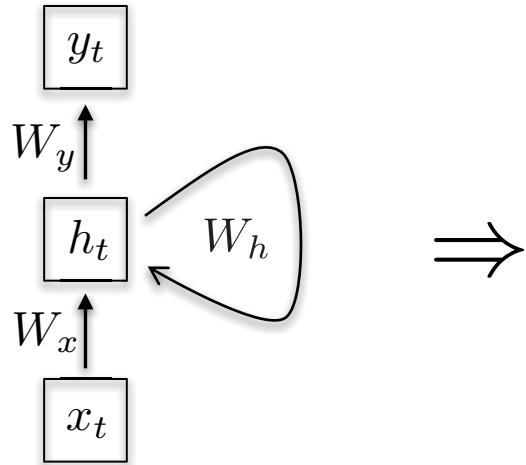


$x_t$  : Input at t  
 $h_t$  : Memory vector at t  
 $y_t$  : Probability of next word at t

$W_x$  : Embedding matrix  
 $W_h$  : Memory matrix  
 $W_y$  : Decoding matrix

# Vanilla RNN

- Iterative scheme:



$$h_t = \tanh(W_h h_t + W_x x_t)$$

$$y_t = \text{softmax}(W_y h_t)$$

$$\text{Softmax} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} \frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3}} \\ \frac{e^{x_2}}{e^{x_1} + e^{x_2} + e^{x_3}} \\ \frac{e^{x_3}}{e^{x_1} + e^{x_2} + e^{x_3}} \end{pmatrix}$$

- Major shortcoming:

VRNN do not learn long-term dependencies

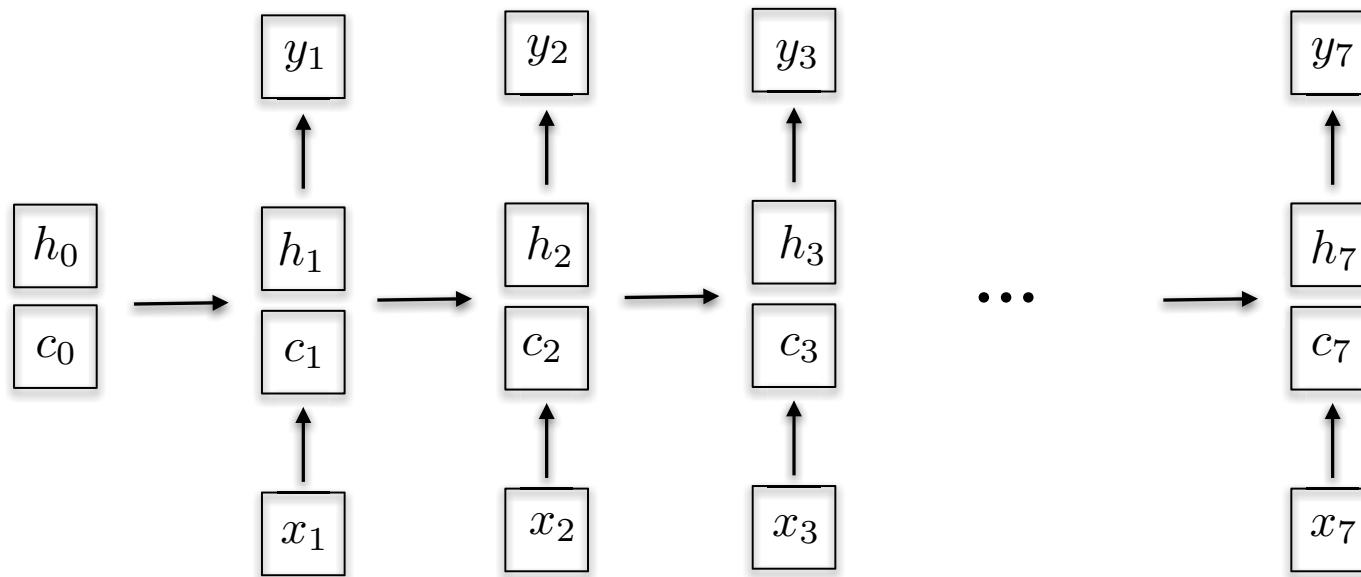
*Vanishing gradient problem (largest eigenvalue of  $W_h$  is less than 1)*

$$\text{Ex. Softmax} \begin{pmatrix} -1.5 \\ 3 \\ 1 \end{pmatrix} = \begin{pmatrix} .01 \\ .87 \\ .12 \end{pmatrix}$$

# Long Short-Term Memory Networks

[Hochreiter-Schmidhuber'97]

- LSTM architecture



$x_t$  : Input at  $t$

$h_t$  : *Short-term* memory vector at  $t$

$c_t$  : *Long-term* memory vector at  $t$

$y_t$  : Probability of next word at  $t$

# Non-Linear Dynamics

- Iterative scheme:

