

CE7454 : Deep Learning for Data Science

Lecture 14: Graph Neural Networks

Semester 1 2019/20

Xavier Bresson

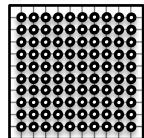
School of Computer Science and Engineering
Data Science and AI Research Centre
Nanyang Technological University (NTU), Singapore



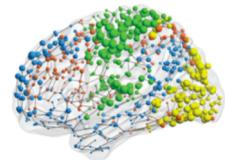
NANYANG
TECHNOLOGICAL
UNIVERSITY
SINGAPORE

Outline

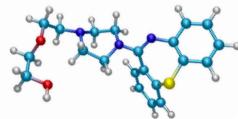
- Part 1: Traditional ConvNets



- Part 2: Spectral Graph ConvNets



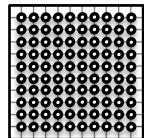
- Part 3: Spatial Graph ConvNets



- Conclusion

Outline

- Part 1: Traditional ConvNets



- Part 2: Spectral Graph ConvNets



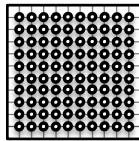
- Part 3: Spatial Graph ConvNets



- Conclusion

Outline

- Part 1: Traditional ConvNets
 - Architecture
 - Graph Structured Data



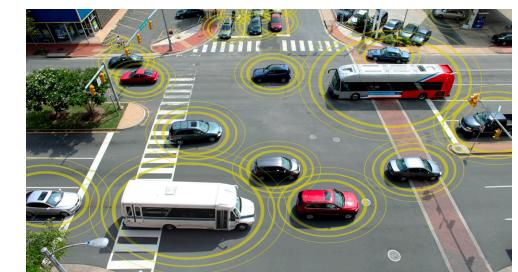
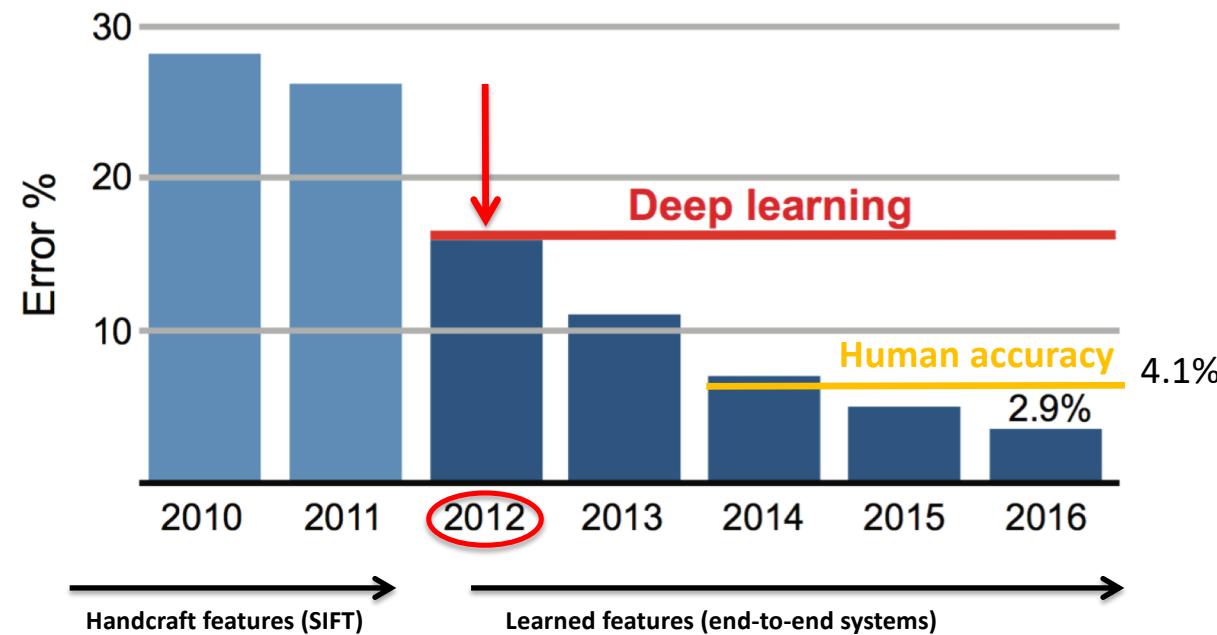
ConvNets

- A breakthrough in **image recognition** :

LeCun, Bottou, Bengio, Haffner 1998

Krizhevsky, Sutskever Hinton, 2012

IMAGENET



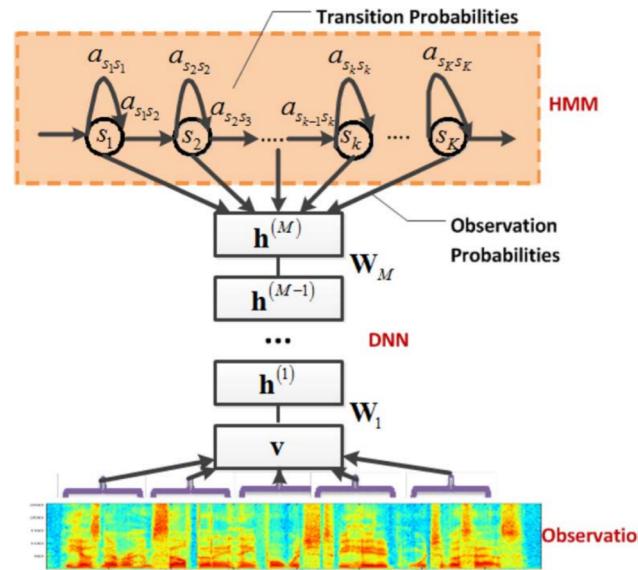
ConvNets

- A breakthrough in speech recognition :

Dahl, Yu, Deng, Acero, 2010

Hinton, Deng, Yu, Dahl et al. 2012

Acoustic model	Recog \ WER	RT03S FSH	Hub5 SWB
Traditional features	1-pass -adapt	27.4	23.6
Deep Learning	1-pass -adapt	18.5	16.1



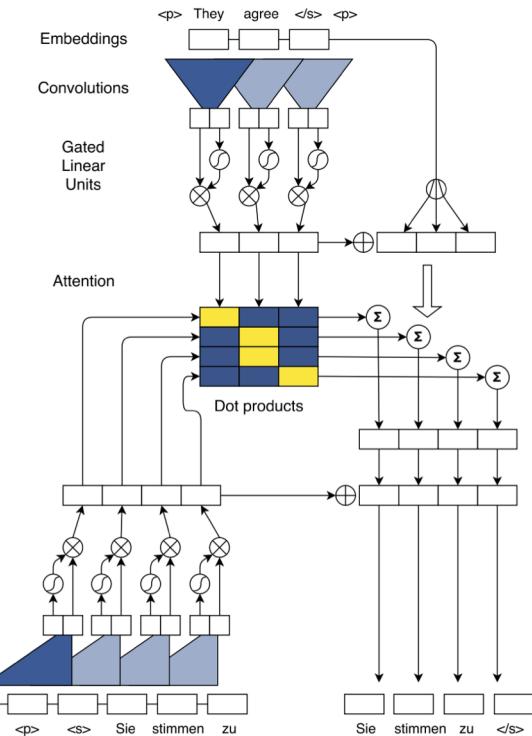
ConvNets

- A breakthrough in **Machine Translation** :
- Gehring, Auli, Grangier, Yarats, Dauphin, 2017

WMT'14 English-German	BLEU
Wu et al. (2016) GNMT	26.20
Wu et al. (2016) GNMT + RL	26.30
ConvS2S	26.43

WMT'14 English-French	BLEU
Zhou et al. (2016)	40.4
Wu et al. (2016) GNMT	40.35
Wu et al. (2016) GNMT + RL	41.16
ConvS2S	41.44
ConvS2S (10 models)	41.62

	BLEU	Time (s)
GNMT GPU (K80)	31.20	3,028
GNMT CPU 88 cores	31.20	1,322
GNMT TPU	31.21	384
ConvS2S GPU (K40) $b = 1$	33.45	327
ConvS2S GPU (M40) $b = 1$	33.45	221
ConvS2S GPU (GTX-1080ti) $b = 1$	33.45	142
ConvS2S CPU 48 cores $b = 1$	33.45	142
ConvS2S GPU (K40) $b = 5$	34.10	587
ConvS2S CPU 48 cores $b = 5$	34.10	482
ConvS2S GPU (M40) $b = 5$	34.10	406
ConvS2S GPU (GTX-1080ti) $b = 5$	34.10	256

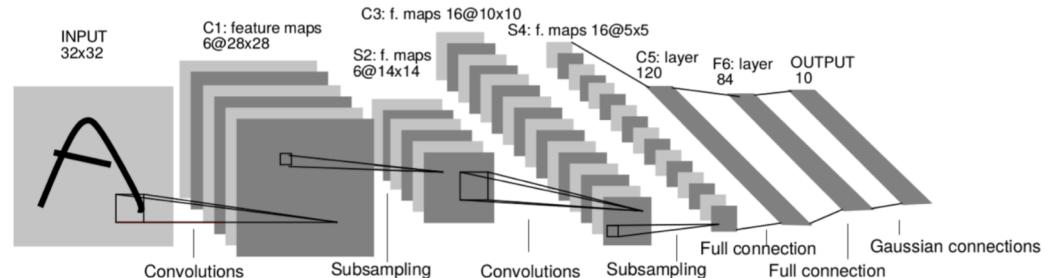


Oliver Schröbe [DE] Ich verstehe nicht, weshalb sich alle plötzlich darüber aufregen, dass wir Steuergelder nach Griechenland "verschwendeten". Als wären Steuergelder bisher sinnvoll angelegt worden... May 5 at 10:06am via TweetDeck · Comment · Like · Translate

Oliver Schröbe [DE] I do not understand why all of a sudden get e about the fact that we are wasting taxpayers' money to Greece. W tax money would have been useful to create ... May 5 at 10:06am via TweetDeck · Comment · Like · Untranslate

ConvNets

- An architecture for **high-dimensional learning** :

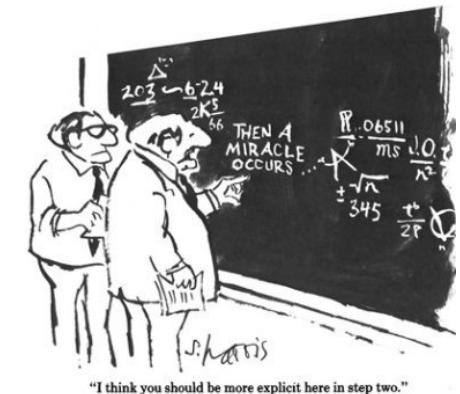


- **Curse of dimensionality** :

$$\text{dim(image)} = 1024 \times 1024 \approx 10^6$$

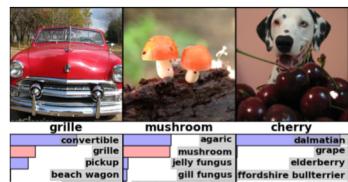
For $N=10$ samples/dim $\Rightarrow 10^{1,000,000}$ points

- ConvNets are **powerful** to solve high-dimensional learning problems.



ConvNets

- Main assumption :
 - Data (image, video, sound) is **compositional**, it is formed of patterns that are:
 - Local (Hubel-Wiesel 1962)
 - Stationary (shared patterns)
 - Hierarchical (multi-scale)
 - ConvNets **leverage the compositionality structure** :
 - They extract compositional features and feed them to classifier, recommender, etc (end-to-end).



Computer Vision

```
    static inline void audit_dupe_lsm_fixup(struct audit_field *sf)
{
    if (sf->type == SF_TYPE_LSM)
        audit_dupe_lsm_fixup(sf->lsm_sf, GFP_KERNEL);
    return;
}

/*-----*/
void __audit_lsm_fixup(struct audit_field *sf, const char *msg)
{
    if (sf->type == SF_TYPE_LSM)
        __audit_lsm_fixup(sf->lsm_sf, msg);
}
```

NLP



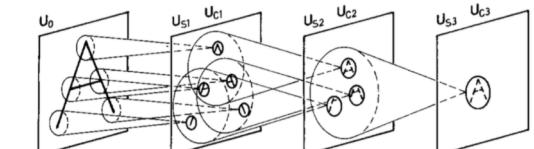
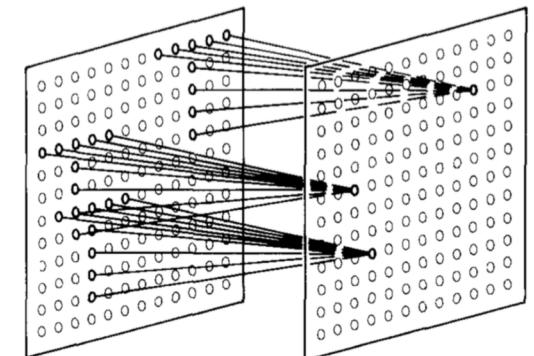
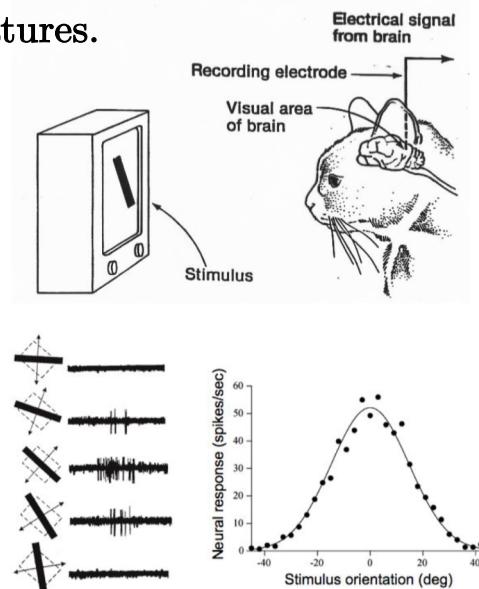
Speech



Games

Key Properties

- **Locality :**
 - Property inspired by the human visual cortex system.
- **Local receptive fields (Hubel, Wiesel 1962) :**
 - Activate in the presence of **local** features.



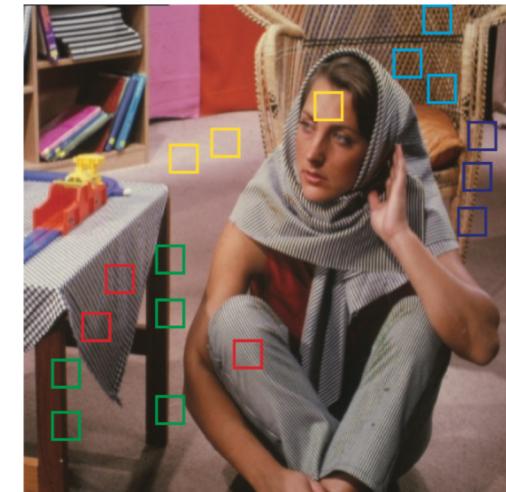
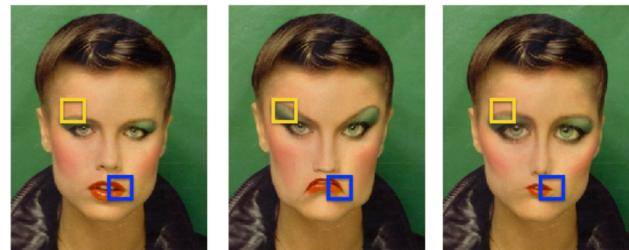
Neocognitron
Fukushima 1980

Key Properties

- Stationarity \Leftrightarrow Translation invariance
 - Global invariance

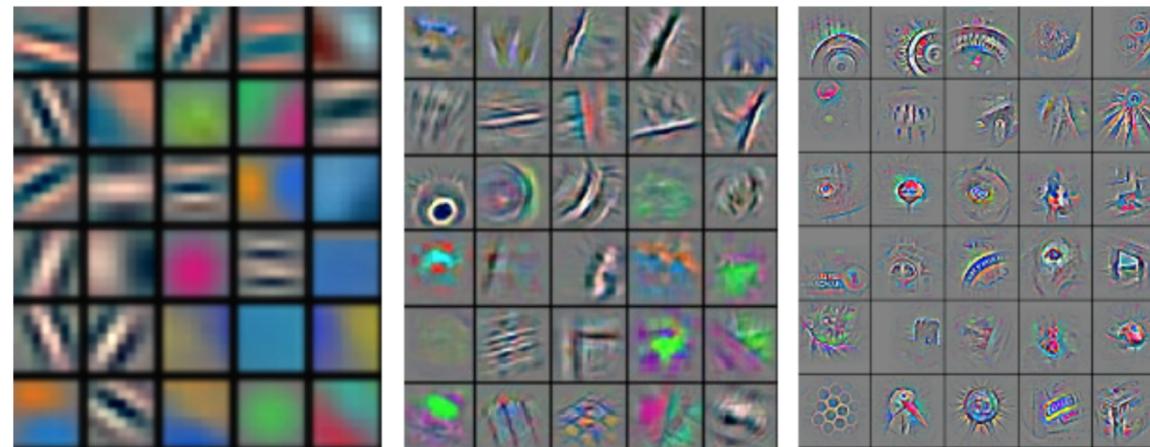


- Local stationarity \Leftrightarrow Similar patches are shared across the data domain
 - Local invariance, essential for intra-class variations



Key Properties

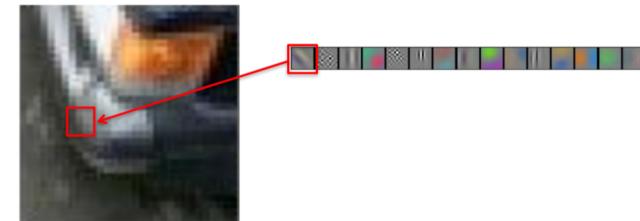
- Multi-scale :
 - Simple structures combine to compose slightly more abstract structures, and so on, in a hierarchical way.
- Inspired by brain visual primary cortex (V1 and V2 neurons).



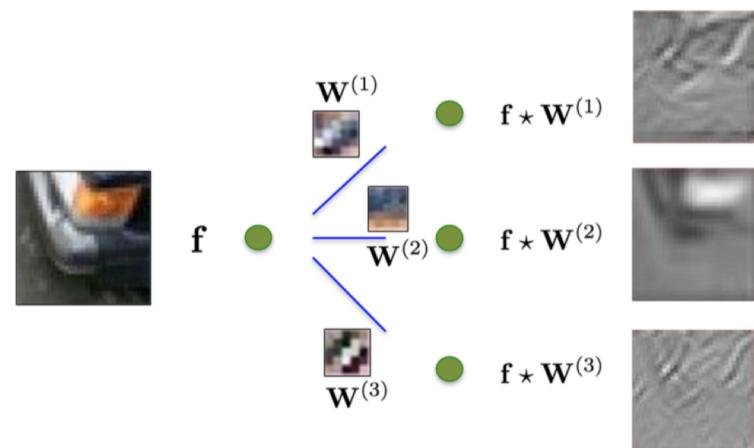
Features learned by ConvNet become increasingly more complex at deeper layers
(Zeiler, Fergus 2013)

Learning Complexity

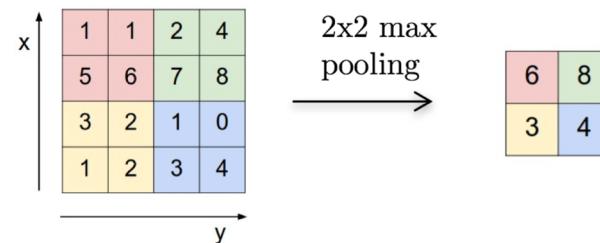
- **Locality :**
 - Compact support kernels $\Rightarrow O(1)$ parameters per filter.



- **Stationarity :**
 - Convolutional operators $\Rightarrow O(n \log n)$ in general (FFT) and $O(n)$ for compact kernels.

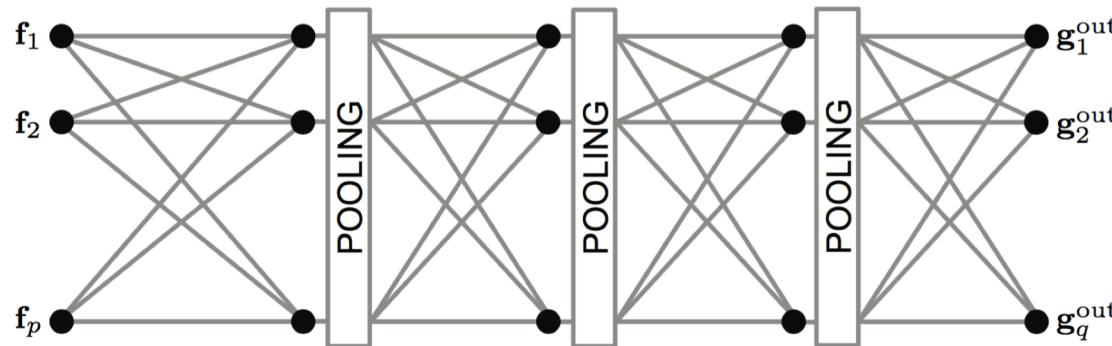


- **Multi-scale :**
 - Downsampling + pooling $\Rightarrow O(n)$



Compositional layers

$$\begin{aligned}\mathbf{f}_l &= l\text{-th image feature (R,G,B channels), } \dim(\mathbf{f}_l) = n \times 1 \\ \mathbf{g}_l^{(k)} &= l\text{-th feature map, } \dim(\mathbf{g}_l^{(k)}) = n_l^{(k)} \times 1\end{aligned}$$



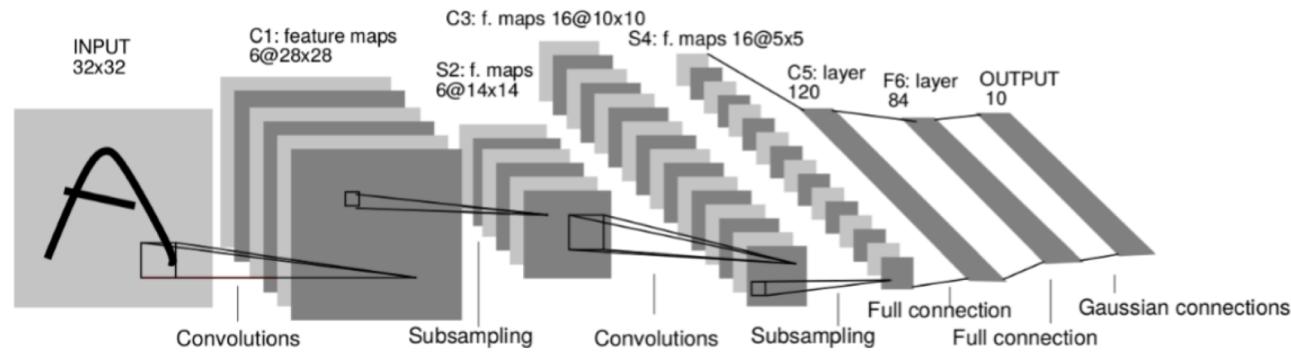
Compositional features consist of multiple convolutional + pooling layers.

Convolutional layer $\mathbf{g}_l^{(k)} = \xi \left(\sum_{l'=1}^{q_{k-1}} \mathbf{W}_{l,l'}^{(k)} \star \xi \left(\sum_{l'=1}^{q_{k-2}} \mathbf{W}_{l,l'}^{(k-1)} \star \xi \left(\dots \mathbf{f}_{l'} \right) \right) \right)$

Activation, e.g. $\xi(x) = \max\{x, 0\}$ rectified linear unit (ReLU)

Pooling $\mathbf{g}_l^{(k)}(x) = \|\mathbf{g}_l^{(k-1)}(x') : x' \in \mathcal{N}(x)\|_p \quad p = 1, 2, \text{ or } \infty$

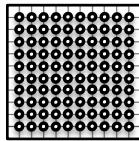
ConvNets



- ☺ Filters localized in space (**locality**)
- ☺ Convolutional filters (**stationarity**)
- ☺ Multiple layers (**multi-scale**)
- ☺ $O(1)$ parameters per filter (independent of input image size n)
- ☺ $O(n)$ complexity per layer (filtering done in the spatial domain)

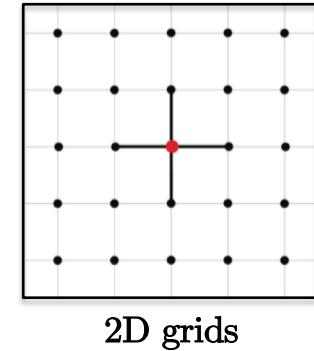
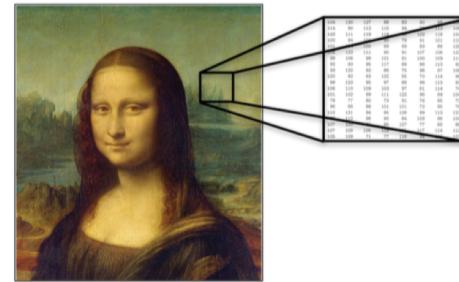
Outline

- Part 1: Traditional ConvNets
 - Architecture
 - Graph Structured Data



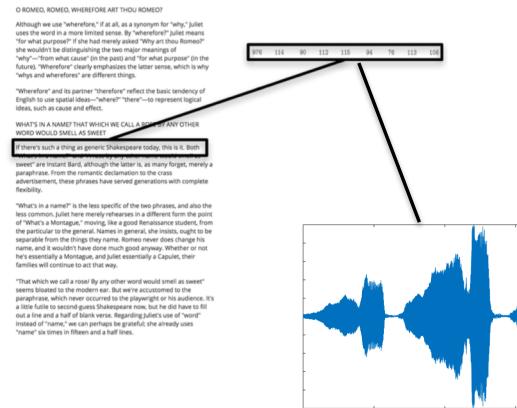
Data Domain

- Image, volume, video lie on
2D, 3D, 2D+1 Euclidean domains



2D grids

- Sentence, word, sound lie on
1D Euclidean domain



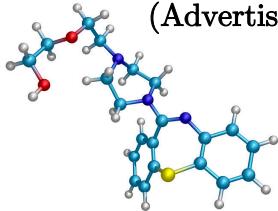
1D grid

- These domains have strong regular spatial structures.
 - All ConvNet operations are mathematically well defined and fast (convolution, pooling).

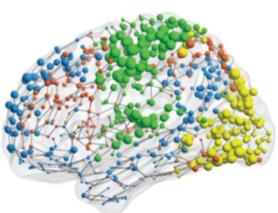
Graph Structured Data



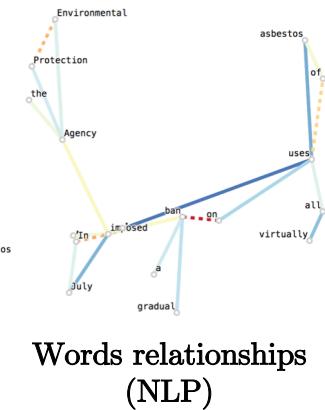
Social networks
(Advertisement)



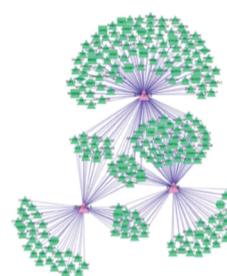
Drug/Material
molecules
(Chemistry)



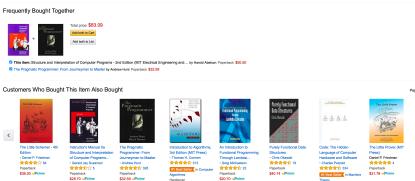
Brain
connectivity
(Neuroscience)



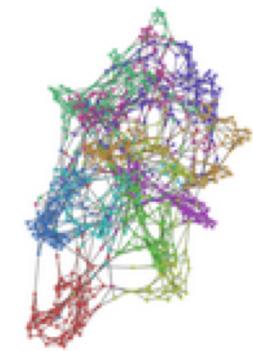
Words relationships
(NLP)



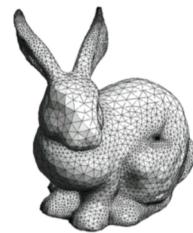
Gene Regulatory
Network



Recommender
systems (Amazon,
Netflix)



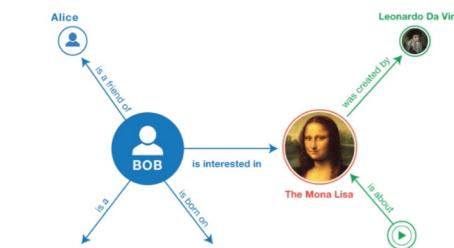
Graphs/
Networks



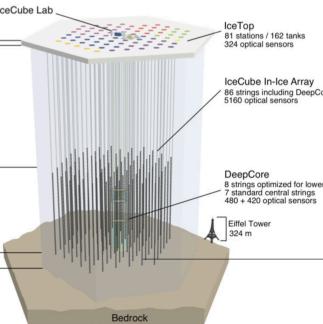
3D Meshes
(Computer Graphics)



Transportation
networks



Knowledge graph
(Causality)



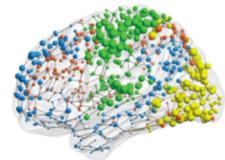
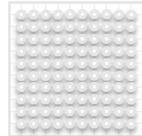
Neutrino
detection (High-
energy Physics)

New Challenges

- How to extend ConvNets to graph structured data ?
- Assumption :
 - Graph structured data is compositional :
 - Locally stationary and manifest hierarchical structures.
- How to define compositionality on graphs ? (convolution/pooling on graphs)
- How to make them fast? (linear complexity)

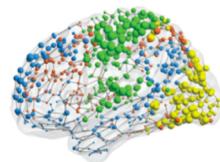
Outline

- Part 1: Traditional ConvNets
- Part 2: Spectral Graph ConvNets
- Part 3: Spatial Graph ConvNets
- Conclusion



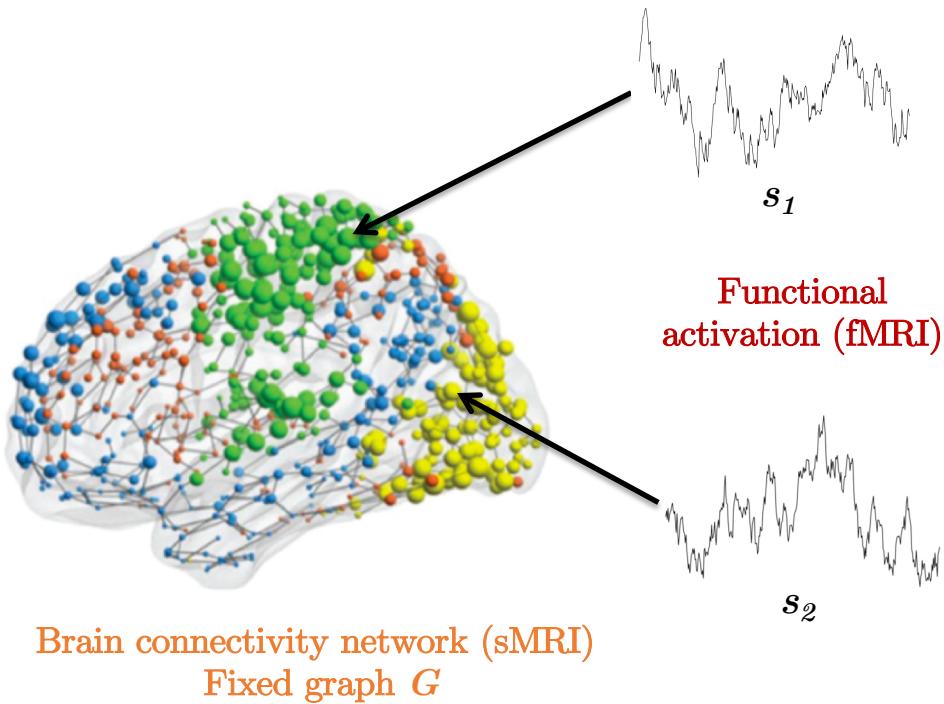
Outline

- Part 2: Spectral Graph ConvNets
 - Problem Setting
 - Graph Theory
 - Spectral Graph Convolution
 - Graph Pooling
 - Spectral GCNs
 - Applications

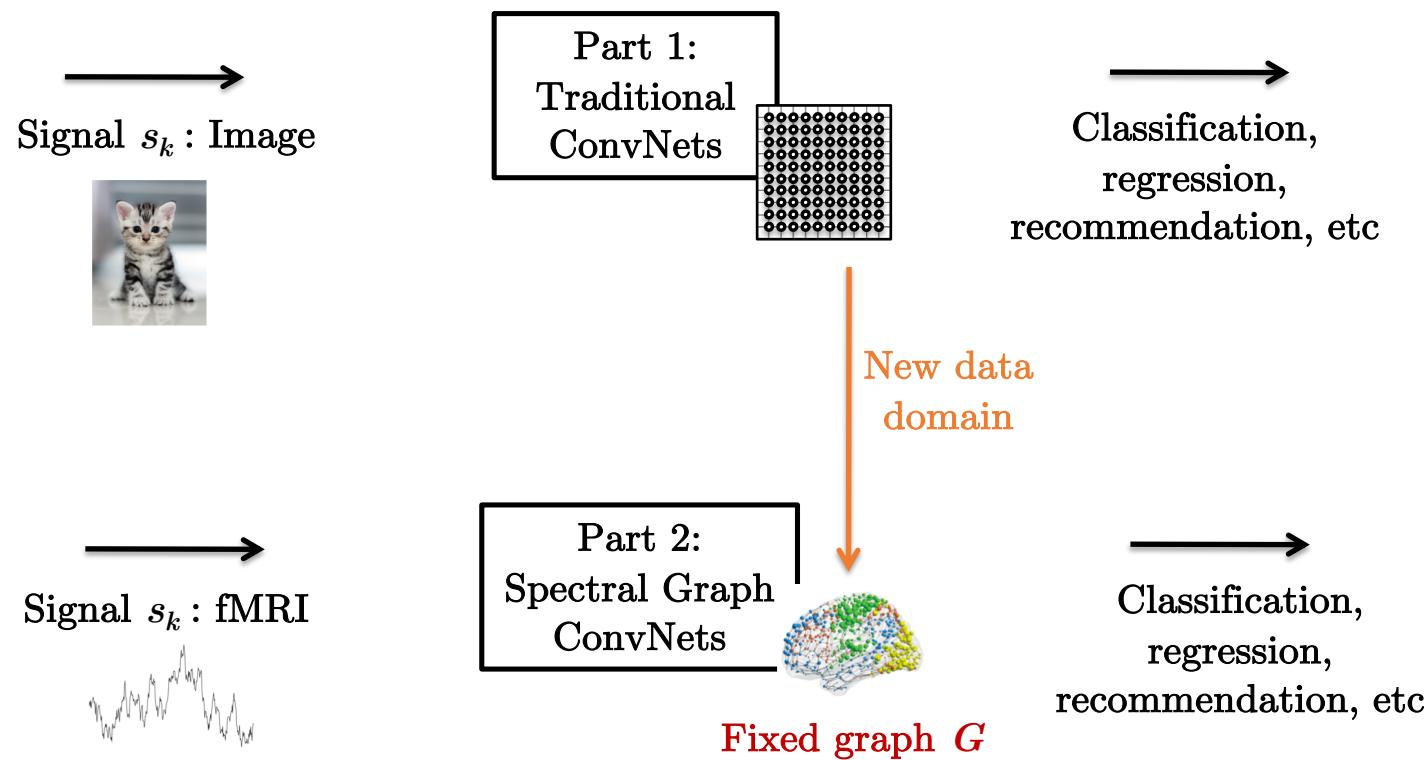


Problem Setting

- Given **fixed** arbitrary graph(s) G , and a set of signals s_k on G to be analyzed :



Spectral ConvNets

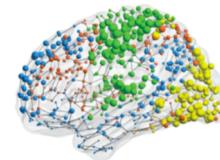


Generalization challenges

- Convolution and down-sampling (for pooling) must be generalized from Euclidean grid domains to graphs.
- How ?
 - Spectral graph theory allows to redefine convolution in the context of graphs with Fourier analysis.
 - Graph theory provide graph clustering techniques to reformulate down-sampling for graphs.

Outline

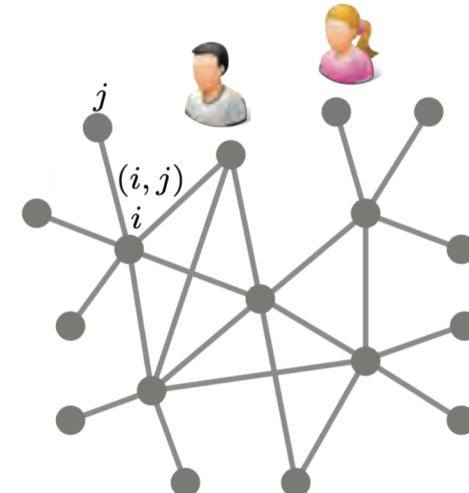
- Part 2: Spectral Graph ConvNets
 - Problem Setting
 - Graph Theory
 - Spectral Graph Convolution
 - Graph Pooling
 - Spectral GCNs
 - Applications



Graphs

- Graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$
- Vertices $\mathcal{V} = \{1, \dots, n\}$
- Edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$
- Vertex weights $b_i > 0$ for $i \in \mathcal{V}$
- Edge weights $a_{ij} \geq 0$ for $(i, j) \in \mathcal{E}$
- Vertex fields $L^2(\mathcal{V}) = \{f : \mathcal{V} \rightarrow \mathbb{R}^h\}$
Represented as $\mathbf{f} = (f_1, \dots, f_n)$
- Hilbert space with inner product

$$\langle f, g \rangle_{L^2(\mathcal{V})} = \sum_{i \in \mathcal{V}} a_i f_i g_i$$



Graph Laplacian

- **Laplacian operator** $\Delta : L^2(\mathcal{V}) \rightarrow L^2(\mathcal{V})$

$$(\Delta f)_i = \frac{1}{b_i} \sum_{j:(i,j) \in \mathcal{E}} a_{ij} (f_i - f_j)$$

difference between f and its local average (2nd derivative on graphs)

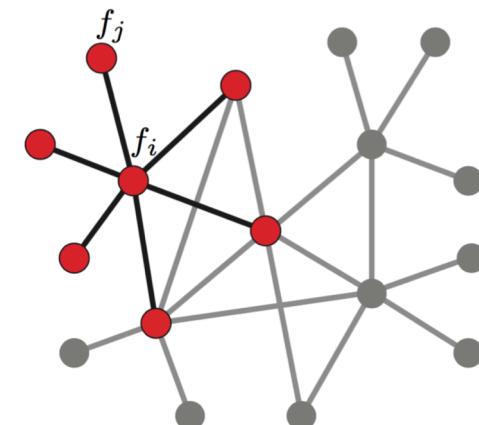
- **Core operator** in spectral graph theory.
- Represented as a **positive semi-definite** $n \times n$ matrix

- **Un-normalized** Laplacian $\Delta = \mathbf{D} - \mathbf{A}$

- **Normalized** Laplacian $\Delta = \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$

- **Random walk** Laplacian $\Delta = \mathbf{I} - \mathbf{D}^{-1} \mathbf{A}$

where $\mathbf{A} = (a_{ij})$ and $\mathbf{D} = \text{diag}(\sum_{j \neq i} a_{ij})$



Spectral Decomposition

- A Laplacian of a graph of n vertices admits n eigenvectors:

$$\Delta \phi_k = \lambda_k \phi_k, \quad k = 1, 2, \dots$$

- Eigenvectors are **real** and **orthonormal** ($\langle \phi_k, \phi_{k'} \rangle_{L^2(\mathcal{V})} = \delta_{kk'}$)
 - Eigenvalues are **non-negative** ($0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$)
 - Laplacian eigenvectors are also called **Fourier basis functions/modes**.
-
- Eigen-decomposition of graph Laplacian:

$$\Delta = \Phi^T \Lambda \Phi$$

where $\Phi = (\phi_1, \dots, \phi_n)$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$

Interpretation

- Find the **smoothest** orthonormal basis $\Phi = (\phi_1, \dots, \phi_n)$ on a graph

$$\begin{aligned} \min_{\phi_k} E_{\text{Dir}}(\phi_k) \quad & \text{s.t.} \quad \|\phi_k\| = 1, \quad k = 2, 3, \dots, n \\ & \phi_k \perp \text{span}\{\phi_1, \dots, \phi_{k-1}\} \end{aligned}$$

where E_{Dir} is the **Dirichlet** energy = measure of smoothness of a function

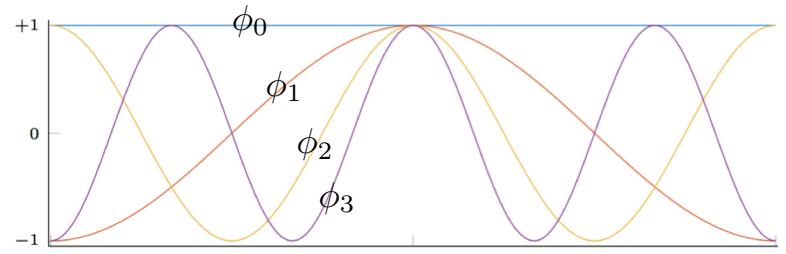
$$E_{\text{Dir}}(\mathbf{f}) = \mathbf{f}^\top \Delta \mathbf{f}$$

- Solution :
 - Smallest n Laplacian eigenvectors :

$$\begin{aligned} \min_{\Phi \in \mathbb{R}^{n \times n}} \underbrace{\text{trace}(\Phi^\top \Delta \Phi)}_{\|\Phi\|_{\mathcal{G}} \text{ Dirichlet norm}} \quad & \text{s.t.} \quad \Phi^\top \Phi = \mathbf{I} \end{aligned}$$

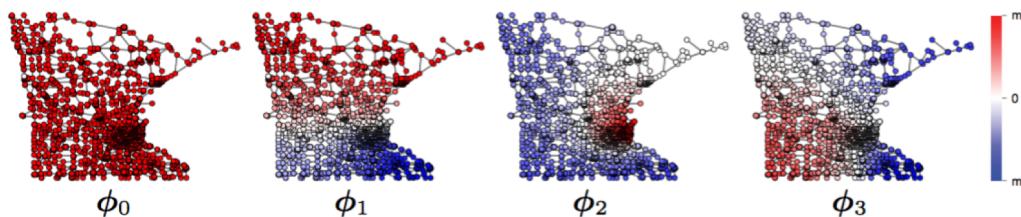
Fourier Modes

- Euclidean domain :



First eigenvectors of 1D Euclidean Laplacian = standard Fourier basis

- Graph domain :



First Laplacian eigenvectors of a graph

Laplacian eigenvectors related to graph geometry!