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(z_t)$$

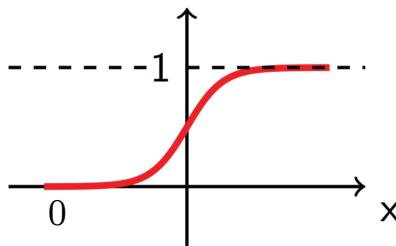
$$h_t = o_t \odot \tanh(c_t)$$

$$y_t = \text{softmax}(W_y h_t)$$

with	$i_t = \sigma(W_h^i h_t + W_x^i x_t)$	Input gate
	$f_t = \sigma(W_h^f h_t + W_x^f x_t)$	Forget gate ( $f=0,1$ )
	$o_t = \sigma(W_h^o h_t + W_x^o x_t)$	Output gate
	$z_t = \sigma(W_h^z h_t + W_x^z x_t)$	

Sigmoid function:

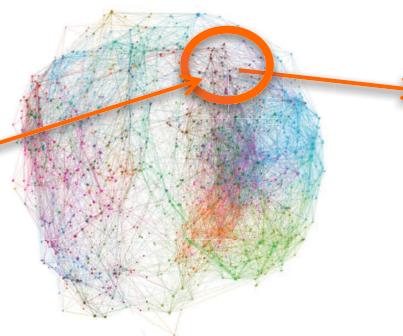
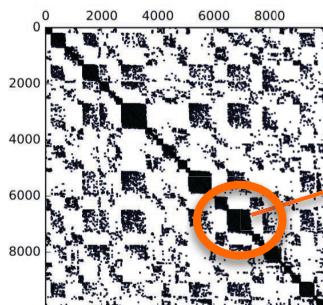
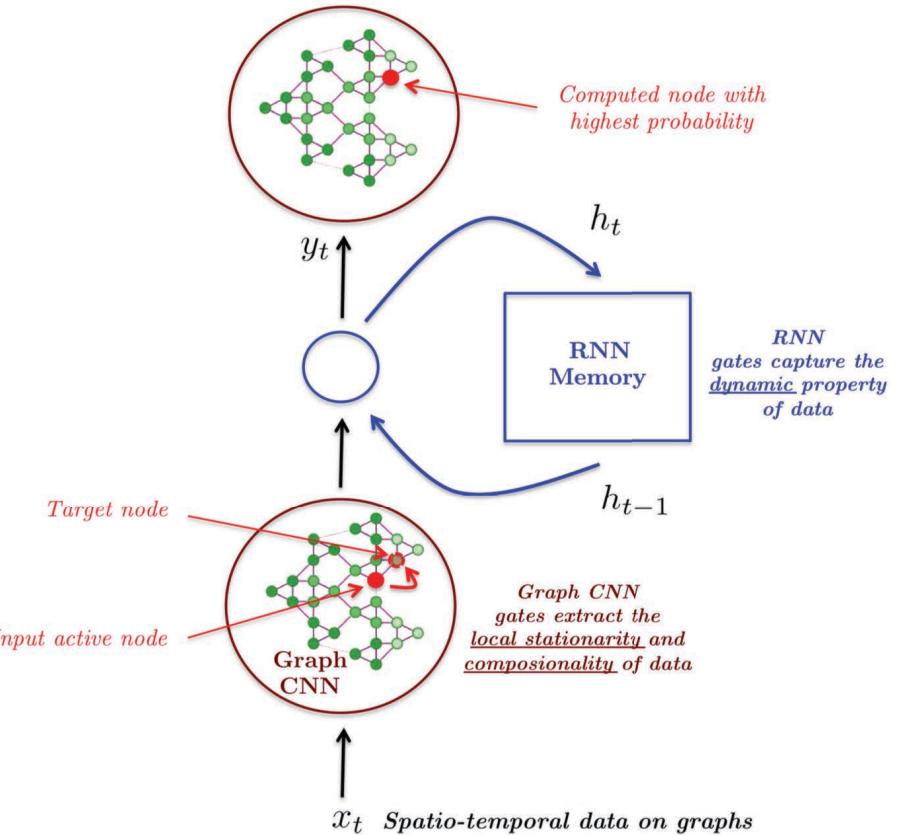
$$\sigma(x) = \frac{1}{1+e^{-x}}$$



# LSTM on Graphs

[Seo-Defferrard-Vandergheynst-B'16]

- Combine CNN + LSTM:
  - CNN analyzes static information from data; finds multiscale local stationary patterns.
  - LSTM identifies data dynamics; non-linear temporal properties of data.
- Graph CNN: Observe that words form a space of topics that can be modeled by graphs of word relationships.