(s.a. communities, hubs, etc)

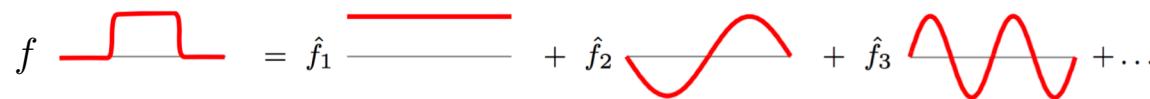
Spectral Graph Clustering^[1]

[1] Von Luxburg, A tutorial on spectral clustering, 2007

Euclidean Fourier analysis

- A function $f : [-\pi, \pi] \rightarrow \mathbb{R}$ can be written as **Fourier series** :

$$f(x) = \sum_{k \geq 0} \underbrace{\frac{1}{2\pi} \int_{-\pi}^{\pi} f(x') e^{-ikx'} dx'}_{\hat{f}_k = \langle f, e^{-ikx} \rangle_{L^2([-\pi, \pi])}} e^{-ikx}$$



- Fourier basis e^{-ikx} = **Laplace-Beltrami eigenfunctions** :

$$-\Delta \phi_k = k^2 \phi_k$$

$$\begin{cases} \phi_k &= \text{Fourier mode} \\ k &= \text{frequency of Fourier mode} \end{cases}$$

Fourier analysis on graphs

- A function $f : \mathcal{V} \rightarrow \mathbb{R}$ can be written as Fourier series^[1]:

$$f_i = \sum_{k=1}^n \underbrace{\langle f, \phi_k \rangle_{L^2(\mathcal{V})}}_{\hat{f}_k} \phi_{k,i}$$

- \hat{f}_k is the k -th graph Fourier coefficient.
- In matrix-vector notation, with the $n \times n$ Fourier matrix $\Phi = [\phi_1, \dots, \phi_n]$

↑ $\hat{f} = \Phi^\top f$ and $f = \Phi \hat{f}$ ↑
 Fourier transform Inverse Fourier transform

- Graph Fourier basis ϕ_k = Laplacian eigenvectors :

$$\Delta \phi_k = \lambda_k \phi_k \quad \text{with} \quad \begin{cases} \phi_k &= \text{graph Fourier mode} \\ \lambda_k &= \text{(square) frequency} \end{cases}$$

[1] K. Hammond, P. Vandergheynst, R. Gribonval, Wavelets on graphs via spectral graph theory, 2011

Convolution in Continuous Euclidean Space

- Given two functions $f, g : [-\pi, \pi] \rightarrow \mathbb{R}$ their **convolution** is a function

$$(f \star g)(x) = \int_{-\pi}^{\pi} f(x')g(x - x')dx'$$

- Shift-invariance :** $f(x - x_0) \star g(x) = (f \star g)(x - x_0)$

- Convolution theorem :**

- Convolution can be computed in the Fourier domain as

$$\widehat{(f \star g)} = \hat{f} \cdot \hat{g}$$

- Efficient computation :**

- FFT complexity is $O(n \log n)$

Convolution in Discrete Euclidean Space

- **Convolution** of two vectors $\mathbf{f} = (f_1, \dots, f_n)^\top$ and $\mathbf{g} = (g_1, \dots, g_n)^\top$

$$(\mathbf{f} \star \mathbf{g})_i = \sum_m g_{(i-m) \text{ mod } n} \cdot f_m$$

$$\mathbf{f} \star \mathbf{g} = \underbrace{\begin{bmatrix} g_1 & g_2 & \dots & \dots & g_n \\ g_n & g_1 & g_2 & \dots & g_{n-1} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ g_3 & g_4 & \dots & g_1 & g_2 \\ g_2 & g_3 & \dots & \dots & g_1 \end{bmatrix}}_{\text{Circulant matrix}} \begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix}$$

diagonalised by Fourier basis (Toeplitz)

$$= \Phi \begin{bmatrix} \hat{g}_1 & & \\ & \ddots & \\ & & \hat{g}_n \end{bmatrix} \Phi^\top \mathbf{f} = \Phi (\Phi^\top \mathbf{g} \circ \Phi^\top \mathbf{f})$$

Convolution on Graphs

- Spectral convolution of $f, g \in L^2(\mathcal{V})$ can be defined by analogy^[1]:

$$(f \star g)_i = \underbrace{\sum_{k \geq 1} \underbrace{\langle f, \phi_k \rangle_{L^2(\mathcal{V})} \langle g, \phi_k \rangle_{L^2(\mathcal{V})}}_{\text{product in the Fourier domain}}}_{\text{inverse Fourier transform}} \phi_{k,i}$$

- In matrix-vector notation

$$\begin{aligned} \mathbf{f} \star \mathbf{g} &= \Phi (\Phi^\top \mathbf{g} \circ \Phi^\top \mathbf{f}) \\ &= \underbrace{\Phi \text{diag}(\hat{g}_1, \dots, \hat{g}_n) \Phi^\top}_{\mathbf{G}} \mathbf{f} \\ &= \Phi \hat{g}(\Lambda) \Phi^\top \mathbf{f} = \hat{g}(\Phi \Lambda \Phi^\top) \mathbf{f} = \hat{g}(\Delta) \mathbf{f} \end{aligned}$$

- Not shift-invariant (\mathbf{G} has no circulant structure)
- Filter coefficients depend on basis ϕ_1, \dots, ϕ_n
- Expensive computation (no FFT): $O(n^2)$

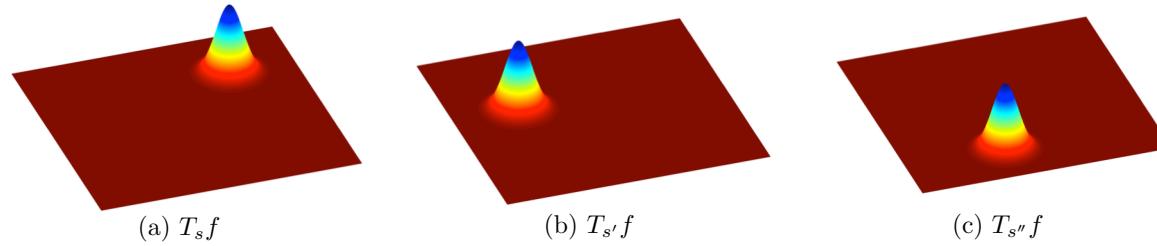
[1] K. Hammond, P. Vandergheynst, R. Gribonval, Wavelets on graphs via spectral graph theory, 2011

No Shift Invariance on Graphs

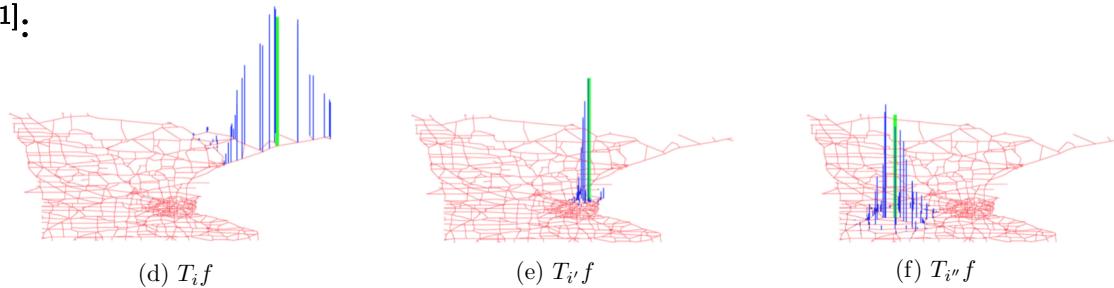
- A signal f on graph can be translated to vertex i as follows :

$$T_i \mathbf{f} = \mathbf{f} \star \delta_i$$

- Euclidean domain :

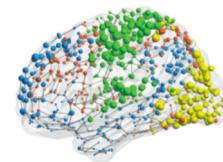


- Graph domain^[1]:



Outline

- Part 2: Spectral Graph ConvNets
 - Problem Setting
 - Graph Theory
 - Spectral Graph Convolution
 - Graph Pooling
 - Spectral GCNs
 - Applications

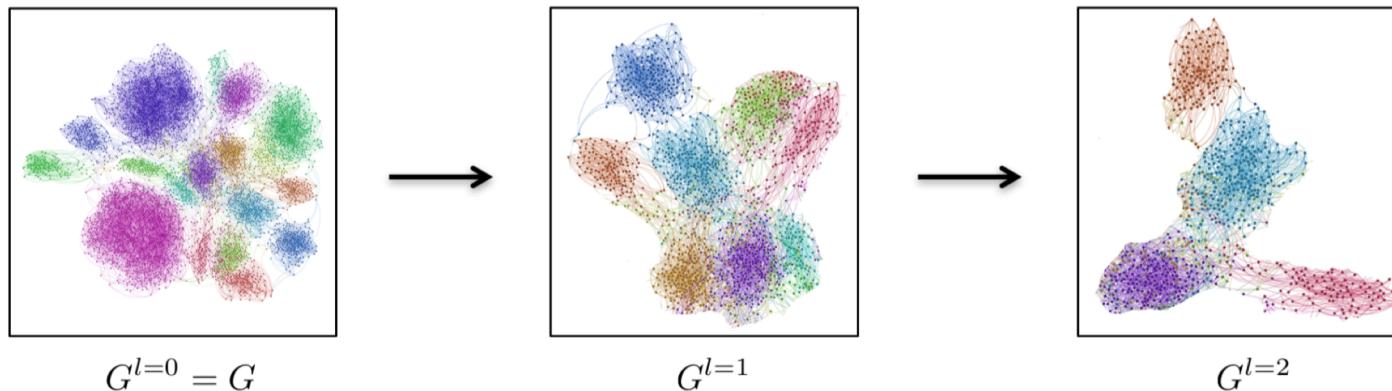


Graph Down-Sampling

- Goals :
 - Pool similar local features with max/average pooling.
 - Create **invariance** to **local** geometric deformations.
 - Series of pooling layers create **invariance** to **global** geometric deformations.
- Challenges :
 - Design a **multi-scale coarsening** algorithm that preserves **non-linear** graph structures.
 - How to make graph pooling fast?

Graph Down-Sampling

- Graph **down-sampling** \Leftrightarrow graph **coarsening** \Leftrightarrow graph **partitioning** :
 - Decompose G into meaningful clusters.



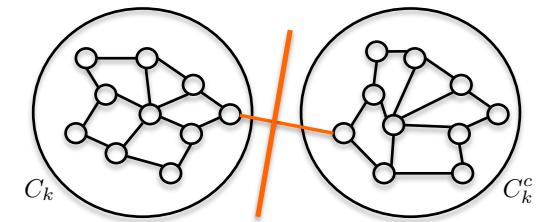
- Graph partitioning is NP-hard \Rightarrow Approximation

Balanced Cuts

- Powerful combinatorial graph partitioning models :

- Normalized cut (Shi, Malik, 2000)

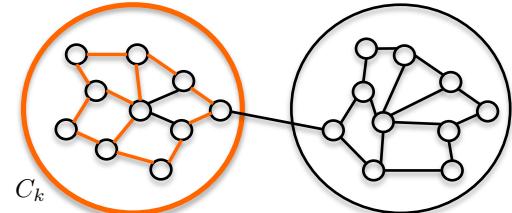
$$\min_{C_1, \dots, C_K} \sum_{k=1}^K \frac{\text{Cut}(C_k, C_k^c)}{\text{Vol}(C_k)}$$



Partitioning by min edge cuts.

- Normalized association

$$\max_{C_1, \dots, C_K} \sum_{k=1}^K \frac{\text{Assoc}(C_k)}{\text{Vol}(C_k)}$$



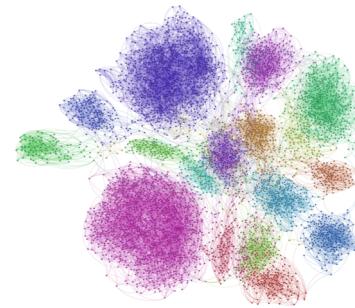
Partitioning by max vertex matching.

where $\text{Cut}(A, B) := \sum_{i \in A, j \in B} a_{ij}$, $\text{Assoc}(A) := \sum_{i \in A, i \in B} a_{ij}$,
 $\text{Vol}(A) := \sum_{i \in A, j \in B} d_i$, and $d_i := \sum_{j \in V} a_{ij}$.

- Both models are equivalent, but lead to different algorithms.

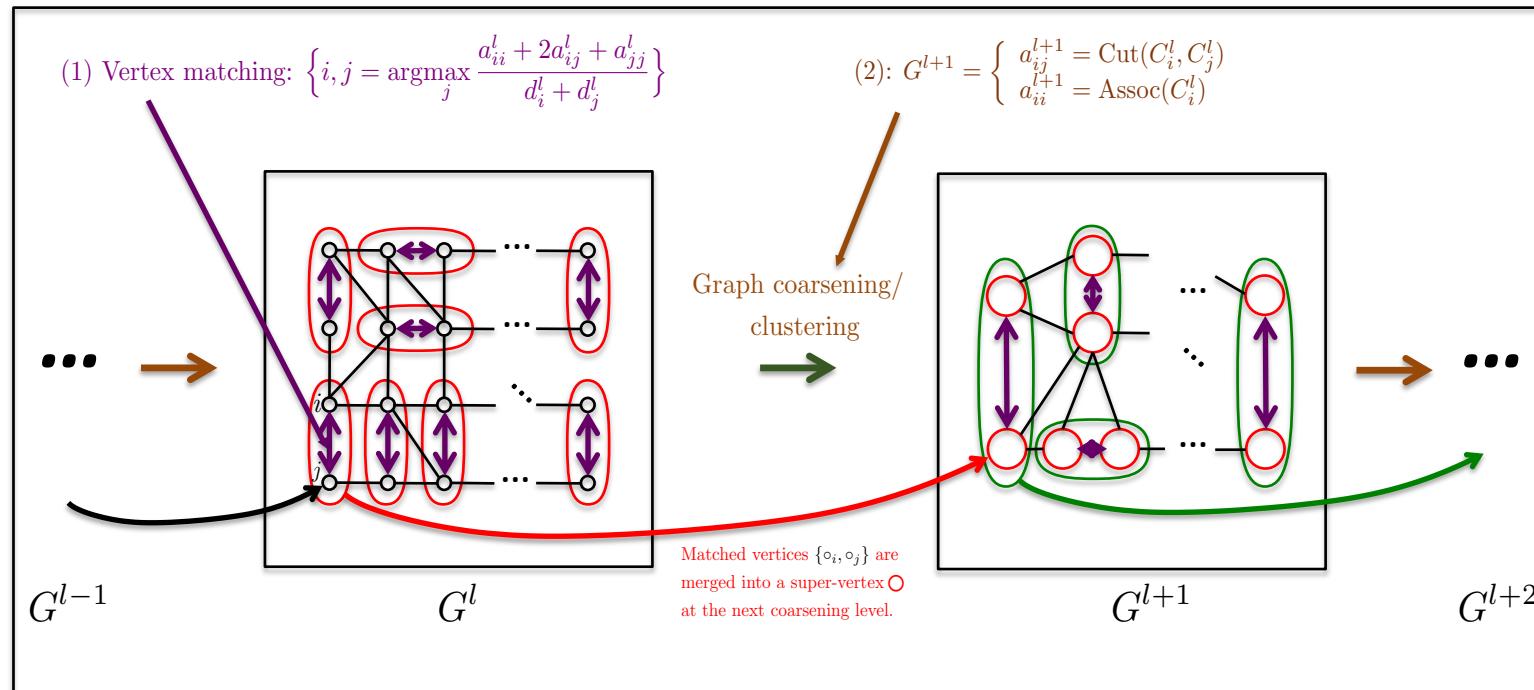
Balanced Cuts

- Balanced cuts are **NP-hard** :
 - Most popular approximation techniques focus on **linear spectral relaxation** (eigenproblem with global solution).
- Graph geometry are generally **not linear** :
 - **Graclus** (Dhillon, Guan, Kulis 2007) algorithm computes non-linear clusters that locally maximize the Normalized Association.
- Graclus algorithm offers a control of the **coarsening ratio of ≈ 2** (like image grid) using **heavy-edge matching (HEM)** by Karypis, Kumar 1995.



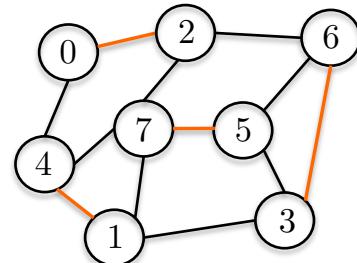
Heavy-Edge Matching (HEM)

- HEM proceeds by **two** successive steps (that guarantees a local solution of Norm Assoc):
 - **Vertex matching**
 - **Graph coarsening**

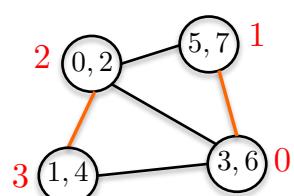


Unstructured Pooling

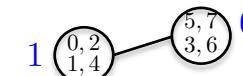
- Sequence of coarsened graphs produced by HEM :



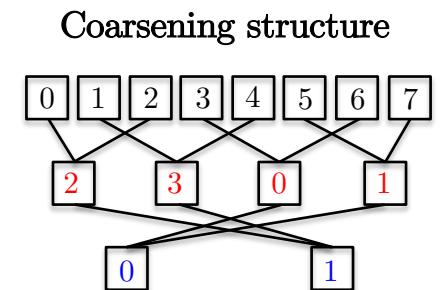
$$G^{l=0} = G$$



$$G^{l=1}$$



$$G^{l=2}$$

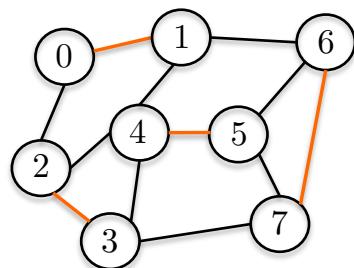


Unstructured arrangement of nodes

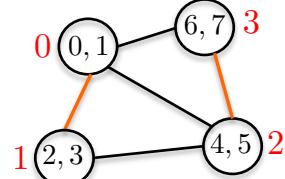
- Stores a table of indices for graph and all its coarsened versions.
- Computationally inefficient

Fast Graph Pooling^[1]

- Structured pooling :
 - Arrangement of the node indexing such that adjacent nodes are hierarchically merged at the next coarser level.



$$G^{l=0} = G$$

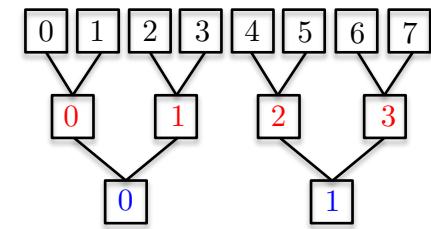


$$G^{l=1}$$



$$G^{l=2}$$

Coarsening structure



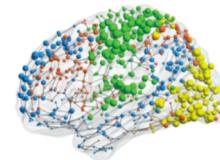
Binary tree arrangement of nodes

☺ As efficient as 1D-Euclidean grid pooling !

[1] Defferrard, Bresson, Vandergheynst, Convolutional neural networks on graphs with fast localized spectral filtering, 2016

Outline

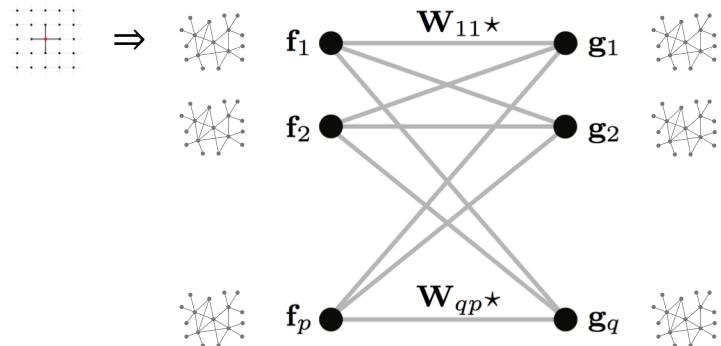
- Part 2: Spectral Graph ConvNets
 - Problem Setting
 - Graph Theory
 - Spectral Graph Convolution
 - Graph Pooling
 - Spectral GCNs
 - Applications



Vanilla Spectral GCNs^[1]

- Graph convolutional layer :

$$\begin{aligned}\mathbf{f}_l &= l\text{-th data feature on graphs, } \dim(\mathbf{f}_l) = n \times 1 \\ \mathbf{g}_l &= l\text{-th feature map, } \dim(\mathbf{g}_l) = n \times 1\end{aligned}$$



Conv. layer $\mathbf{g}_l = \xi \left(\sum_{l'=1}^p \mathbf{W}_{l,l'} \star \mathbf{f}_{l'} \right)$

Activation, e.g. $\xi(x) = \max\{x, 0\}$ rectified linear unit (ReLU)

[1] Bruna, Zaremba, Szlam, LeCun 2014

Spectral Graph Convolution

- Convolutional layer in the spatial domain:

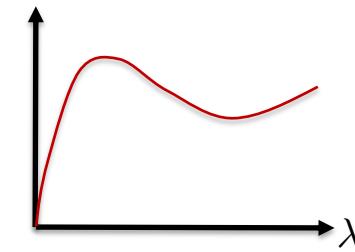
$$\mathbf{g}_l = \xi \left(\sum_{l'=1}^p \mathbf{W}_{l,l'} \star \mathbf{f}_{l'} \right),$$



where $\mathbf{W}_{l,l'}$ = matrix of graph spatial filter,

can also be expressed in the spectral domain (using $\mathbf{g} \star \mathbf{f} = \Phi \hat{\mathbf{g}}(\Lambda) \Phi^\top \mathbf{f}$)

$$\mathbf{g}_l = \xi \left(\sum_{l'=1}^p \Phi \hat{\mathbf{W}}_{l,l'} \Phi^\top \mathbf{f}_{l'} \right),$$



where $\hat{\mathbf{W}}_{l,l'} = n \times n$ diagonal matrix of graph spectral filter.

We will denote the spectral filter without the hat symbol, i.e. $\mathbf{W}_{l,l'}$

Vanilla Spectral GCNs^[1]

- Series of spectral convolutional layers :

$$\mathbf{g}_l^{(k)} = \xi \left(\sum_{l'=1}^{q^{(k-1)}} \Phi \mathbf{W}_{l,l'}^{(k)} \Phi^\top \mathbf{g}_{l'}^{(k-1)} \right),$$

with spectral coefficients $\mathbf{W}_{l,l'}^{(k)}$ to be learned at each layer.

- ☺ First spectral graph CNN architecture
- ☹ No guarantee of spatial localization of filters
- ☹ $O(n)$ parameters per layer
- ☹ $O(n^2)$ computation of forward and inverse Fourier transforms ϕ, ϕ^\top (no FFT on graphs)
- ☹ Filters are basis-dependent \Rightarrow does not generalize across graphs

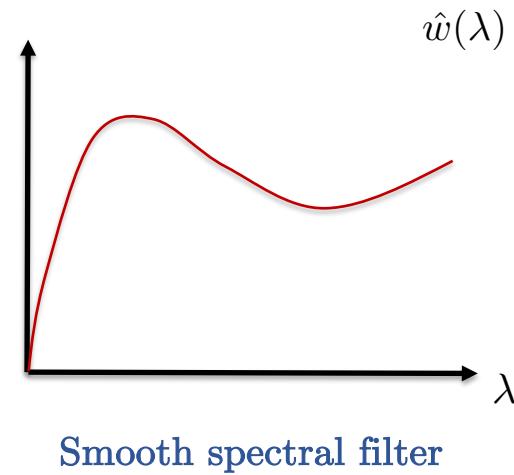
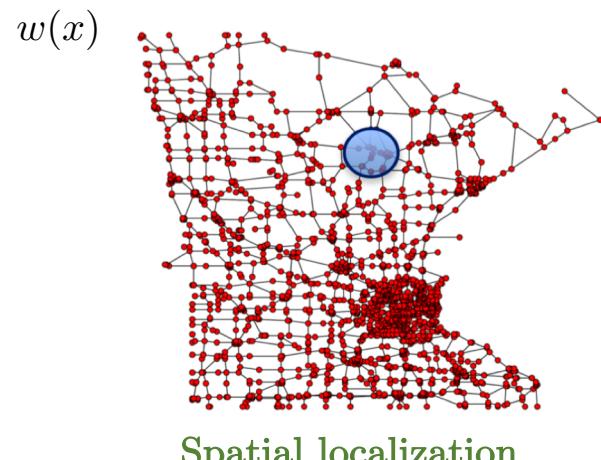
[1] Bruna, Zaremba, Szlam, LeCun 2014

Spatial Localization^[1]

- Parseval's identity

$$\int_{-\infty}^{+\infty} |x|^{2k} |w(x)|^2 dx = \int_{-\infty}^{+\infty} \left| \frac{\partial^k \hat{w}(\lambda)}{\partial \lambda^k} \right|^2 d\lambda$$

⇒ Localization in space = smoothness in frequency domain



[1] Henaff, Bruna, LeCun 2015

Smooth Spectral Filter

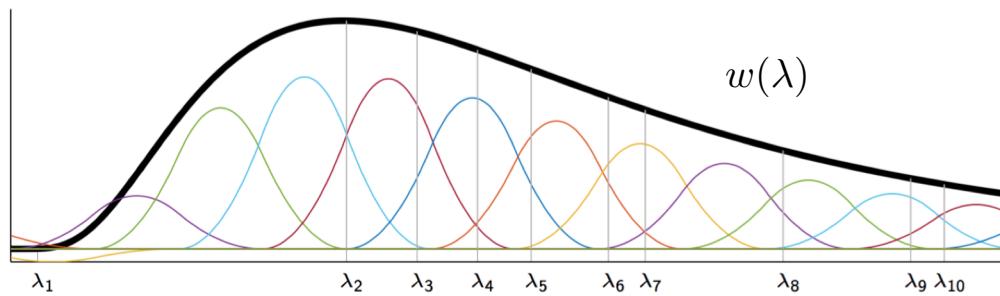
- Parametrize a smooth spectral filter function $w(\lambda)$ with a linear combination of smooth kernel functions $\beta_1(\lambda), \dots, \beta_r(\lambda)$, here splines^[1]:

$$w_{\alpha}(\lambda) = \sum_{j=1}^r \alpha_j \beta_j(\lambda)$$

\Downarrow

$$w_{\alpha}(\lambda_i) = \sum_{j=1}^r \alpha_j \beta_j(\lambda_i) = (\mathbf{B}\alpha)_i \quad \Rightarrow \quad \mathbf{W} = \text{Diag}(\mathbf{B}\alpha)$$

Splines
Coefficients (learned by backpropagation)



[1] Henaff, Bruna, LeCun 2015

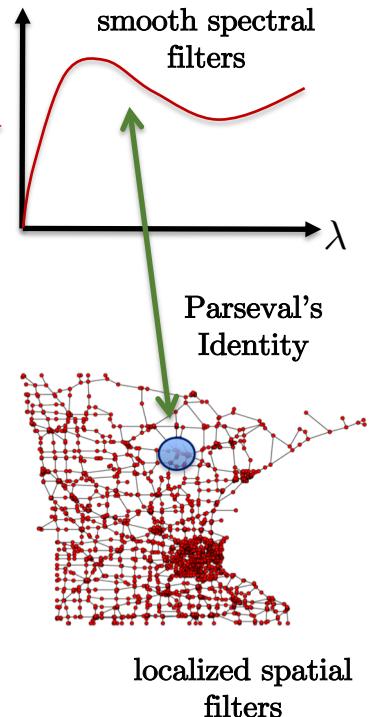
SplineGCNs^[1,2]

- Series of graph convolutional layers:

$$\mathbf{g}_l^{(k)} = \xi \left(\sum_{l'=1}^{q^{(k-1)}} \Phi \mathbf{W}_{l,l'}^{(k)} \Phi^\top \mathbf{g}_{l'}^{(k-1)} \right),$$

with **smooth spectral filters** at each layer.

- ☺ Fast-decaying filters in space
- ☺ $O(1)$ parameters per layer
- ☺ $O(n^2)$ computation of forward and inverse Fourier transforms
 ϕ, ϕ^\top (no FFT on graphs)
- ☺ Filters are **basis-dependent** \Rightarrow does not generalize across graphs



[1] Bruna, Zaremba, Szlam, LeCun, Spectral Networks and Locally Connected Networks on Graphs, 2014
[2] Henaff, Bruna, LeCun, Deep Convolutional Networks on Graph-Structured Data, 2015

Polynomial Spectral Filters^[1]

- Represent spectral smooth functions with **polynomials of Laplacians** :

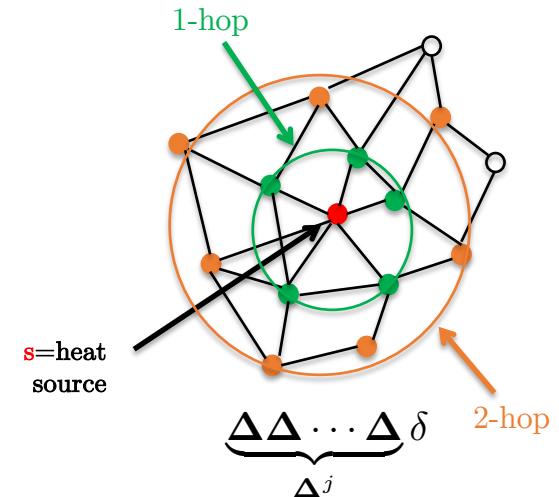
$$w_{\alpha}(\lambda) = \sum_{j=0}^r \alpha_j \lambda^j$$

where $\alpha = (\alpha_1, \dots, \alpha_r)^\top$ is the vector of filter parameters.

- **Convolutional layer** : Apply spectral filter to feature signal f

$$w_{\alpha}(\Delta)f = \sum_{j=0}^r \alpha_j \Delta^j f$$

- Key observation: Each Laplacian operation increases the support of a function by 1-hop \Rightarrow Exact control the size of Laplacian-based filters.



Linear Complexity

- Application of the filter to a feature signal f :

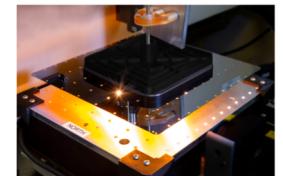
$$w_{\alpha}(\Delta)f = \sum_{j=0}^r \alpha_j \Delta^j f$$

- Denote $\mathbf{X}_0 = f$ and define $\mathbf{X}_1 = \Delta \mathbf{X}_0 = \Delta f$ and the sequence $\mathbf{X}_j = \Delta \mathbf{X}_{j-1}$

$$w_{\alpha}(\Delta)f = \sum_{j=0}^r \alpha_j \mathbf{X}_j$$

- Two important observations:

1. *No* computation of the eigen-decomposition of Laplacian (ϕ, Λ) is required.
 2. Observe that the sequence $\{\mathbf{X}_j\}$ is generated by **multiplication of a sparse matrix Δ and a vector** \Rightarrow Complexity is $O(E.r) = O(n)$ for sparse (real-world) graphs.
- Graph convolutional layers are **GPU** friendly (but not yet optimized or soon to be).



Aug. 2019

PolyGCNs^[1]

- Series of spectral convolutional layers with **monomials of Laplacian** $1, x, x^2, x^3, \dots$:

$$w_{\alpha}(\Delta)\mathbf{f} = \sum_{j=0}^r \alpha_j \Delta^j \mathbf{f}$$

- Properties :

- ☺ Filters are **exactly localized** in r -hops support
- ☺ **O(1)** parameters per layer
- ☺ No eigen-decomposition/ $\phi, \phi^\top \Rightarrow O(n)$ computational complexity (assuming sparsely-connected graphs)
- ☹ Unstable under coefficients perturbation (hard to optimize)
- ☹ Filters are **basis-dependent** \Rightarrow does not generalize across graphs

Chebyshev Polynomials

- Graph convolution with **non-orthogonal** monomial basis $1, x, x^2, x^3, \dots$

$$w_{\alpha}(\Delta)\mathbf{f} = \sum_{j=0}^r \alpha_j \Delta^j \mathbf{f}$$

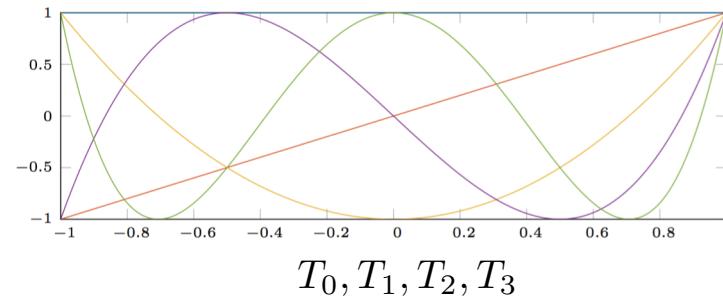
- Graph convolution with **orthogonal Chebyshev polynomials** :

$$w_{\alpha}(\tilde{\Delta})\mathbf{f} = \sum_{j=0}^r \alpha_j T_j(\tilde{\Delta})\mathbf{f}$$

$\tilde{\Delta} = 2\lambda_n^{-1}\Delta - \mathbf{I}$
Scaled Laplacian

- **Orthonormal** on $L^2([-1, +1])$ w.r.t. $\langle f, g \rangle = \int_{-1}^{+1} f(\tilde{\lambda})g(\tilde{\lambda}) \frac{d\tilde{\lambda}}{\sqrt{1-\tilde{\lambda}^2}}$

- **Stable under perturbation** of coefficients



ChebGCNs^[1]

- Application of the Chebyshev filter to a feature signal \mathbf{f} :

with

$$w_{\alpha}(\tilde{\Delta})\mathbf{f} = \sum_{j=0}^r \alpha_j T_j(\tilde{\Delta})\mathbf{f} = \sum_{j=0}^r \alpha_j \mathbf{X}^{(j)}$$

$$\begin{aligned}\mathbf{X}^{(j)} &= T_j(\tilde{\Delta})\mathbf{f} \\ &= 2\tilde{\Delta}\mathbf{X}^{(j-1)} - \mathbf{X}^{(j-2)}, \quad \mathbf{X}^{(0)} = \mathbf{f}, \quad \mathbf{X}^{(1)} = \tilde{\Delta}\mathbf{f}\end{aligned}$$

- ☺ Filters are exactly localized in r -hops support
- ☺ $O(1)$ parameters per layer
- ☺ No eigen-decomposition/ $\phi, \phi^\top \Rightarrow O(n)$ computational complexity (assuming sparsely-connected graphs)
- ☺ Stable under coefficients perturbation
- ☹ Filters are basis-dependent \Rightarrow does not generalize across graphs

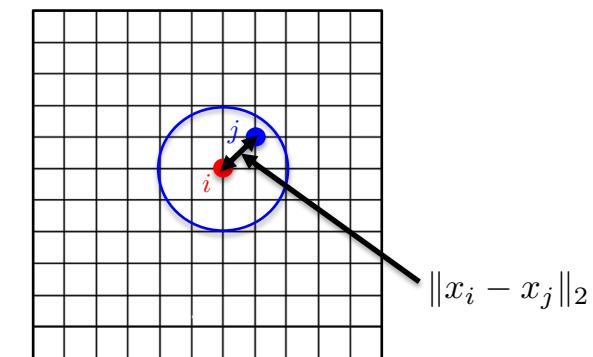
Numerical Experiments

- ChebGCNs for **Image Classification**.
- MNIST : 70,000 images represented on a **2D grid of size 28x28** (dim=28²=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

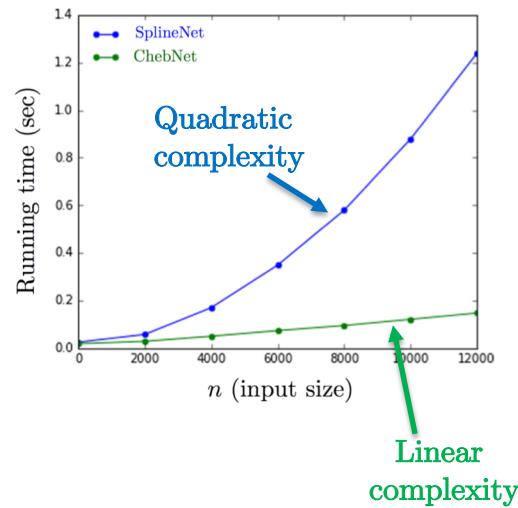
- **Graph :**
 - A k-NN graph ($k=8$) of the Euclidean grid :

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

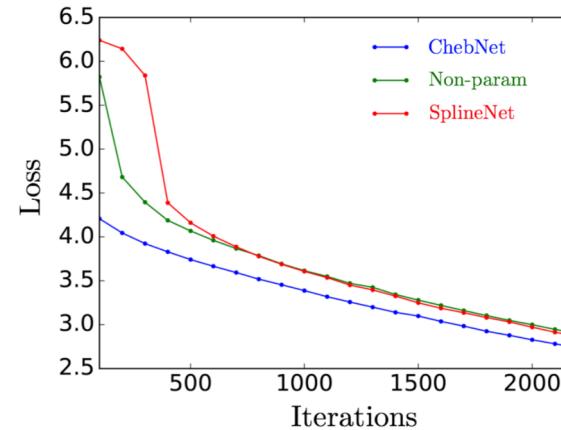


Numerical Experiments

- Running time

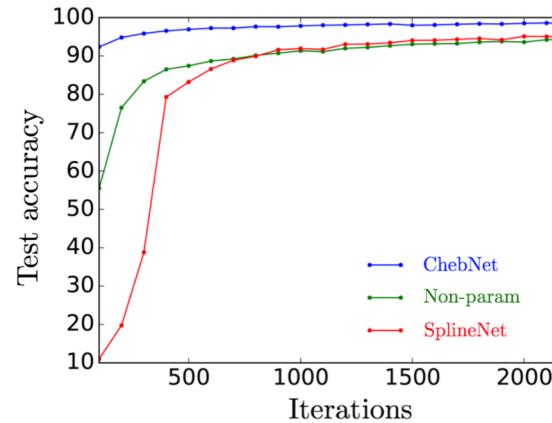


- Optimization



- Accuracy

Model	Order	Accuracy
LeNet5	-	99.33%
SplineNet	25	97.75%
ChebNet	25	99.14%



Lab 01

- Traditional LeNet5 ConvNets
- LeNet5 ChebGCNs

The image shows two Jupyter notebook interfaces side-by-side. Both notebooks have a title bar indicating they are "Trusted" and use "Python 3".

Left Notebook (01_standard_ConvNets):

- Title:** Lab 01 : Traditional LeNet5 ConvNets - demo
- Abstract:** LeNet5 Convolutional Neural Networks
Gradient-based learning applied to document recognition
Y LeCun, L Bottou, Y Bengio, P Haffner
Proceedings of the IEEE 86 (11), 2278-2324
- Code:** A single cell (In [1]) containing Python code to import torch, set CUDA device order, and initialize the manual seed.

```
import torch
from torch.autograd import Variable
import torch.nn.functional as F
import torch.nn as nn
import collections
import time
import numpy as np

import os
os.environ["CUDA_DEVICE_ORDER"]="PCI_BUS_ID"
os.environ["CUDA_VISIBLE_DEVICES"]="0"

if torch.cuda.is_available():
    print('cuda available')
    dtypeFloat = torch.cuda.FloatTensor
    dtypeLong = torch.cuda.LongTensor
    torch.cuda.manual_seed(1)
else:
    print('cuda not available')
    dtypeFloat = torch.FloatTensor
    dtypeLong = torch.LongTensor
    torch.manual_seed(1)
```

Right Notebook (02_ChebGCNs):

- Title:** Lab 01 : LeNet5 ChebGCNs - demo
- Abstract:** Spectral Graph ConvNets
Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering
M Defferrard, X Bresson, P Vandergheynst
Advances in Neural Information Processing Systems, 3844-3852, 2016
ArXiv preprint: arXiv:1606.09375
- Code:** A single cell (In [1]) containing Python code to import torch, set CUDA device order, and initialize the manual seed.

```
import torch
from torch.autograd import Variable
import torch.nn.functional as F
import torch.nn as nn
import collections
import time
import numpy as np

import sys
sys.path.insert(0, 'lib')
load_ext autoreload
autoreload 2

import os
os.environ["CUDA_DEVICE_ORDER"]="PCI_BUS_ID"
os.environ["CUDA_VISIBLE_DEVICES"]="0"

if torch.cuda.is_available():
    print('cuda available')
    dtypeFloat = torch.cuda.FloatTensor
    dtypeLong = torch.cuda.LongTensor
    torch.cuda.manual_seed(1)
else:
    print('cuda not available')
    dtypeFloat = torch.FloatTensor
    dtypeLong = torch.LongTensor
    torch.manual_seed(1)
```

Lab 01

- ChebGCNs for **your graph** problems :

A is the **adjacency matrix** of your graph.

Graph Adjacency Matrix

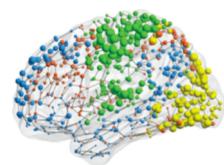
```
: from lib.grid_graph import grid_graph
from lib.coarsening import coarsen
from lib.coarsening import lmax_L
from lib.coarsening import perm_data
from lib.coarsening import rescale_L

# Construct graph
t_start = time.time()
grid_side = 28
number_edges = 8
metric = 'euclidean'

##### YOUR GRAPH ADJACENCY MATRIX HERE #####
A = grid_graph(grid_side,number_edges,metric) # create graph of Euclidean grid
##### YOUR GRAPH ADJACENCY MATRIX HERE #####
# Compute coarsened graphs
coarsening_levels = 4
L, perm = coarsen(A, coarsening_levels)
```

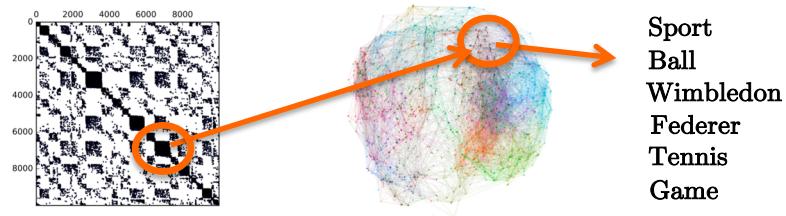
Outline

- Part 2: Spectral Graph ConvNets
 - Problem Setting
 - Graph Theory
 - Spectral Graph Convolution
 - Graph Pooling
 - Spectral GCNs
 - Applications

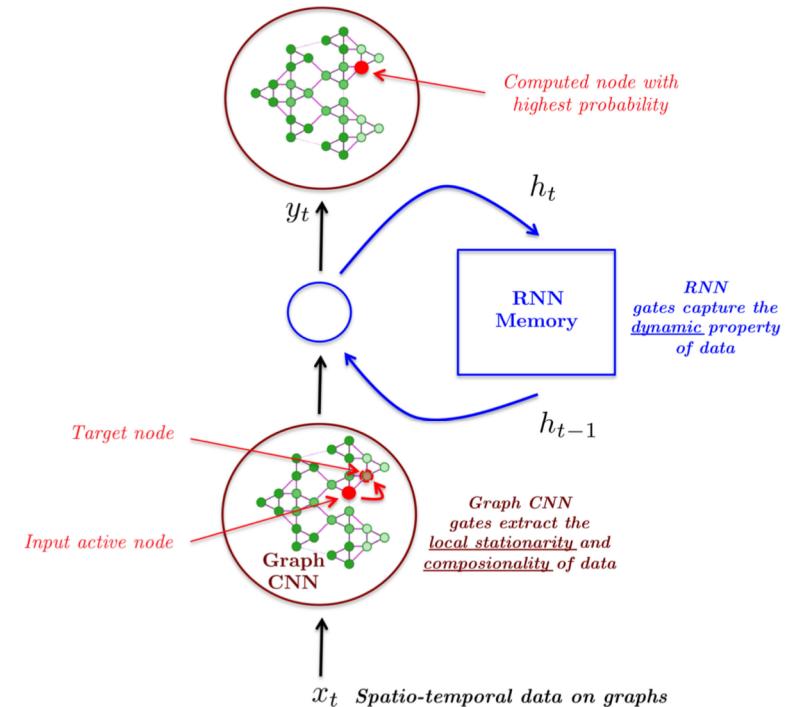


Temporal Predictions on Graphs^[1]

- We combine ChebGCNs + LSTM :
 - ChebGCNs analyzes **static information** from data on graphs (find multi-scale local stationary patterns).
 - Words form a **space of topics** that are encoded with graphs of relationships.