Sport  
Ball  
Wimbledon  
Federer  
Tennis  
Game

# Iterative scheme

$$c_t = f_t \odot c_{t-1} + i_t \odot \tanh(z_t)$$

$$h_t = o_t \odot \tanh(c_t)$$

$$y_t = \text{softmax}(W_y h_t)$$

with

$$i_t = \sigma(W_h^i h_t + W_c^i x_t^{cnn})$$

$$f_t = \sigma(W_h^f h_t + W_c^f x_t^{cnn})$$

$$o_t = \sigma(W_h^o h_t + W_c^o x_t^{cnn})$$

$$z_t = \sigma(W_h^z h_t + W_c^z x_t^{cnn})$$

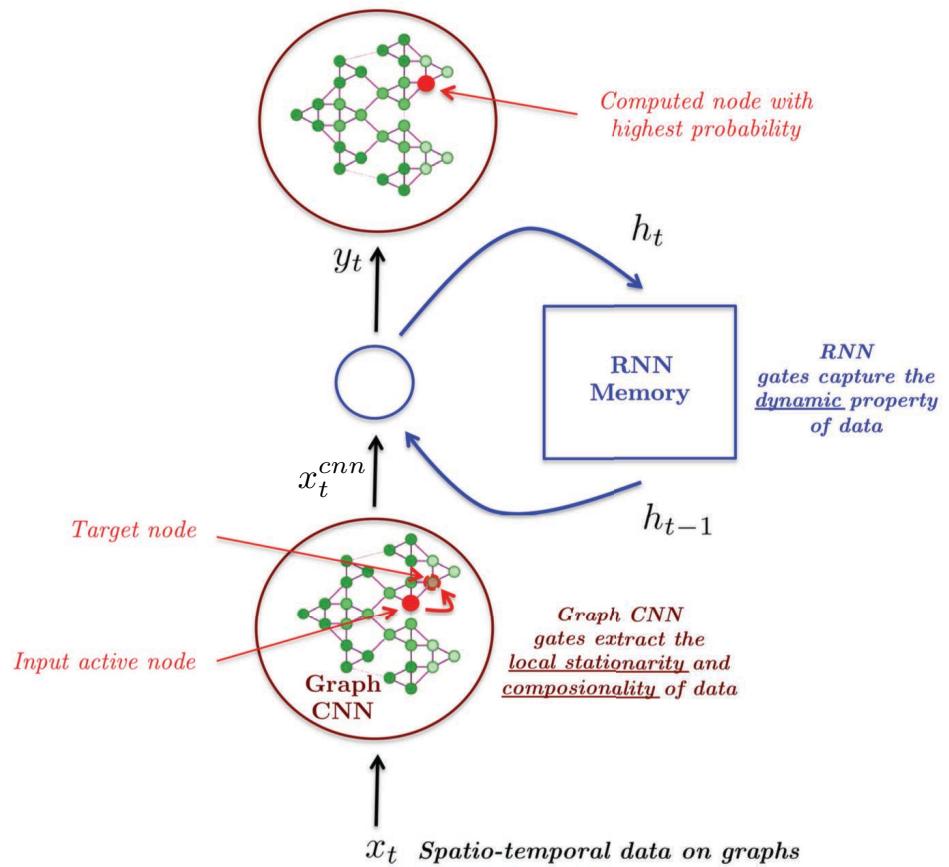
$x_t$  : Input at t

$h_t$  : Short-term memory vector at t

$c_t$  : Long-term memory vector at t

$y_t$  : Probability of next word at t

$x_t^{cnn}$  : Output of CNN at t



# Numerical Experiments

- Benchmark Penn Treebank Corpus:
  - Length of training set = 1 million words.
  - Length of test set = 100,000 words.
- Note: Experiments on T=20 time steps of LSTM.

Architecture	Representation	Parameters	Train Perplexity	Test Perplexity
Zaremba et al. (2014) code <sup>6</sup>	embedding	681,800	36.96	117.29
Zaremba et al. (2014) code <sup>6</sup>	one-hot	34,011,600	53.89	118.82
LSTM	embedding	681,800	48.38	120.90
LSTM	one-hot	34,011,600	54.41	120.16
LSTM, dropout	one-hot	34,011,600	145.59	112.98
GCRN-M1	one-hot	42,011,602	18.49	177.14
GCRN-M1, dropout	one-hot	42,011,602	114.29	<b>98.67</b>

Table 1. Performance (perplexity) on the PTB dataset.