- LSTM identifies **dynamical information** (non-linear temporal properties of data).
- We consider the basic NLP task of **Word Prediction**.



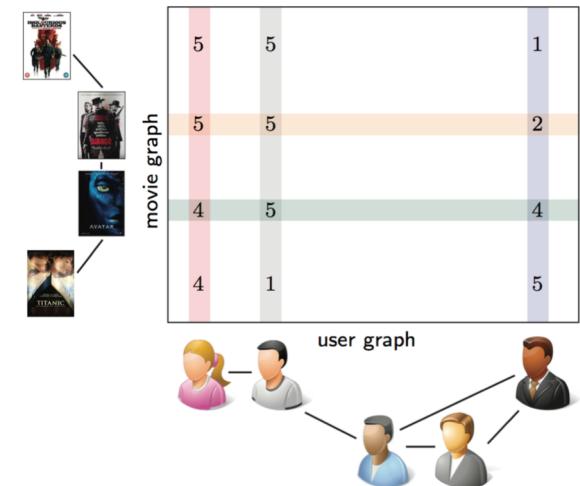
[1] Y Seo, M Defferrard, P Vandergheynst, X Bresson, Structured sequence modeling with graph convolutional recurrent networks, International Conference on Neural Information Processing, 2018

Recommender Systems

- Tackling the Netflix challenge.

- Matrix Completion on Multiple Graphs^[1] :

$$\min_{\mathbf{X} \in \mathbb{R}^{m \times n}} \|\mathbf{X}\|_* + \mu \|\Omega \circ (\mathbf{X} - \mathbf{A})\|_F^2 \\ + \mu_c \underbrace{\text{tr}(\mathbf{X} \Delta_c \mathbf{X}^\top)}_{\|\mathbf{X}\|_{G_c}^2} + \mu_r \underbrace{\text{tr}(\mathbf{X}^\top \Delta_r \mathbf{X})}_{\|\mathbf{X}\|_{G_r}^2}$$



- Multi-graph convolutional layer^[2] :

$$\mathbf{Y}_l^{(k)} = \xi \left(\sum_{l'=1}^p \sum_{j,j'=0}^r \theta_{jj' ll'} T_j(\tilde{\Delta}_r) \mathbf{Y}_{l'}^{(k-1)} T_{j'}(\tilde{\Delta}_c) \right)$$

Chebyshev
 2D-spectral filters

$l = 1, \dots, q$ #output features
 $l' = 1, \dots, p$ #input features

[1] V Kalofolias, X Bresson, M Bronstein, P Vandergheynst, Matrix completion on graphs, NeurIPS'14

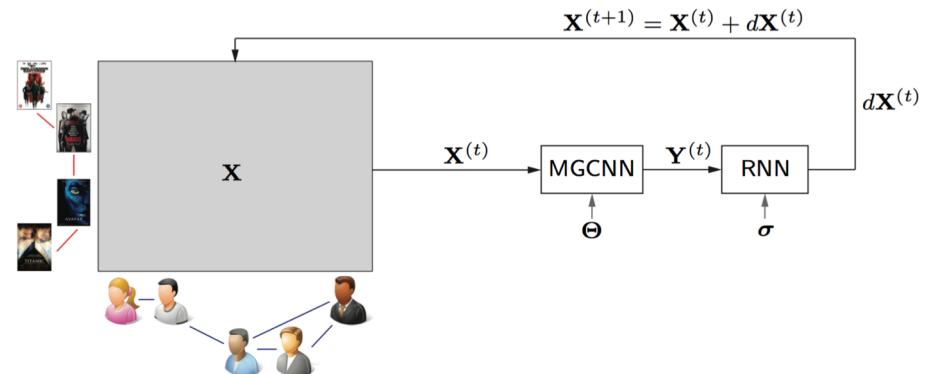
[2] F Monti, M Bronstein, X Bresson, Geometric matrix completion with recurrent multi-graph neural networks, NeurIPS'17

Recommender Systems^[1]

- We cast the completion problem as **non-linear diffusion** process.
 - Less parameters** to learn (less overfitting).
 - Favorable for **highly sparse entries** (better diffusion).
 - Allows to **learn temporal dynamics** of entries if such information is available.

Method	MovieLens ¹	Flixster ²	Douban ³	Yahoo ⁴
IMC ⁵	1.653	–	–	–
GMC ⁶	0.996	–	–	–
MC ⁷	0.973	–	–	–
GRALS ⁸	0.945	1.245	0.833	38.042
sRGCNN (Cheb)⁹	0.929	0.926	0.801	22.415
sRGCNN (Cayley)¹⁰	0.922	–	–	–

Data: ¹Miller et al. 2003; ²Jamali, Ester 2010; ³Ma et al. 2011; ⁴Dror et al. 2012
 Methods: ⁵Jain, Dhillon 2013; ⁶Kalofolias et al. 2014; ⁷Candès, Recht 2012; ⁸Rao et al. 2015; ⁹Monti, Bresson, Bronstein 2017; ¹⁰Levie et al. 2017



Recurrent multigraph CNN (RMCNN) architecture
for matrix completion

$$\min_{\Theta, \sigma} \|\mathbf{X}_{\Theta, \sigma}^{(T)}\|_{\mathcal{G}_r}^2 + \|\mathbf{X}_{\Theta, \sigma}^{(T)}\|_{\mathcal{G}_c}^2 + \frac{\mu}{2} \|\Omega \circ (\mathbf{X}_{\Theta, \sigma}^{(T)} - \mathbf{Y})\|_F^2$$

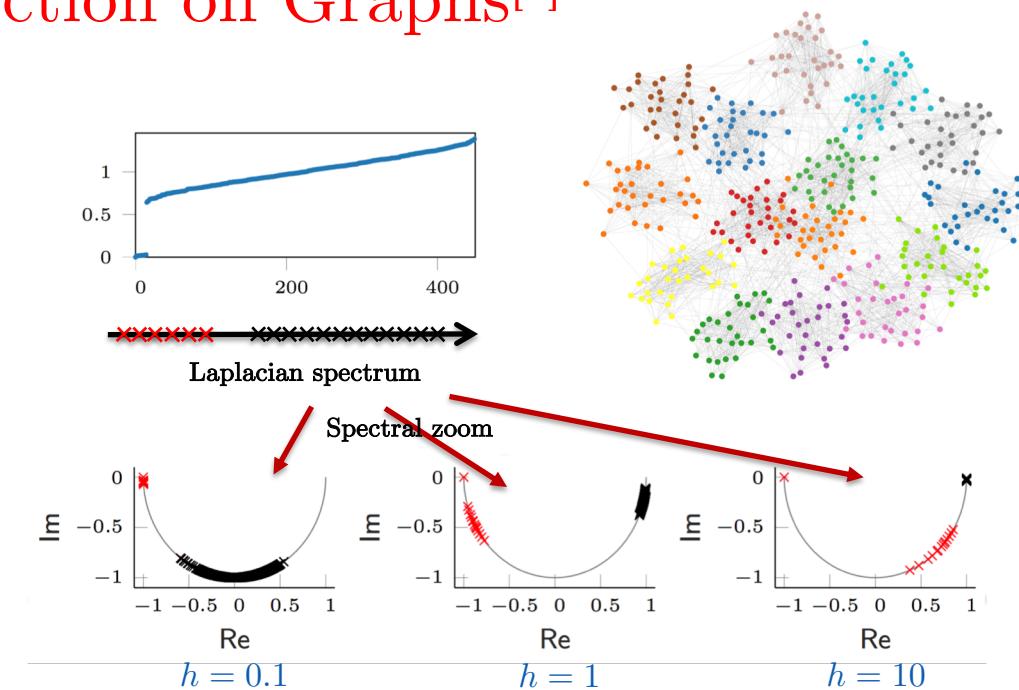
[1] F Monti, M Bronstein, X Bresson, Geometric matrix completion with recurrent multi-graph neural networks, NeurIPS'17

Community Detection on Graphs^[1]

- ChebGCNs are **unstable** to produce filters with frequency bands of interest.
- We represent spectral filters with **Cayley rationals** of order r :

$$w_{\mathbf{c},h}(\lambda) = c_0 + 2\operatorname{Re} \left\{ \sum_{j=1}^r c_j (h\lambda - i)^j (h\lambda + i)^{-j} \right\}$$

Coefficients
(bypropagation)
Spectral zoom
(hyper-parameter)

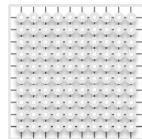


- ☺ CayleyGCNs w/ **same properties** of ChebGCNs in addition of
- ☺ **Spectral zoom** property allowing to better **localize** in frequency.
- ☺ **Richer class of filters** than Chebyshev for the same order.

[1] R Levie, F Monti, X Bresson, MM Bronstein, Cayleynets: Graph convolutional neural networks with complex rational spectral filters, IEEE Transactions on Signal Processing, 2018

Outline

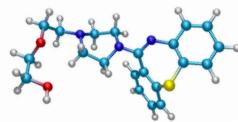
- Part 1: Traditional ConvNets



- Part 2: Spectral Graph ConvNets



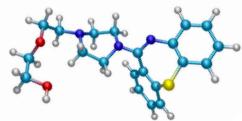
- Part 3: Spatial Graph ConvNets



- Conclusion

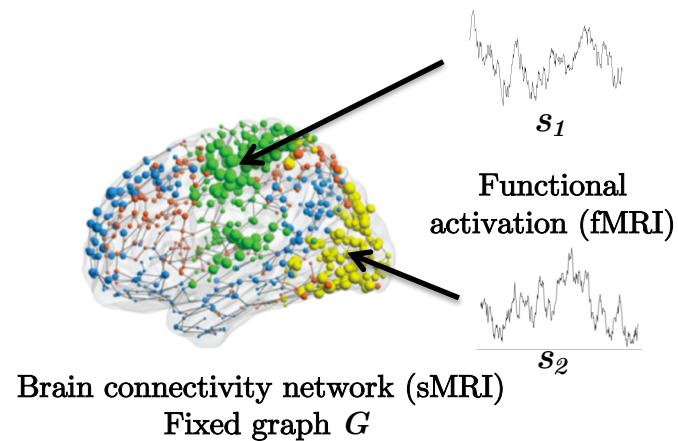
Outline

- Part 3: Spatial Graph ConvNets
 - Problem Setting
 - Spatial Graph Architectures
 - Quantum Chemistry
 - Operations Research
 - DGL

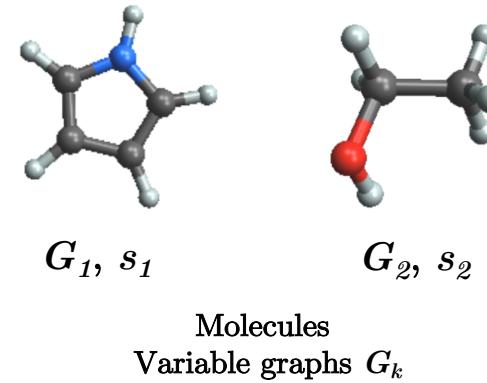


Problem Setting

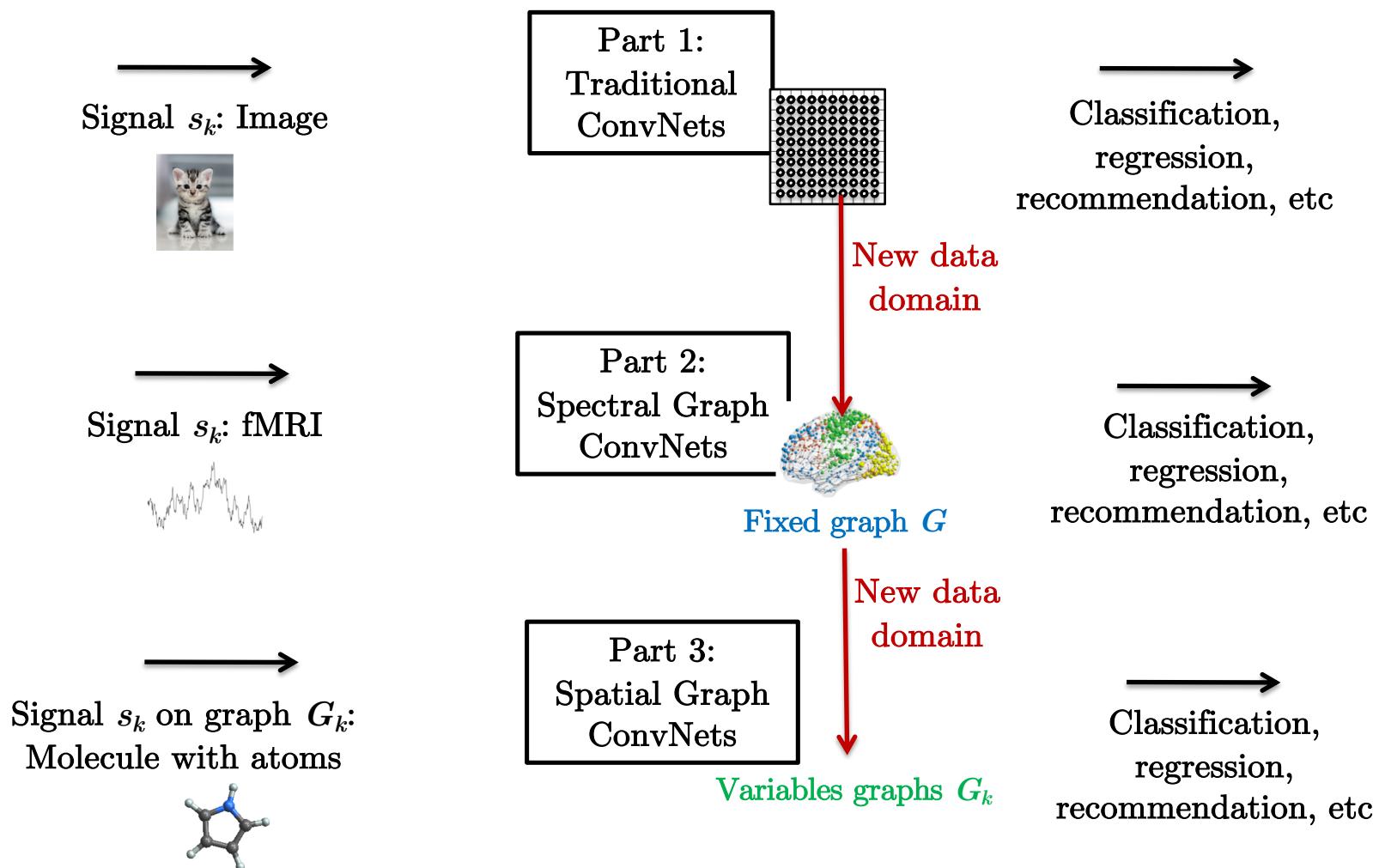
- Spectral Graph ConvNets :
 - Given **fixed** graph(s) G , and a set of signals s_k on G to be analyzed.



- Spatial Graph ConvNets :
 - Given a set of **variable** graphs G_k and signals s_k on G_k to be analyzed.

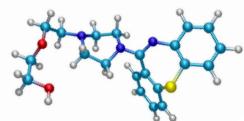


Graph ConvNet Architectures



Outline

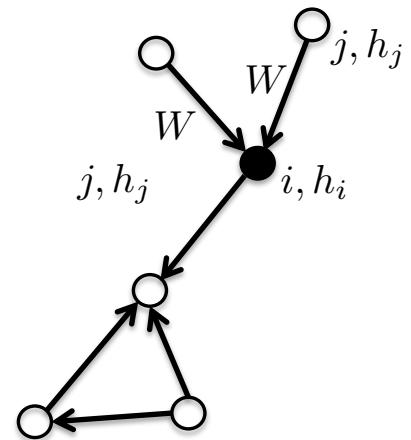
- Part 3: Spatial Graph ConvNets
 - Problem Setting
 - Spatial Graph Architectures
 - Quantum Chemistry
 - Operations Research
 - DGL



Spatial GNNs^[1]

- Minimal inner structures to design spatial GNNs :
 - Invariant by vertex re-indexing (no graph matching is required)
 - Locality/local reception field (only neighbors are considered)
 - Weight sharing (convolutional operations)
 - Independence w.r.t. graph size

$$h_i^{\ell+1} = f_{\text{GNN}}(h_i^\ell, \{h_j^\ell : j \rightarrow i\})$$



[1] Scarselli, Gori, Tsoi, Hagenbuchner, Monfardini, The Graph Neural Network Model, 2009
[2] Li, Tarlow, Brockschmidt, Zemel, Gated Graph Sequence Neural Networks, 2015

Graph RNNs

- Vanilla Graph RNNs with Multi-Layer Perceptron (MLP)^[1] :

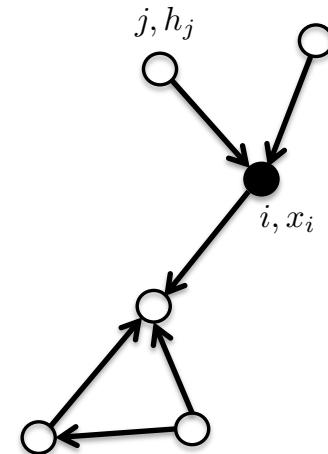
$$h_i = \sum_{j \rightarrow i} \mathcal{C}_{\text{G-MLP}}(x_i, h_j) = \sum_{j \rightarrow i} A \sigma(B \sigma(Ux_i + Vh_j))$$

- Graph GRUs^[2,3] (Gated Recurrent Units) :

$$h_i = \mathcal{C}_{\text{G-GRU}}(x_i, \sum_{j \rightarrow i} h_j)$$

Fixed-point iterative scheme required
(unlike NLP, no directed acyclic graph) :

$$\begin{aligned}\bar{h}_i^t &= \sum_{j \rightarrow i} h_j^t, \quad h_i^{t=0} = x_i \\ z_i^{t+1} &= \sigma(U_z h_i^t + V_z \bar{h}_i^t) \\ r_i^{t+1} &= \sigma(U_r h_i^t + V_r \bar{h}_i^t) \\ \tilde{h}_i^{t+1} &= \tanh(U_h(h_i^t \odot r_i^{t+1}) + V_h \bar{h}_i^t) \\ h_i^{t+1} &= (1 - z_i^{t+1}) \odot h_i^t + z_i^{t+1} \odot \tilde{h}_i^{t+1}\end{aligned}$$



[1] Scarselli, Gori, Tsoi, Hagenbuchner, Monfardini, The Graph Neural Network Model, 2009

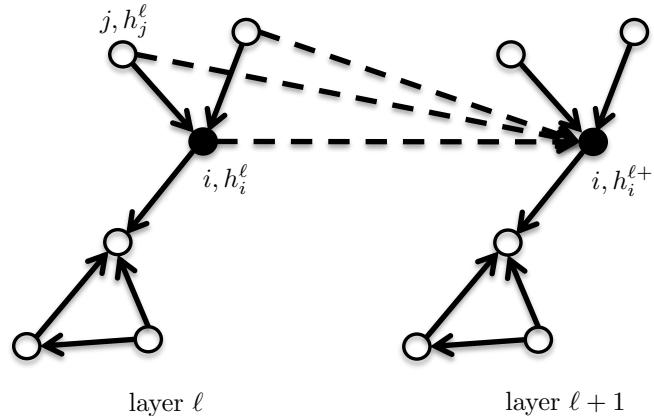
[2] Li, Tarlow, Brockschmidt, Zemel, Gated Graph Sequence Neural Networks, 2015

[3] Cho, Merriennboer, Gulcehre, Bahdanau, Bougares, Schwenk, Bengio, Learning Phrase Representations using RNN for Statistical Machine Translation, 2014

Vanilla Graph ConvNets

- Vanilla GCNs^[1,2] (with ReLU) and GraphSAGE^[3] (with max) :

$$\begin{aligned} h_i^{\ell+1} &= \mathcal{C}_{\text{G-VCN}}\left(h_i^\ell, \sum_{j \rightarrow i} h_j^\ell\right), \quad h_i^{\ell=0} = x_i \\ &= \text{ReLU}\left(U^\ell h_i^\ell + V^\ell \sum_{j \rightarrow i} h_j^\ell\right), \quad h_i^{\ell=0} = x_i \end{aligned}$$