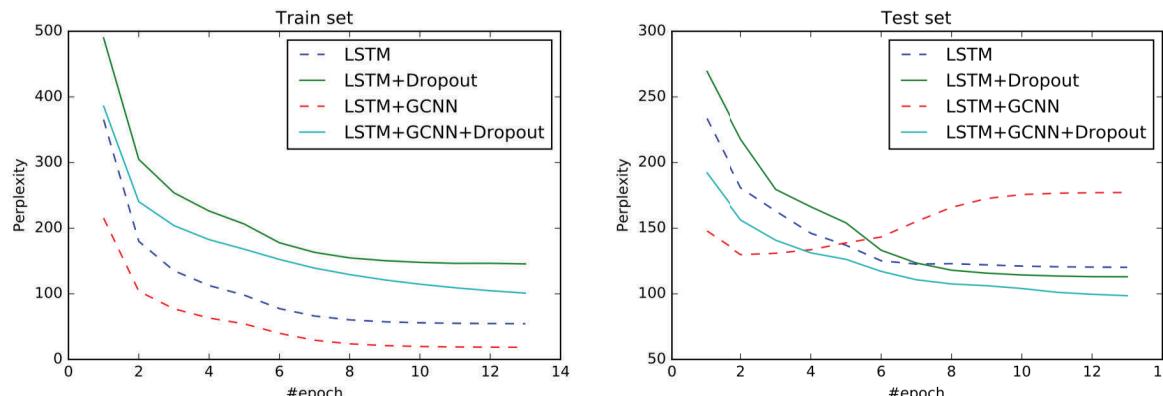


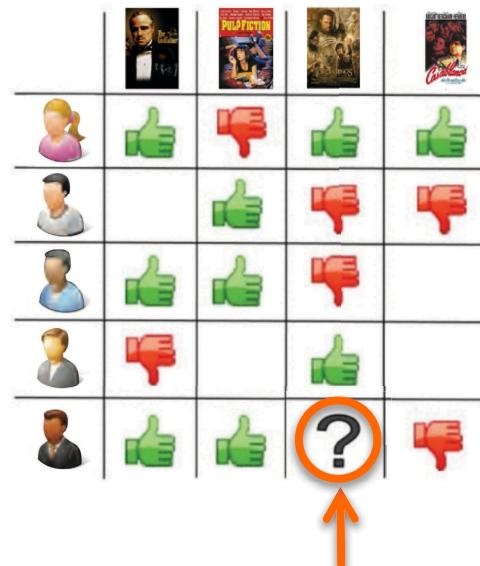
Figure 1. Learning dynamic of LSTM with and without graph structure.

# Outline

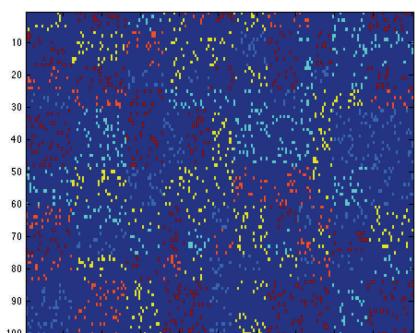
- Data Science and Deep Learning
  - *A brief history*
  - *What makes DL great*
- Convolutional Neural Networks
  - *Architecture*
- Convolutional Neural Networks on Graphs
  - *Non-Euclidean data*
  - *Graph convolution, graph downsampling and pooling*
  - *Numerical experiments*
- Natural Language Processing
- **Recommender Systems**
- Conclusion

# Recommender Systems

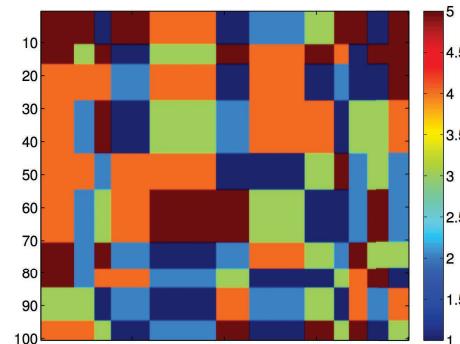
- **Goal:** Predict user ratings for items based on previous ratings.  
Ex. Netflix



- **Formalization:** Given a few ratings/entries  $M_{ij}$  of item  $j$  and user  $i$ , find a low-rank matrix  $X$  that best fits the ratings.



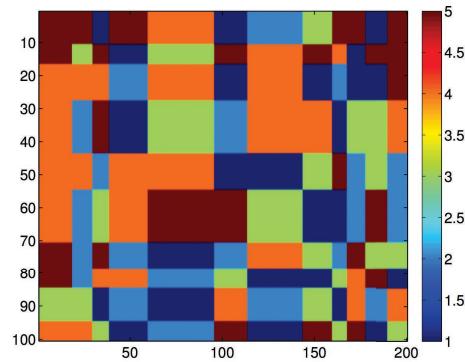
*Recommendation*  
=   
*Matrix completion*



# Low-Rank Assumption

- Assumption: Rating matrix  $X$  is low-rank.

*Definition:* A low-rank matrix has many columns and rows that are linearly dependent.



- For Netflix movie recommendation:

- There exist communities of users who rate movies the same way.*
- There are groups of movies that receive same ratings.*

Same valid assumptions for Amazon (users, products), LinkedIn (users, jobs), Facebook (users, ads), etc.

# Low-Rank Matrix Completion [Candes-Recht'09]

- Formalization:

$$\min_X \text{rank}(X) \text{ s.t. } \underbrace{X_{ij} = M_{ij}, \forall ij \in \Omega}_{\text{No noise in measurements}} \quad \xrightarrow{\hspace{1cm}}$$

↓

$$\min_X \text{rank}(X) + \frac{\lambda}{2} \|\Omega \circ (X - M)\|_F^2$$

$\underbrace{\hspace{1cm}}_{\text{Robustness w.r.t. noise in measurements}}$

- Solution: Combinatorial problem is NP-hard (because of rank).

# Relaxation

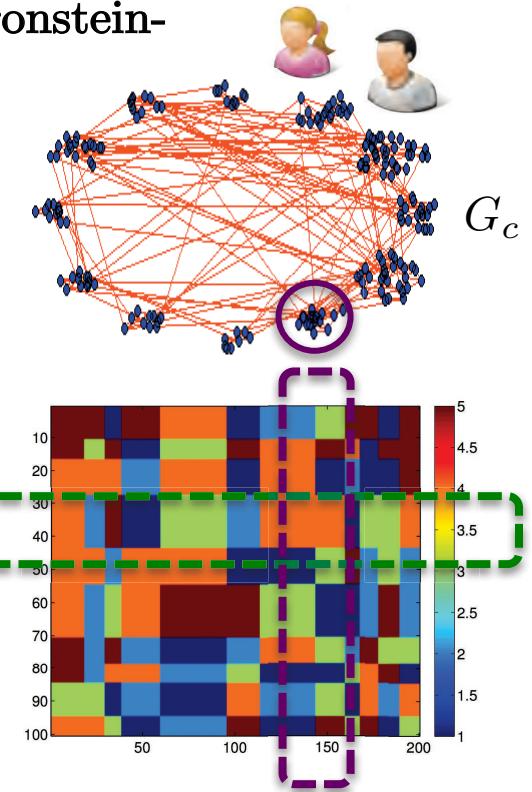
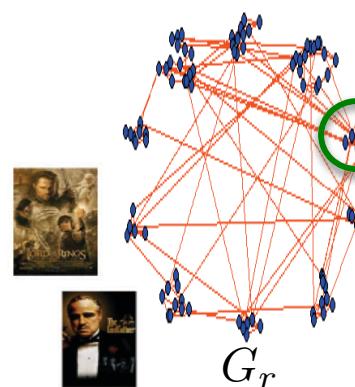
- There exists two classes of rank relaxation:
  1. Nuclear norm relaxation  $\Rightarrow$  **Collaborative filtering recom systems (Netflix)**
  2. Graph relaxation  $\Rightarrow$  **Content filtering recom systems (Amazon)**
- Graph relaxation with graph Dirichlet [Kalofolias-B-Bronstein-Vandergheynst'14]:

$$\min_X \frac{1}{2} \|X\|_{G_c}^2 + \frac{1}{2} \|X\|_{G_r}^2 + \frac{\lambda}{2} \|\Omega \circ (X - M)\|_F^2$$

with  $\|X\|_{G_r}^2 = \text{tr}(X^T L_r X)$

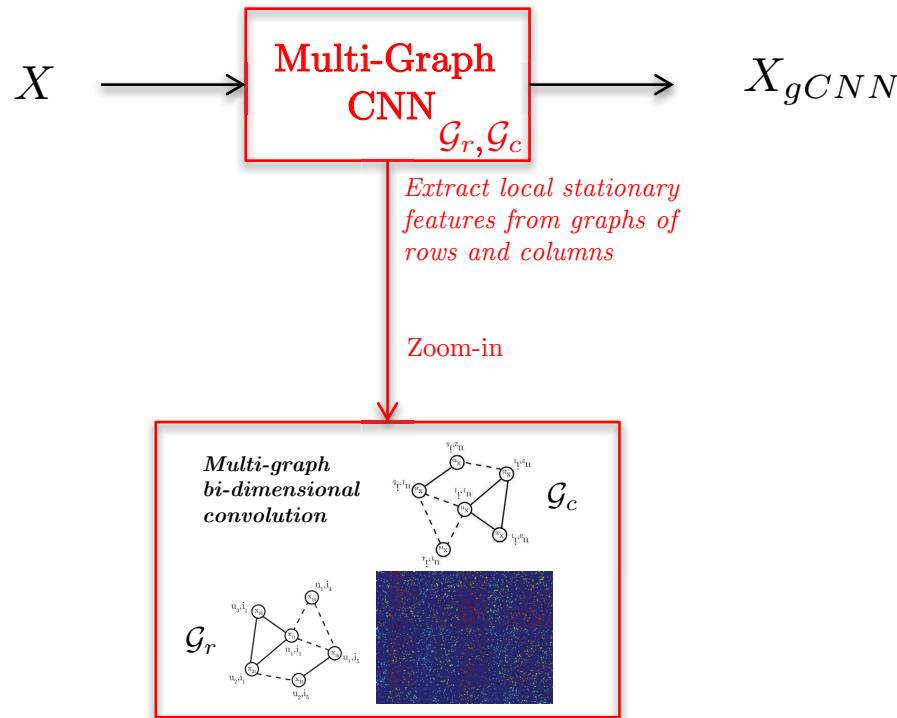
$$\|X\|_{G_c}^2 = \text{tr}(X L_c X^T)$$