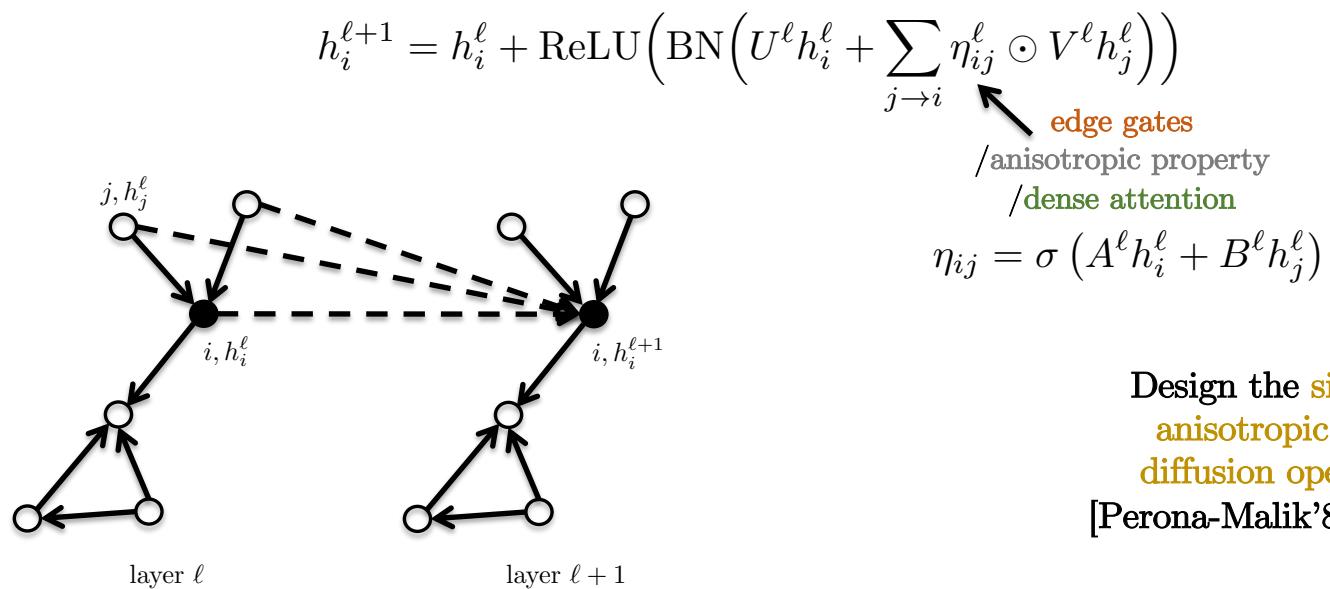
[1] Sukhbaatar, Szlam, Fergus, Learning Multiagent Communication with Backpropagation, 2016

[2] Kipf, Welling, Semi-Supervised Classification with Graph Convolutional Networks, 2017

[3] Hamilton, Ying, Leskovec, Inductive representation learning on large graphs, 2017

GatedGCNs^[1]

- Graph ConvNets architecture^[2,3,4] augmented with edge gating mechanism^[5], residuality^[6] and batch normalization^[7]:



[1] Bresson, Laurent, Residual gated graph convnets, 2017

[2] Sukhbaatar, Szlam, Fergus, Learning Multiagent Communication with Backpropagation, 2016

[3] Kipf, Welling, Semi-Supervised Classification with Graph Convolutional Networks, 2017

[4] Hamilton, Ying, Leskovec, Inductive representation learning on large graphs, 2017

[5] Marcheggiani, Titov, Encoding Sentences with Graph Convolutional Networks for Semantic Role Labeling, 2017

[6] He, Zhang, Ren, Sun, Deep Residual Learning for Image Recognition, 2016

[7] Ioffe, Szegedy, Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift, 2015

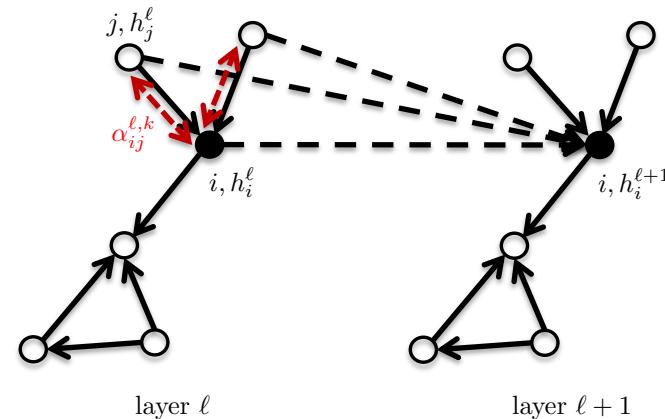
Graph Attention Networks^[1]

- **Attention mechanism** in 1-hop neighborhood :

$$h_i^{\ell+1} = \left(\left\|_{k=1}^K \sigma \left(\sum_{j \rightarrow i} \alpha_{ij}^{\ell,k} V^{\ell,k} h_j^\ell \right) \right) \right) W^\ell$$

$$\begin{aligned} \alpha_{ij}^{\ell,k} &= \text{softmax}_{1\text{-hop}} \left((Q^{\ell,k} h_i^\ell)^T (K^{\ell,k} h_j^\ell) \right) \\ &= \frac{e^{(Q^{\ell,k} h_i^\ell)^T (K^{\ell,k} h_j^\ell)}}{\sum_{j' \rightarrow i} e^{(Q^{\ell,k} h_i^\ell)^T (K^{\ell,k} h_{j'}^\ell)}} \end{aligned}$$

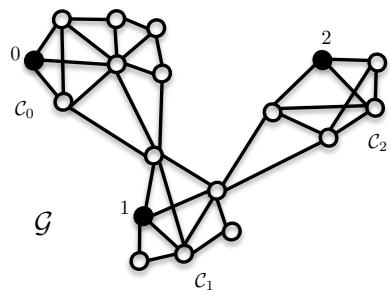
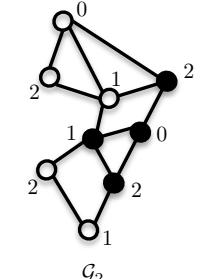
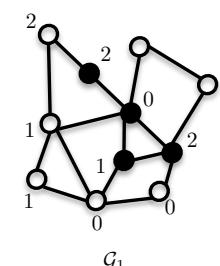
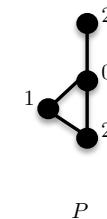
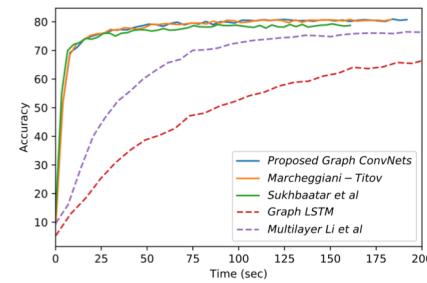
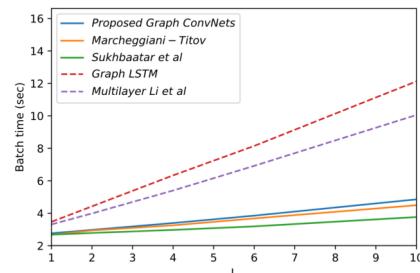
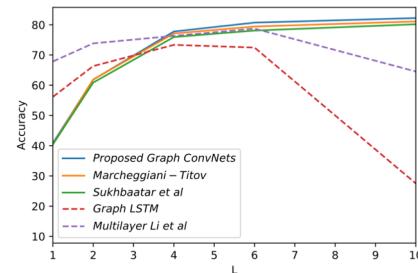
**Self-attention mechanism
in 1-hop neighborhood
/sparse attention**



[1] Velickovic, Cucurull, Casanova, Romero, Lio, Bengio, Graph Attention Networks, 2018

GraphRNNs vs GCNs^[1]

- Numerical comparison on two basic graph problems:
 - Sub-graph matching/recognition^[2]
 - Semi-supervised classification



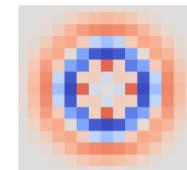
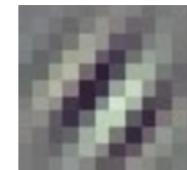
- GCNs offer better performances than GraphRNNs (as they can stack many layers).

[1] Bresson, Laurent, Residual gated graph convnets, 2017

[2] Scarselli, Gori, Tsoi, Hagenbuchner, Monfardini, The Graph Neural Network Model, 2009

Anisotropy vs Isotropy

- Standard ConvNets produce **anisotropic** filters because Euclidean grids have directional structure.
- Graph ConvNets compute **isotropic** filters because there is no notion of directions on arbitrary graphs.
- How to get anisotropy back in GNNs ?
 - Edge gates^[1]/attention mechanism^[2] information to treat neighbors differently.
 - Differentiate graph edges^[3] (e.g. different connections between atoms)



[1] Bresson, Laurent, Residual gated graph convnets, 2017

[2] Velickovic, Cucurull, Casanova, Romero, Lio, Bengio, Graph Attention Networks, 2018

[3] Gilmer, Schoenholz, Riley, Vinyals, Dahl, Neural message passing for quantum chemistry, 2017

Lab 02

- GatedGCNs for Sub-Graph Recognition
- GatedGCNs for Graph Classification

jupyter 01_GatedGCNs_subgraph_recognition Last Checkpoint: a few seconds ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 Logout

In [1]:

```
Residual Gated Graph ConvNets
X Bresson, T Laurent
ICLR 2017
ArXiv preprint: arXiv:1711.07553

import torch
from torch.autograd import Variable
import torch.nn.functional as F
import torch.nn as nn
import pdb
import time
import numpy as np
import pickle

if torch.cuda.is_available():
    print('cuda available')
    dtypeFloat = torch.cuda.FloatTensor
    dtypeLong = torch.cuda.LongTensor
    #torch.cuda.manual_seed(1)
else:
    print('cuda not available')
    dtypeFloat = torch.FloatTensor
    dtypeLong = torch.LongTensor
    #torch.manual_seed(1)

# import files in folder util
import sys
sys.path.insert(0, 'util/')
import block
import graph_generator as g

from sklearn.metrics import confusion_matrix
```

jupyter 02_GatedGCNs_graph_clustering Last Checkpoint: 37 minutes ago (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 Logout

In [1]:

```
Residual Gated Graph ConvNets
X Bresson, T Laurent
ICLR 2017
ArXiv preprint: arXiv:1711.07553

import torch
from torch.autograd import Variable
import torch.nn.functional as F
import torch.nn as nn
import pdb
import time
import numpy as np
import pickle

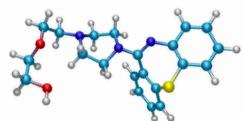
if torch.cuda.is_available():
    print('cuda available')
    dtypeFloat = torch.cuda.FloatTensor
    dtypeLong = torch.cuda.LongTensor
    #torch.cuda.manual_seed(1)
else:
    print('cuda not available')
    dtypeFloat = torch.FloatTensor
    dtypeLong = torch.LongTensor
    #torch.manual_seed(1)

# import files in folder util
import sys
sys.path.insert(0, 'util/')
import block
import graph_generator as g

from sklearn.metrics import confusion_matrix
```

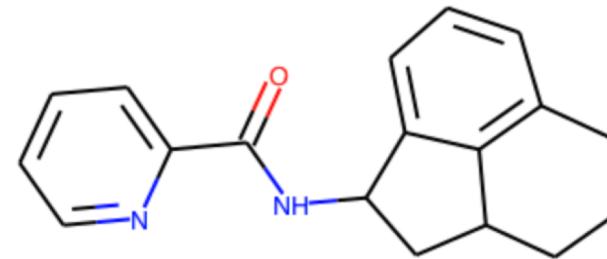
Outline

- Part 3: Spatial Graph ConvNets
 - Problem Setting
 - Spatial Graph Architectures
 - Quantum Chemistry
 - Operations Research
 - DGL



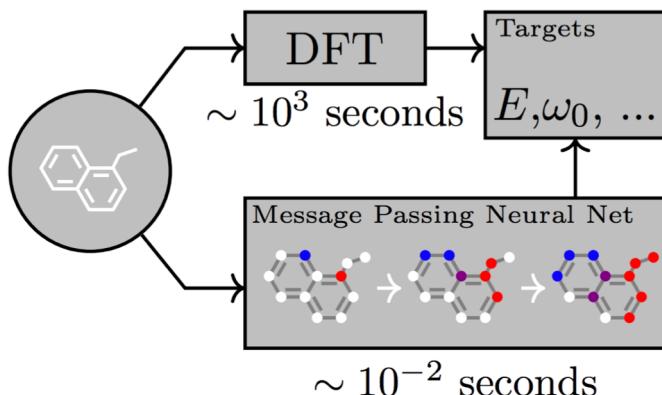
Quantum Chemistry

- Molecules are graphs where vertices are atoms (e.g. H, C, N, Cl) and edges are types of bonds between atoms (e.g. single, double, aromatic).
- Tasks are :
 - Computing chemical properties :
 - Solubility
 - Octanol-water partition coefficient
 - Synthetic accessibility
 - Generating new molecules with specific optimized properties



Molecule Regression

- Duvenaud-et-al NeurIPS'15 :
 - Learn molecule fingerprints
- Gilmer-Vinyals-Dahl-et.al 2017 (Google) :
 - Use **graph NNs** to learn 13 important chemical properties.
 - Standard estimation is done by **DFT** (approximation of coupled Schrodinger PDEs).

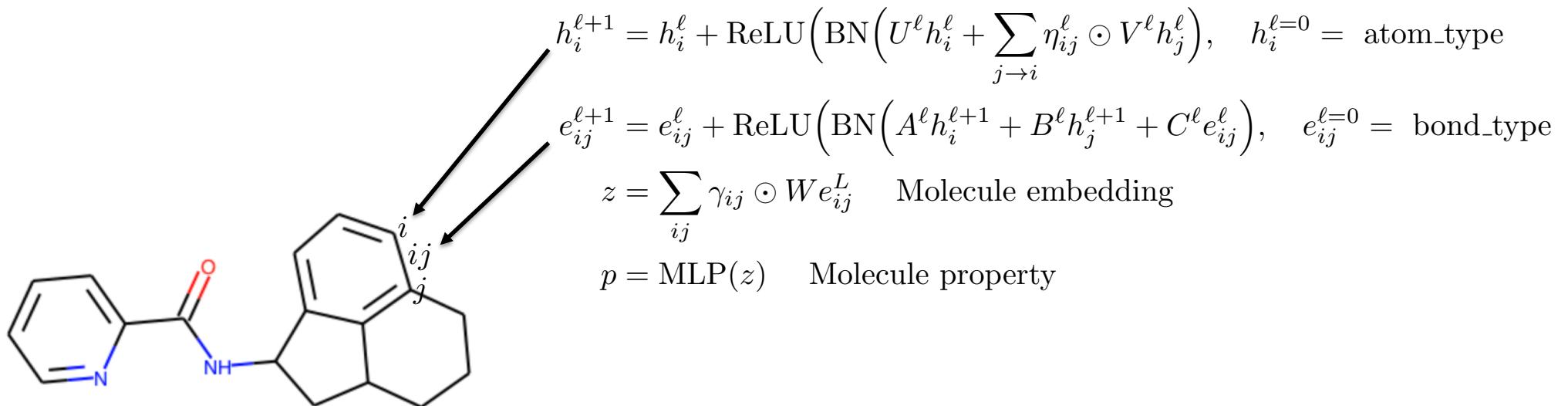


Chemical Properties	enn-s2s	SchNet	Message Passing
	Gilmer-etal'16 Google	Schutt-etal'17 TU Berlin	Jørgensen-etal'18 U. Danemark
alpha (Bohr ³)	0.092	0.235	0.077
Cv (cal/molK)	0.040	0.033	0.032
G (meV)	19	14	12.200
Gap (meV)	69	63	58.000
H (meV)	17	14	11.300
HOMO (meV)	43	41	36.700
LUMO (meV)	37	34	30.800
R ² (Bohr ²)	0.180	0.073	0.072
U (meV)	19	19	10.600
U0 (meV)	19	14	10.500
Zpve (meV)	1.5	1.7	1.490
mu (Debye)	0.030	0.033	0.029

Table 7: Mean absolute error of predictions for different target properties of the QM9 dataset using 110k molecules.

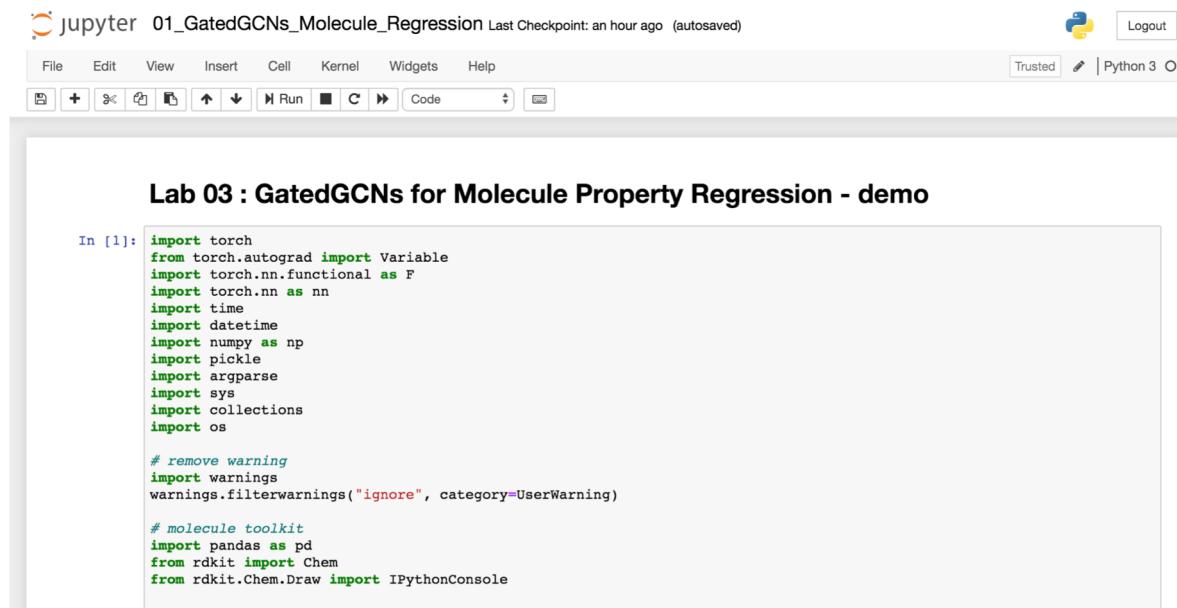
Molecule Property Regression with GatedGCNs

- Two steps :
 - Vertex, edge and graph embedding layers
 - Regression layer



Lab 03

- GatedGCNs for Molecule Property Regression



The screenshot shows a Jupyter Notebook interface with the following details:

- Title Bar:** jupyter 01_GatedGCNs_Molecule_Regression Last Checkpoint: an hour ago (autosaved)
- Toolbar:** File, Edit, View, Insert, Cell, Kernel, Widgets, Help, Trusted, Python 3
- Code Cell:** In [1]:

```
import torch
from torch.autograd import Variable
import torch.nn.functional as F
import torch.nn as nn
import time
import datetime
import numpy as np
import pickle
import argparse
import sys
import collections
import os

# remove warning
import warnings
warnings.filterwarnings("ignore", category=UserWarning)

# molecule toolkit
import pandas as pd
from rdkit import Chem
from rdkit.Chem.Draw import IPythonConsole
```
- Logout Button:** Logout

AI For Materials and Drugs Design

- Workshop “Generative and Reinforcement Learning with Physics” at IPAM-UCLA, Sept 2019
 - Apply Deep Learning to generate new materials and drugs (big impact in the coming years)
 - Self-driving lab to create these materials and drugs quicker and less expensively (takes a decade and costs \$350M-\$2.7B to bring a new drug to market)

<https://fortune.com/2019/09/17/artificial-intelligence-insilico-deepmind>

The screenshot shows the homepage of the workshop website. At the top, there's a navigation bar with 'Workshops' and a sub-menu 'Programs > Workshops > Workshop I: From Passive to Active: Generative and Reinforcement Learning with Physics'. Below the navigation is a large banner with a blue and white abstract design. The main content area has a light gray background. It features a section titled 'Workshop I: From Passive to Active: Generative and Reinforcement Learning with Physics' with a subtitle 'Part of the Long Program Machine Learning for Physics and the Physics of Learning'. The date 'SEPTEMBER 23 - 27, 2019' is listed. Below this is a 'Overview' section with a detailed text about the workshop's goals and constraints, accompanied by a small diagram of a neural network architecture. There are also sections for 'SPEAKER LIST', 'LOGGING', 'SCHEDULE', and 'APPLICATION & REGISTRATION'.

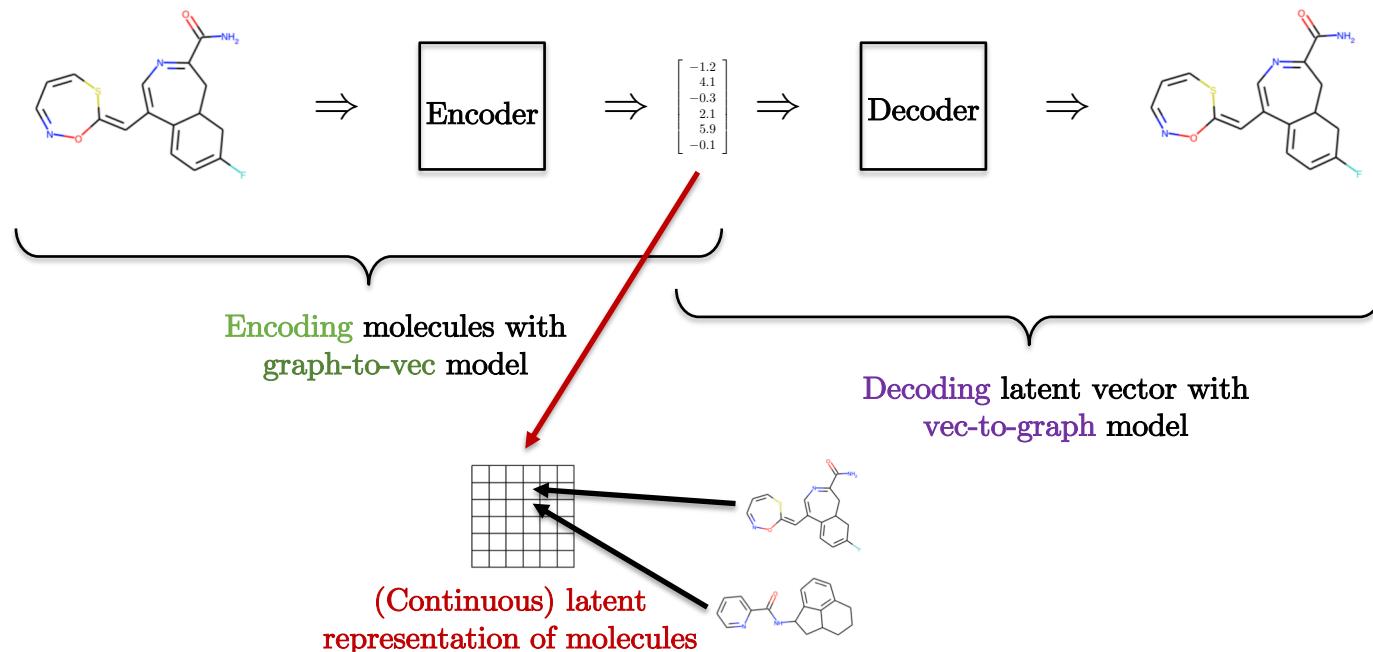
Speaker List

Josh Agar (Lehigh University)
Alain Aspuru-Guzik (University of Toronto)
Joshua Bloom (University of California, Berkeley (UC Berkeley))
Xavier Bresson (Nanyang Technological University, Singapore)
Giuseppe Carleo (Flatiron Institute, a Division of the Simons Foundation, Center for Computational
Juan Felipe Carrasquilla Alvarado (Vector Institute)
Nicholas Chiaromonte (Rice University, Physics and Astronomy)
Nicola De Ceo (University of Amsterdam)
Laurent Drine (Google)
Refael Gomez-Bombarelli (Massachusetts Institute of Technology)
Richard Hennig (University of Florida)
Oleksandr Isayev (Carnegie Mellon University)
Ross King (University of Science and Technology in Manchester (UMIST))
Philip Kohn (Howard University, Graduate School (Physics))
Joshua Yeochi Lin (University of Illinois at Urbana-Champaign, Physics)
Benjamin Nachman (Lawrence Berkeley National Laboratory)
Peter Nonnenmacher (University of Regensburg)
Ankit Patel (Rice University)
Han Pu (Rice University)
Patrick Reichenbach (Google)
Lars Ruthotto (Emory University)
Klaus Schenck (Cornell University)
Phala Shanahan (Massachusetts Institute of Technology)
Ganesh Sivaraman (Argonne National Laboratory, Argonne Leadership Computing Facility)
Tess Smits (Lawrence Berkeley National Laboratory, Physics)
Jascha Sohl-Dickstein (Google)
Alexandre Tkatchenko (University of Luxembourg, Theory)

The screenshot shows a news article from Fortune magazine. The title is 'Why These Two Innovations In Artificial Intelligence Are So Important: Eye on A.I.'. The author is Jeremy Kahn and Jonathan Vanian, and it was published on September 17, 2019. The article discusses two significant AI innovations. Below the main title is a photo of a man speaking at a podium with a microphone, wearing glasses and a dark suit. The background is a blue stage with the word 'FORTUNE' repeated. To the right of the main content, there are sections for 'Most Popular Posts' and 'Fact Checking Trump's Claims During One of the Most Chaotic Weeks in His Presidency'.

Graph Molecular Generation

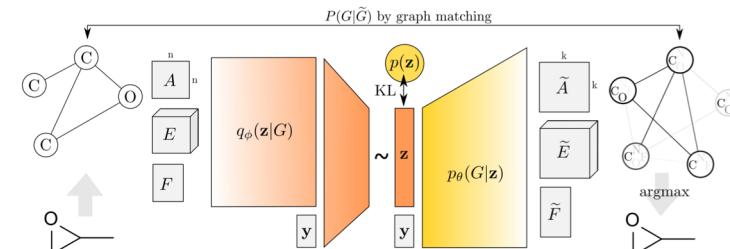
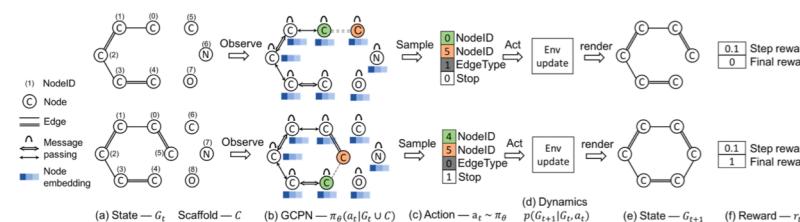
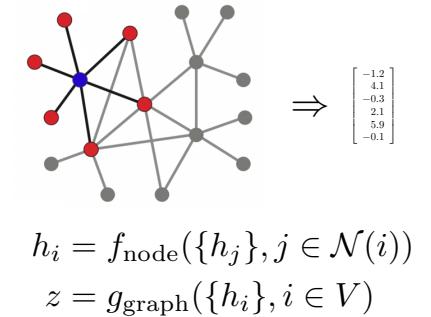
- Goal is to design a **neural network** that can **auto-encode molecules**, generate novel molecules, produce molecules with optimized chemical property.
- Graph VAE Architecture with GatedGCNs^[1] :



[1] Bresson, Laurent, A Two-Step Graph Convolutional Decoder for Molecule Generation, Workshop ML and Physics at NeurIPS'19

Graph Encoder and Decoder

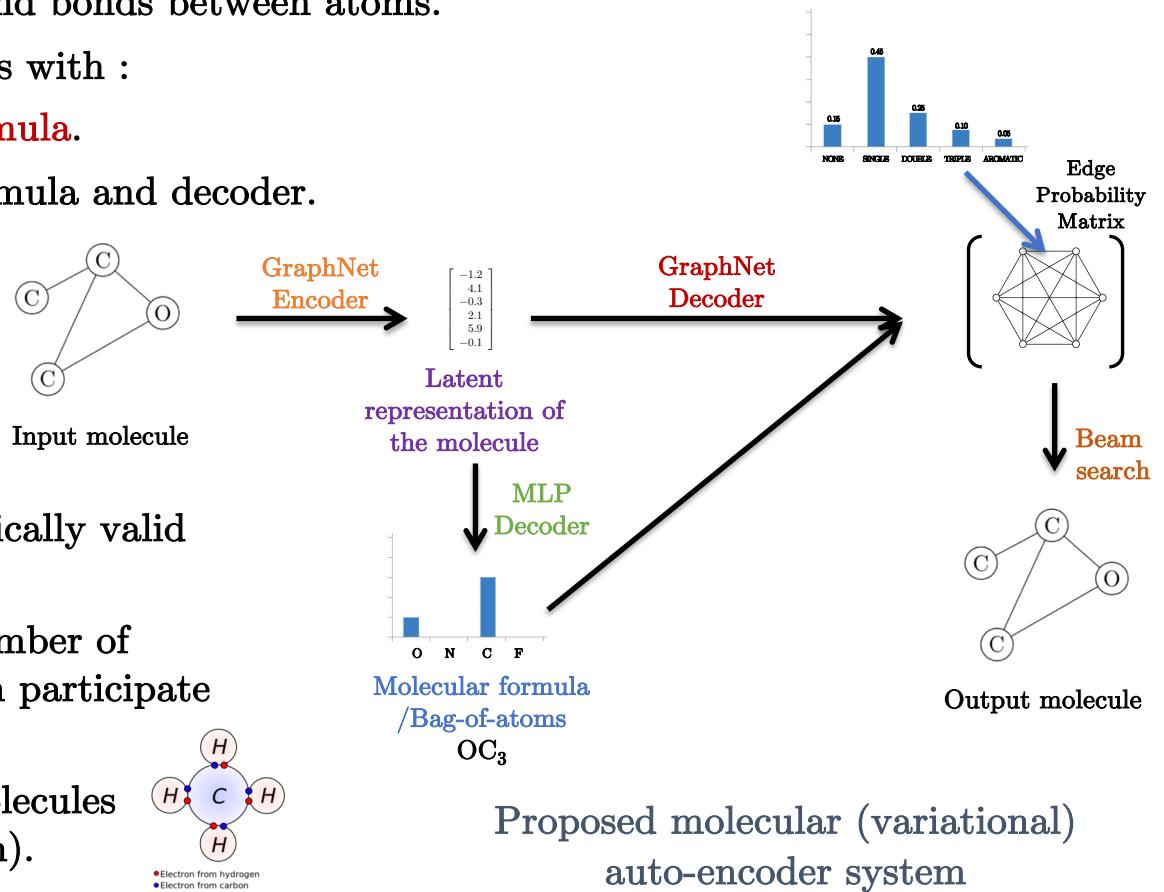
- Graph NNs^[1,2] used to encode molecules to predict molecular properties (1-2 orders of magnitude faster than DFT).
 - Encoding is easy. Decoding is more challenging !
- Two decoding approaches :
 - Auto-regressive models : Sequential generation of molecules (atom-by-atom)^[3,4].
 - One-shot models : Generation of all atoms and bonds in a single pass^[5,6].



- [1] Duvenaud, Maclaurin, Iparraguirre, Bombarell, Hirzel, Aspuru-Guzik, Adams, Convolutional networks on graphs for learning molecular fingerprints, 2015
 [2] Gilmer, Schoenholz, Riley, Vinyals, Dahl, Neural message passing for quantum chemistry, 2017
 [3] Jin, Barzilay, Jaakkola, Junction Tree Variational Autoencoder for Molecular Graph Generation, 2018
 [4] You, Liu, Ying, Pande, Leskovec, Graph convolutional policy network for goal-directed molecular graph generation, 2018
 [5] Simonovsky, Komodakis, GraphVAE: Towards generation of small graphs using variational autoencoders, 2018
 [6] De Cao, Kipf, MolGAN: An implicit generative model for small molecular graphs, 2018

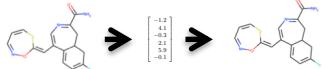
Single-Shot Molecular VAE

- It is hard to generate simultaneously atoms and bonds between atoms.
 - We propose to disentangle these problems with :
 - Atom generation with **molecular formula**.
 - Bond generation** using molecular formula and decoder.
- Advantages:
 - Simple and fast**
(GPU parallelizable) vs
auto-regressive models.
- The one-shot model may not produce a chemically valid molecule:
 - Violation of atom valency** (maximum number of electrons in the outer atom shell that can participate of a chemical bond).
 - We use **beam search** to produce valid molecules (borrowed from NLP sequence generation).



Numerical Experiments

- Molecule reconstruction
 - How many molecules are correctly decoded?



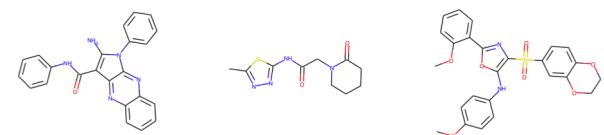
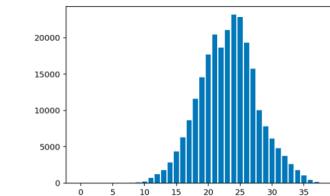
Method	Reconstruction	Validity
CVAE [Gomez-Bombarelli et al., 2016]	44.6%	0.7%
GVAE [Kusner et al., 2017]	53.7%	7.2%
SD-VAE [Dai et al, 2018]	76.2%	43.5%
GraphVAE [Simonovsky, Komodakis, 2018]	-	13.5%
JT-VAE (SL) [Jin et al, 2018]	76.7%	100.0%
GCPN (GAN+RL) [You et al, 2018]	-	-
OURS (VAE+SL)	90.5%	100.0%

Table 1: Percentage of successful reconstruction of 250k ZINC molecules.

- Molecule novelty
 - Beyond memorization – how many molecules sampled from the learned distribution are not in the training set?

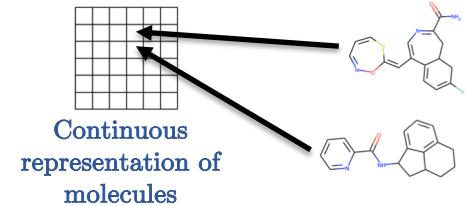
Method	Novelty	Uniqueness
JT-VAE (SL) [Jin et al, 2018]	100.0%	100.0%
GCPN (GAN+RL) [You et al, 2018]	-	-
OURS (VAE+SL)	100.0%	100.0%

Table 2: Sample 5000 molecules from learned prior distribution.



ZINC (250k drug like molecules, up to 38 heavy atoms)

Numerical Experiments



- Molecule optimization
 - How much property improvement can we obtain when optimizing in the latent space?
 - Optimization of solubility constrained by the similarity between the original molecule and the new generated molecule.

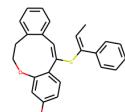
δ	JT-VAE [Jin et al, 2018] (SL)			GCPN [You et al, 2018] (GAN+RL)			OURS (VAE+SL)		
	Improvement	Similarity	Success	Improvement	Similarity	Success	Improvement	Similarity	Success
0.0	1.91 ± 2.04	0.28 ± 0.15	97.5%	4.20 ± 1.28	0.32 ± 0.12	100.0%	5.24 ± 1.55	0.18 ± 0.12	100.0%
0.2	1.68 ± 1.85	0.33 ± 0.13	97.1%	4.12 ± 1.19	0.34 ± 0.11	100.0%	4.29 ± 1.57	0.31 ± 0.12	98.6%
0.4	0.84 ± 1.45	0.51 ± 0.10	83.6%	2.49 ± 1.30	0.47 ± 0.08	100.0%	3.05 ± 1.46	0.51 ± 0.10	84.0%
0.6	0.21 ± 0.71	0.69 ± 0.06	46.4%	0.79 ± 0.63	0.68 ± 0.08	100.0%	2.46 ± 1.27	0.67 ± 0.05	40.1%

Table 7: Molecule optimization results.

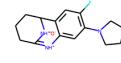
Molecule similarity 0.4



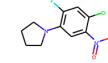
604: -2.94



604: 3.39

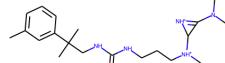


103: -5.40

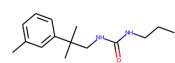


103: 0.88

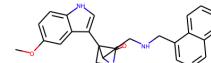
Molecule similarity 0.6



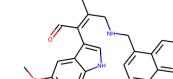
89: -5.64



89: 0.94



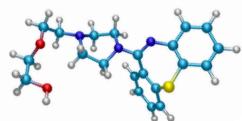
782: -2.57



782: 2.44

Outline

- Part 3: Spatial Graph ConvNets
 - Problem Setting
 - Spatial Graph Architectures
 - Quantum Chemistry
 - Operations Research
 - DGL



Operations Research

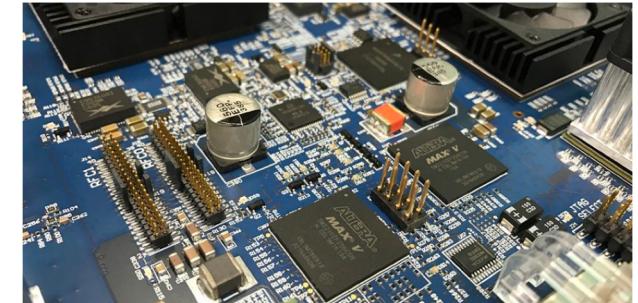
- OR/combinatorial problems s.a. assignment, routing, planning, supply chain, scheduling are used every day in revenue management, transportation, manufacturing, supply chain, public policy, hardware design, etc.



Amazon warehouse
management



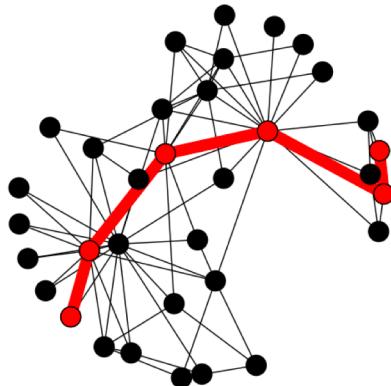
Uber taxis
allocation



Hardware pieces
placement

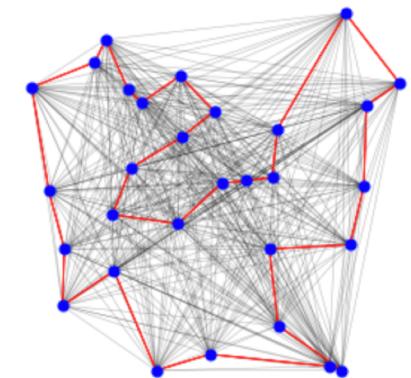
Easy Combinatorial Problems

- Shortest Path Problem :
 - Find a path between two vertices in a graph with shortest distance.
 - Dijkstra's algorithm : Exact solver in polynomial time $O(E+V\log V)$ with Fibonacci heap.
 - Applications : Road Navigator, P2P network, etc



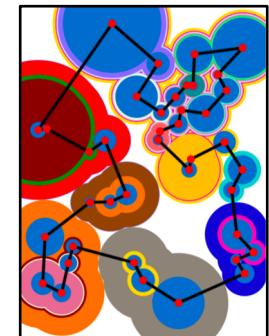
Hard Combinatorial Problems

- Travelling Salesman Problem (TSP) :
 - Find the shortest tour that visits each node and returns to the starting node.
 - NP-hard problem in combinatorial optimization.
 - Exact solution in factorial time $O(V!)$
 - Long history : Von Neumann'51, Dantzig'54, Bellman-Held-Karp'62, Edmonds'67, etc.
 - William Cook (Waterloo U.), MAA Invited Talk on TSP, Jan. 2018 : <https://www.youtube.com/watch?v=q8nQTNvCrjE>
 - Applications : Transportation, scheduling, hardware design, etc.



TSP

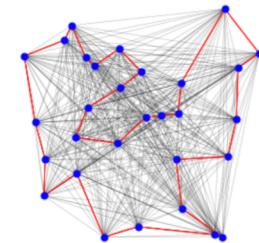
- State-of-the-art solution
 - Concorde^[1] solver developed by Applegate, Bixby, Chvatal, Cook, 2006 (test-of-time).
 - Code available :
 - <http://www.math.uwaterloo.ca/tsp/concorde>
 - Leveraged 80 years of theoretical developments, data structures and heuristics from computer science.
 - Branch-and-Cut approach :
 - Cutting plane algorithms (Dantzig et al., 1954) to iteratively solve linear relaxations of TSP.
 - Branch-and-bound approach to reduce solution search space.



[1] Applegate, Bixby, Chvatal, Cook, 2006

Operations Research and Deep Learning

- Most OR problems are difficult/NP-hard, require heuristic solutions with years of research work and significant specialized knowledge.
 - Popular TSP problem has been studied for more than 80 years w/ von Neumann, Dantzig, etc. Best TSP solver is Concorde^[1] (solve linear relaxations of TSP w/ branch-and-bound).
 - Any small change like adding constraints in TSP would require a new development of the specialized algorithm. OR designs hand-crafted algorithms.
- DL may lead to a potential breakthrough in traditional OR !
 - DL can learn universal high-quality algorithms to solve any OR problem by adapting to any continuous/discrete structure and by continuously learning by reinforcement and optimization of the cost under any constraint^[2,3].



[1] Applegate, Bixby, Chvatal, Cook, 2006

[2] Bengio, Lodi, Prouvost, Machine learning for combinatorial optimization: a methodological tour d'horizon, 2018

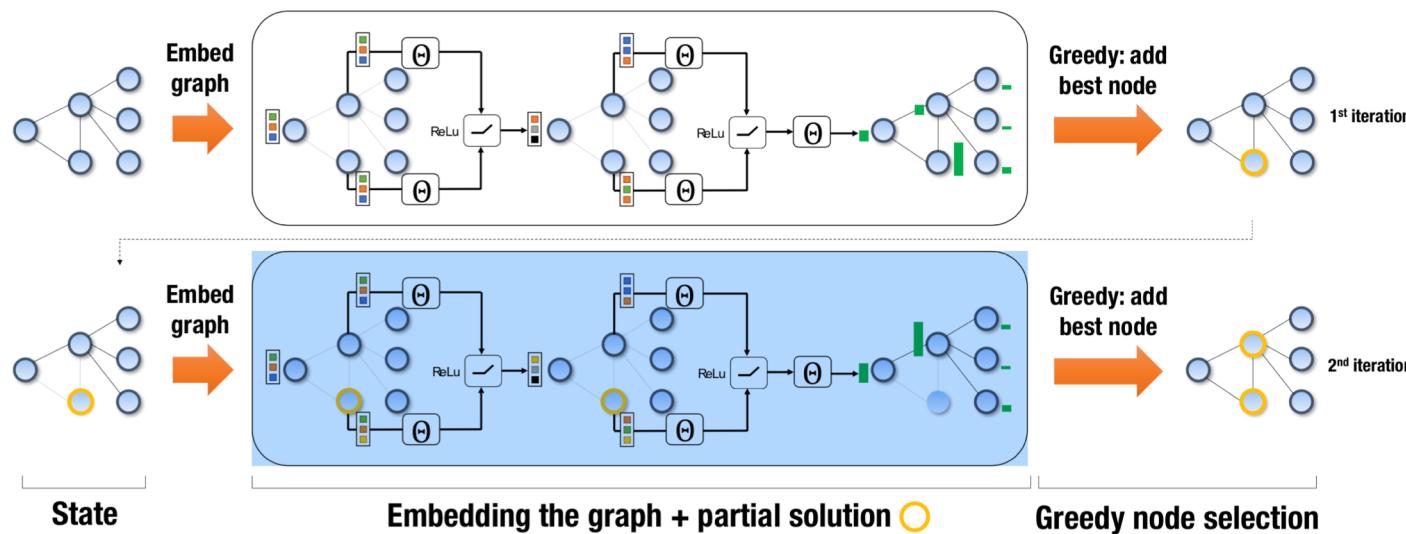
[3] Bello, Pham, V Le, Norouzi, Bengio, Neural combinatorial optimization with reinforcement learning, 2016

Neural Networks for TSP

- Recent efforts to design learning algorithms for TSP (proof-of-concept).
- Two approaches :
 - Auto-regressive models : Sequential generation of TSP solution.
 - One-shot models : Generation of edge probability to be in TSP solution + sampling.