$\Rightarrow$  Diffusion term is designed to force the ratings to be smooth on graphs.



# How Graph CNNs can Help?

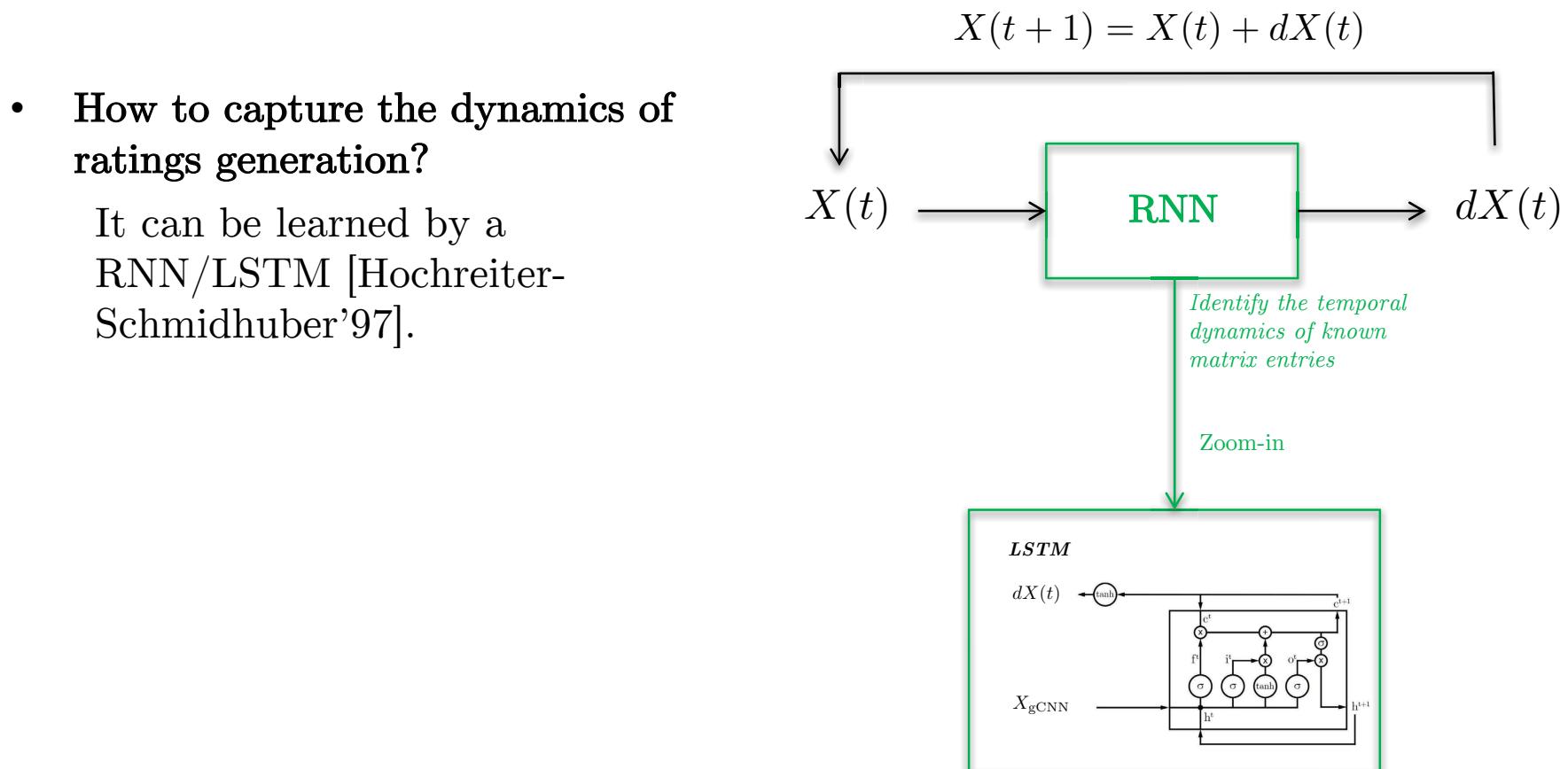
- Graph CNNs can extract highly meaningful patterns on graphs of users and items, and use them for rating predictions.



- Observe that CNNs are designed on multiple graphs.

# Capturing Ratings Dynamics

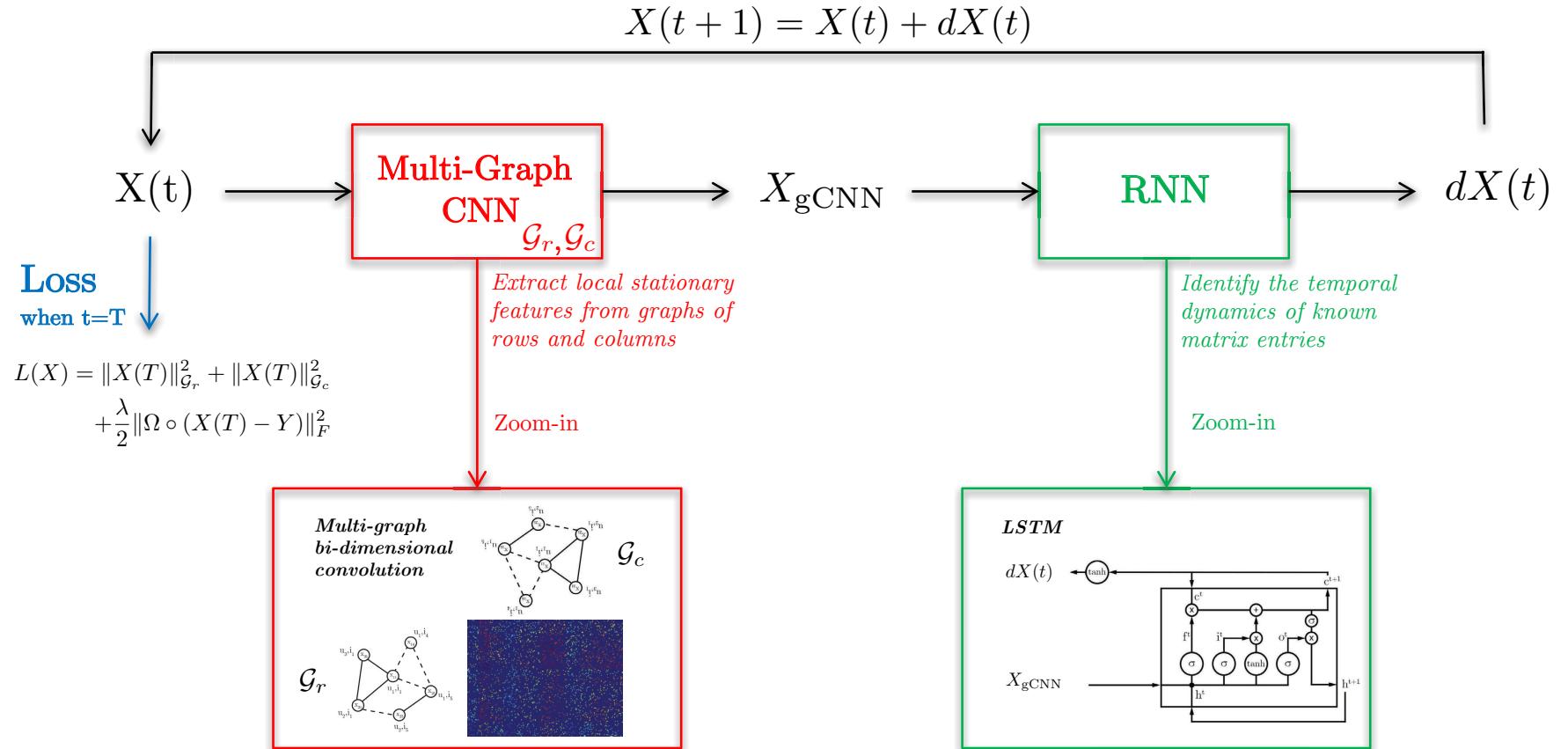
- Assumption: The collected ratings  $M$  are the “snapshot” of a stochastic non-linear diffusion process. For example, movies can go in and out of popularity as triggered by external events such as an actor’s appearance in a new movie or chit-chats on social network [Koren-et-al’09].



# Final Architecture

## [Monti-Bronstein-B'17]

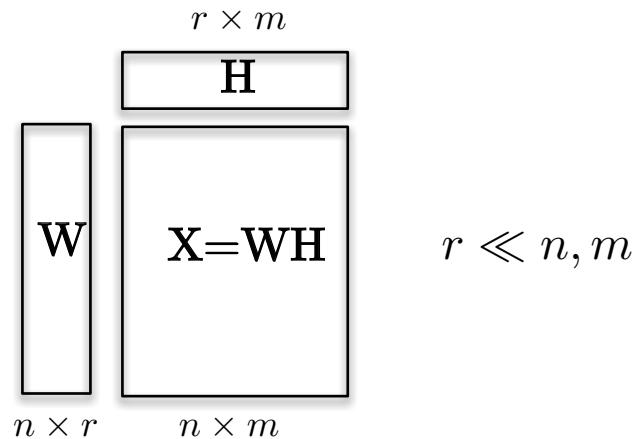
- Proposed recommender system combines gCNN + LSTM:



- If  $X \in \mathbb{R}^{m \times n}$ , then learning complexity is  $O(m.n) \Rightarrow$  Quadratic, no scale up.

# Factorized Models

- Factorized matrices:



- New loss function:

$$\min_X \frac{1}{2} \|W\|_{G_r}^2 + \frac{1}{2} \|H\|_{G_c}^2 + \frac{\lambda}{2} \|\Omega \circ (WH - M)\|_F^2$$

with  $\|W\|_{G_r}^2 = \text{tr}(W^T L_r W)$   
 $\|H\|_{G_c}^2 = \text{tr}(H L_c H^T)$

- If  $X \in \mathbb{R}^{m \times n}$ , then learning complexity is  $O(m+n) \Rightarrow$  Linear

# Numerical Experiments

- We evaluated the system on benchmark datasets.

METHOD	RMSE
GLOBAL MEAN	1.154
USER MEAN	1.063
MOVIE MEAN	1.033
MC (CANDES & RECHT, 2012)	0.973
IMC (JAIN & DHILLON, 2013; XU ET AL., 2013)	1.653
GMC (KALOFOLIAS ET AL., 2014)	0.996
GRALS (RAO ET AL., 2015)	0.945
<b>sRGCNN</b>	<b>0.929</b>

Table 1. Performance (RMS error) on the MovieLens dataset.

METHOD	Flixster	Douban	YahooMusic
GRALS	1.3126 / 1.2447	0.8326	38.0423
<b>sRGCNN</b>	<b>1.1788 / 0.9258</b>	<b>0.8012</b>	<b>22.4149</b>

Table 2. Performance (RMS error) on Flixster, Douban and YahooMusic.

# Outline

- Data Science and Deep Learning
  - *A brief history*
  - *What makes DL great*
- Convolutional Neural Networks
  - *Architecture*
- Convolutional Neural Networks on Graphs
  - *Non-Euclidean data*
  - *Graph convolution, graph downsampling and pooling*
  - *Numerical experiments*
- Natural Language Processing
- Recommender Systems
- Conclusion

# Conclusion

- **Proposed contributions [NIPS'16]:**
  1. *Generalization of CNNs to non-Euclidean data*
  2. *Exact localized filters on graphs*
  3. *Same linear learning complexity as CNNs while being universal to any graph*
  4. *GPU implementation*
- **Several potential applications including**  
Biological networks (gene, brain connectivity)

# Workshop Announcement

## *New Deep Learning Techniques*

*Feb 5-9, 2018*

### *Organizers*

*Xavier Bresson, NTU*

*Michael Bronstein, USI/Intel*

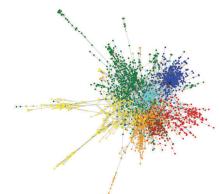
*Joan Bruna, Berkeley*

*Yann LeCun, Facebook FAIR Director*

*Stanley Osher, UCLA*

*Arthur Szlam, Facebook*

The screenshot shows the IPAM website's 'Workshops' section. At the top, there is a navigation bar with links for PROGRAMES, VIDEOS, NEWS, PEOPLE, YOUR VISIT, ABOUT IPAM, DONATE, and CONTACT US. Below the navigation bar, the word 'Workshops' is prominently displayed in white text on a dark blue background. Underneath, a sub-menu shows 'Programs > Workshops > New Deep Learning Techniques'. The main content area features a large, artistic graphic of concentric circles and mathematical symbols like ' $\lambda$ ' and ' $\nabla$ ' on a blue background. Below the graphic, the title 'New Deep Learning Techniques' is shown in bold black text, followed by the date 'FEBRUARY 5 - 9, 2018'. A horizontal navigation bar at the bottom includes links for 'OVERVIEW' (which is highlighted in blue), 'SPEAKER LIST', 'LOGGING', and 'APPLICATION & REGISTRATION'.



<http://www.ipam.ucla.edu/programs/workshops/new-deep-learning-techniques>