Auto-regressive NNs

- Auto-regressive models : Sequential generation of TSP solution^[1,2,3,4].



[1] Vinyals, Fortunato, Jaitly, Pointer networks, 2015

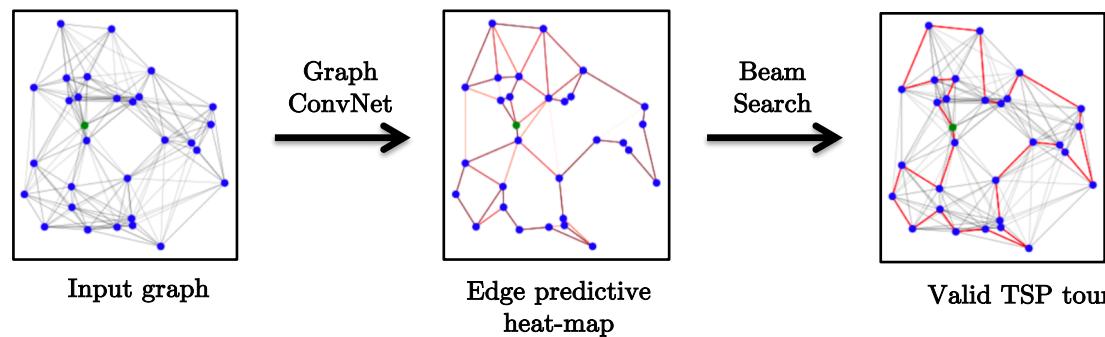
[2] Bello, Pham, V Le, Norouzi, Bengio, Neural combinatorial optimization with reinforcement learning, 2016

[3] Dai, Khalil, Zhang, Dilkina, Song, Learning combinatorial optimization algorithms over graphs, 2017

[4] Kool, van Hoof, Welling, Attention, learn to solve routing problems!, 2019

One-shot NNs

- One-shot models : Generation of edge probability to be in TSP solution + sampling^[1].
 - GatedGCNs to predict edge probability + beam search to sample valid solution^[2]



[1] Nowak, Villar, Bandeira, Bruna, A note on learning algorithms for quadratic assignment with graph neural networks, 2017

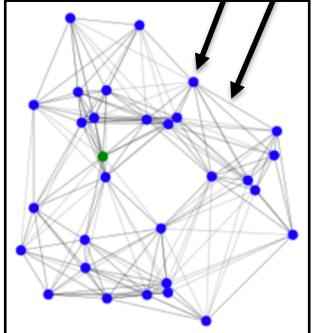
[2] CK Joshi, T Laurent, X Bresson, An Efficient Graph Convolutional Network Technique for the Travelling Salesman Problem, Workshop on GNNs at NeurIPS'19

GatedGCNs for TSP^[1,2]

Node and edge representations

$$\left\{ \begin{array}{l} h_i^{\ell+1} = h_i^\ell + \text{ReLU}\left(\text{BN}\left(W_1^\ell h_i^\ell + \sum_{j \sim i} \eta_{ij}^\ell \odot W_2^\ell h_j^\ell\right)\right) \quad \text{with} \quad \eta_{ij}^\ell = \frac{\sigma(e_{ij}^\ell)}{\sum_{j' \sim i} \sigma(e_{ij'}^\ell) + \varepsilon} \\ e_{ij}^{\ell+1} = e_{ij}^\ell + \text{ReLU}\left(\text{BN}\left(V_1^\ell e_{ij}^\ell + V_2^\ell h_i^\ell + V_3^\ell h_j^\ell\right)\right) \end{array} \right.$$

Dense attention



Edge classifier

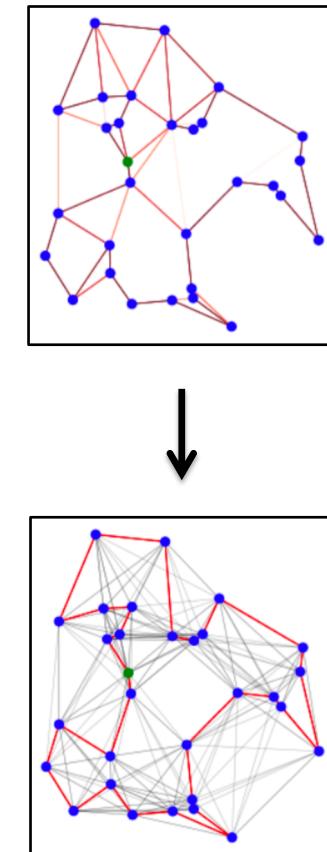
$$\left\{ \begin{array}{l} h_i^{\ell=0} = x_i \in \mathbb{R}^2 \\ e_{ij}^{\ell=0} = \|x_i - x_j\|_2 \\ p_{ij} = \text{sigmoid}(\text{MLP}(e_{ij}^L)) \quad \text{with loss being the logistic regression on edges.} \end{array} \right.$$

[1] Bresson, Laurent, Residual gated graph convnets, ICLR'17

[2] CK Joshi, T Laurent, X Bresson, An Efficient Graph Convolutional Network Technique for the Travelling Salesman Problem, Workshop on GNNs at NeurIPS'19

Beam Search Decoder

- Output of GNN :
 - Edge probability function over the adjacency matrix of tour connections.
- Beam search decoder (NLP) :
 - Starting from a random node, expand b most probable edge connections among the node's neighbors.
 - Keep expanding the top- b partial tours at each stage until all nodes are visited.
 - Final prediction is the tour with the smallest length among the b tours at the end of search.



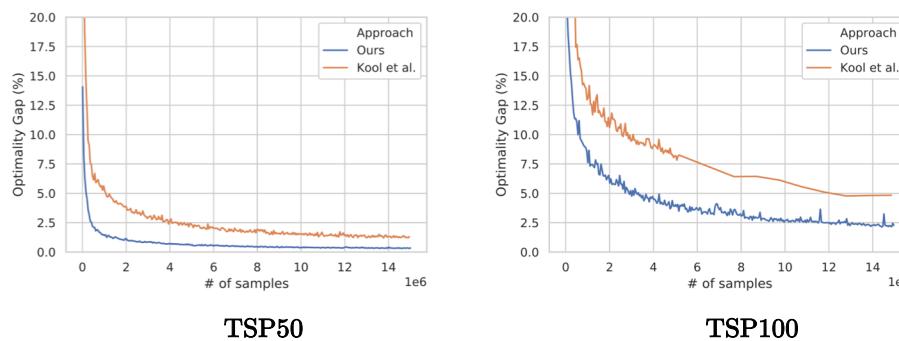
Numerical Experiments

- Our system :
 - Outperforms previous NNs in terms of both optimality gap and evaluation time (graph ConvNet and beam search are parallelizable).
 - But does not outperform Concorde.

$$\text{gap} = \frac{1}{m} \sum_{i=1}^m \left(l_m / \hat{l}_m - 1 \right)$$

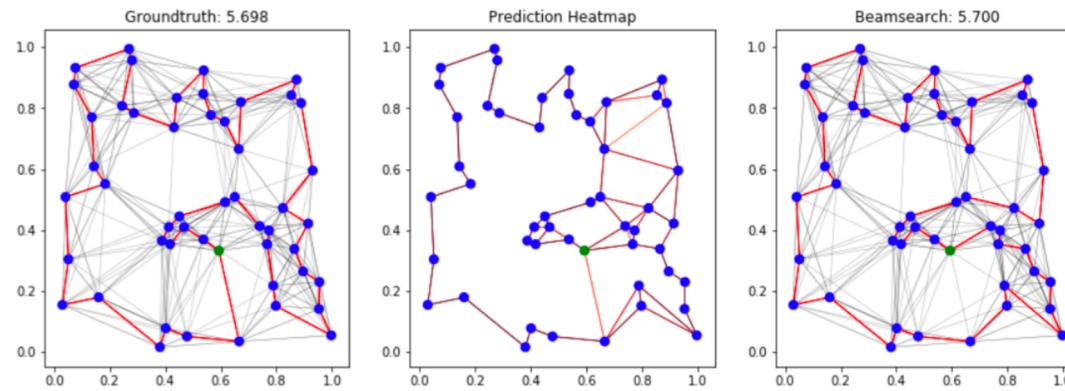
	Method	Type	TSP20			TSP50			TSP100		
			Tour Len.	Opt. Gap.	Time	Tour Len.	Opt. Gap.	Time	Tour Len.	Opt. Gap.	Time
OR techniques	Concorde	Solver	3.84	0.00%	(1m)	5.70	0.00%	(2m)	7.76	0.00%	(3m)
	Gurobi	Solver	3.84	0.00%	(7s)	5.70	0.00%	(2m)	7.76	0.00%	(17m)
	OR Tools	H, S	3.85	0.37%		5.80	1.83%		7.99	2.90%	
NN techniques	GNN [Nowak et al., 2017]	SL, BS	3.93	2.46%		-	-		-	-	
	PtrNet [Bello et al., 2016]	RL, S	-	-		5.75	0.95%		8.00	3.03%	
	GAT [Kool et al., 2019]	RL, S	3.84	0.08%	(5m)	5.73	0.52%	(24m)	7.94	2.26%	(1h)
	GCN (Ours)	SL, BS*	3.84	0.01%	(12m)	5.70	0.01%	(18m)	7.87	1.39%	(40m)

Sample efficiency

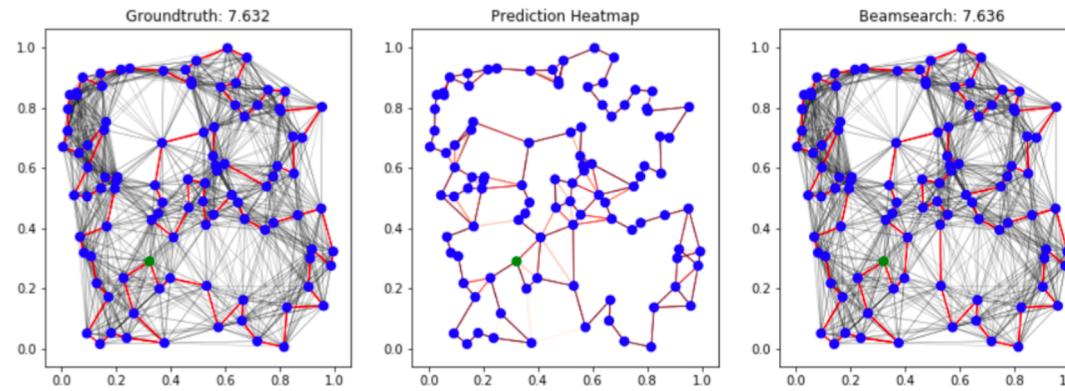


Visualization

TSP50

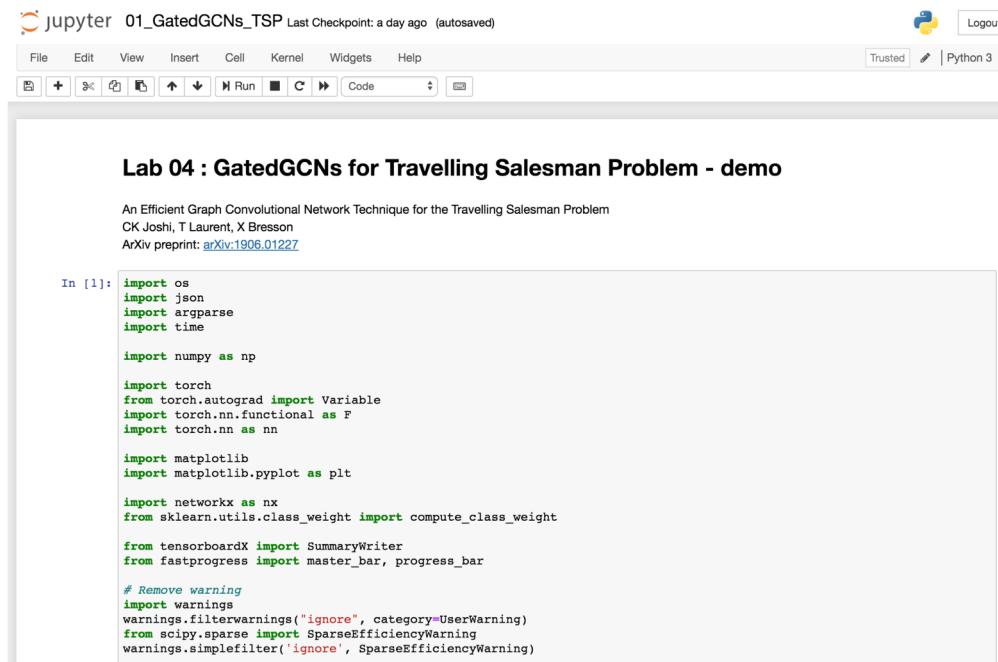


TSP100



Lab 04

- GatedGCNs for Travelling Salesman Problem



The screenshot shows a Jupyter Notebook interface with the following details:

- Title Bar:** jupyter 01_GatedGCNs_TSP Last Checkpoint: a day ago (autosaved)
- Toolbar:** File, Edit, View, Insert, Cell, Kernel, Widgets, Help, Run, Code, Logout, Trusted, Python 3
- Section Title:** Lab 04 : GatedGCNs for Travelling Salesman Problem - demo
- Text:** An Efficient Graph Convolutional Network Technique for the Travelling Salesman Problem
CK Joshi, T Laurent, X Bresson
ArXiv preprint: [arXiv:1906.01227](https://arxiv.org/abs/1906.01227)
- Code Cell (In [1]):**

```
import os
import json
import argparse
import time

import numpy as np

import torch
from torch.autograd import Variable
import torch.nn.functional as F
import torch.nn as nn

import matplotlib
import matplotlib.pyplot as plt

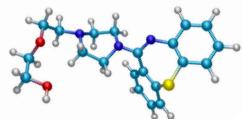
import networkx as nx
from sklearn.utils.class_weight import compute_class_weight

from tensorboardX import SummaryWriter
from fastprogress import master_bar, progress_bar

# Remove warning
import warnings
warnings.filterwarnings("ignore", category=UserWarning)
from scipy.sparse import SparseEfficiencyWarning
warnings.simplefilter('ignore', SparseEfficiencyWarning)
```

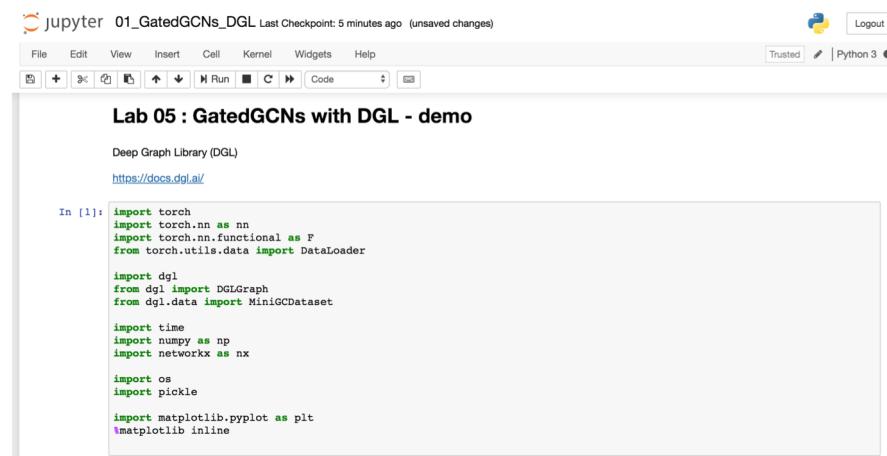
Outline

- Part 3: Spatial Graph ConvNets
 - Problem Setting
 - Spatial Graph Architectures
 - Quantum Chemistry
 - Operations Research
 - DGL



Lab 05

- GatedGCNs with DGL (Deep Graph Library), NYU
<https://docs.dgl.ai>



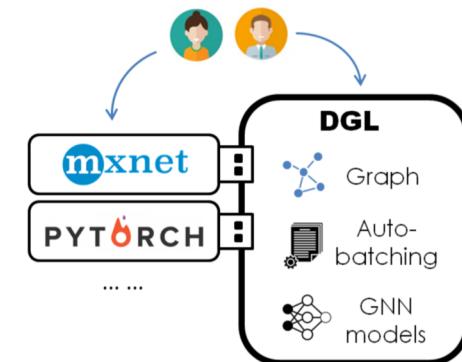
```
jupyter 01_GatedGCNs_DGL Last Checkpoint: 5 minutes ago (unsaved changes)
File Edit View Insert Cell Kernel Widgets Help Logout Trusted Python 3
Lab 05 : GatedGCNs with DGL - demo
Deep Graph Library (DGL)
https://docs.dgl.ai/
In [1]: import torch
import torch.nn as nn
import torch.nn.functional as F
from torch.utils.data import DataLoader

import dgl
from dgl import DGLGraph
from dgl.data import MiniGCNDataset

import time
import numpy as np
import networkx as nx

import os
import pickle

import matplotlib.pyplot as plt
%matplotlib inline
```



Lab 05

$\frac{1}{\sqrt{V_k}}, \frac{1}{\sqrt{E_k}}$ Normalization constants for better optimization

DGL creates a batch of graphs

Create artificial node feature (input degree) and edge feature (value 1)

```
# collate function
def collate(samples):
    graphs, labels = map(list, zip(*samples)) # samples is a list of pairs (graph, label).
    labels = torch.tensor(labels)
    tab_sizes_n = [graphs[i].number_of_nodes() for i in range(len(graphs))] # graph sizes
    tab_snorm_n = [torch.FloatTensor(size,1).fill_(1./float(size)) for size in tab_sizes_n]
    snorm_n = torch.cat(tab_snorm_n).sqrt() # normalization constant for better optimization
    tab_sizes_e = [graphs[i].number_of_edges() for i in range(len(graphs))] # nb of edges
    tab_snorm_e = [torch.FloatTensor(size,1).fill_(1./float(size)) for size in tab_sizes_e]
    snorm_e = torch.cat(tab_snorm_e).sqrt() # normalization constant for better optimization
    batched_graph = dgl.batch(graphs) # batch graphs
    return batched_graph, labels, snorm_n, snorm_e

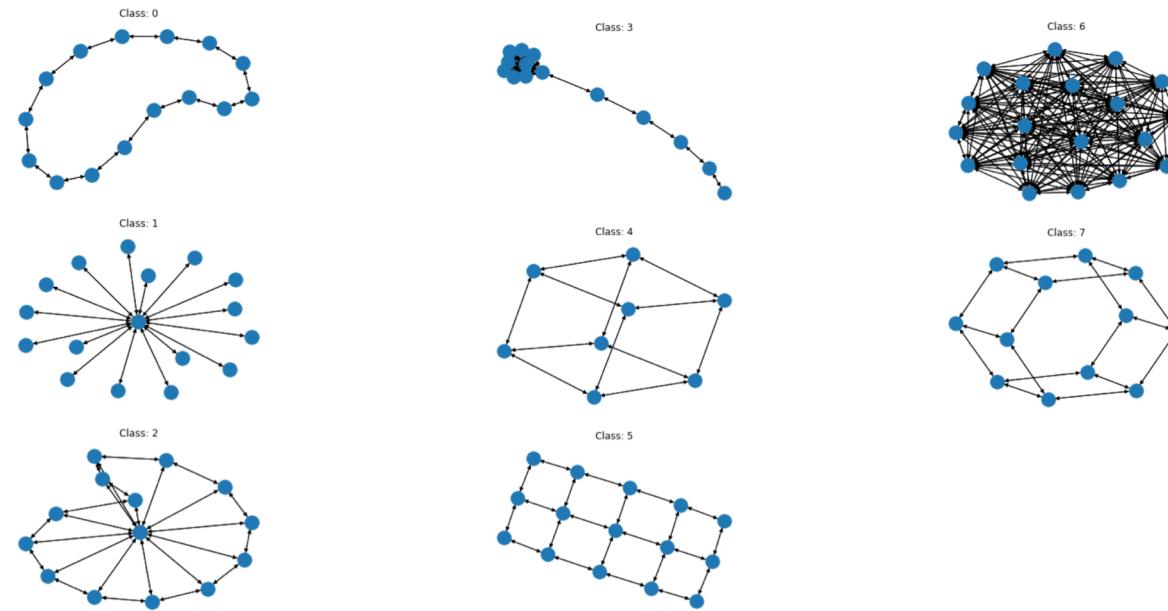
# create artificial data feature (= in degree) for each node
def create_artificial_features(dataset):
    for (graph,_) in dataset:
        graph.ndata['feat'] = graph.in_degrees().view(-1, 1).float()
        graph.edata['feat'] = torch.ones(graph.number_of_edges(),1)
    return dataset

# use artificial graph dataset of DGL
trainset = MiniGCDataset(8, 10, 20)
trainset = create_artificial_features(trainset)
print(trainset[0])

(DGLGraph(num_nodes=13, num_edges=39,
         nndata_schemes={'feat': Scheme(shape=(1,), dtype=torch.float32)},
         edata_schemes={'feat': Scheme(shape=(1,), dtype=torch.float32)}), 0)
```

Lab 05

```
# use artificial graph dataset of DGL  
trainset = MiniGCDataset(8, 10, 20)
```

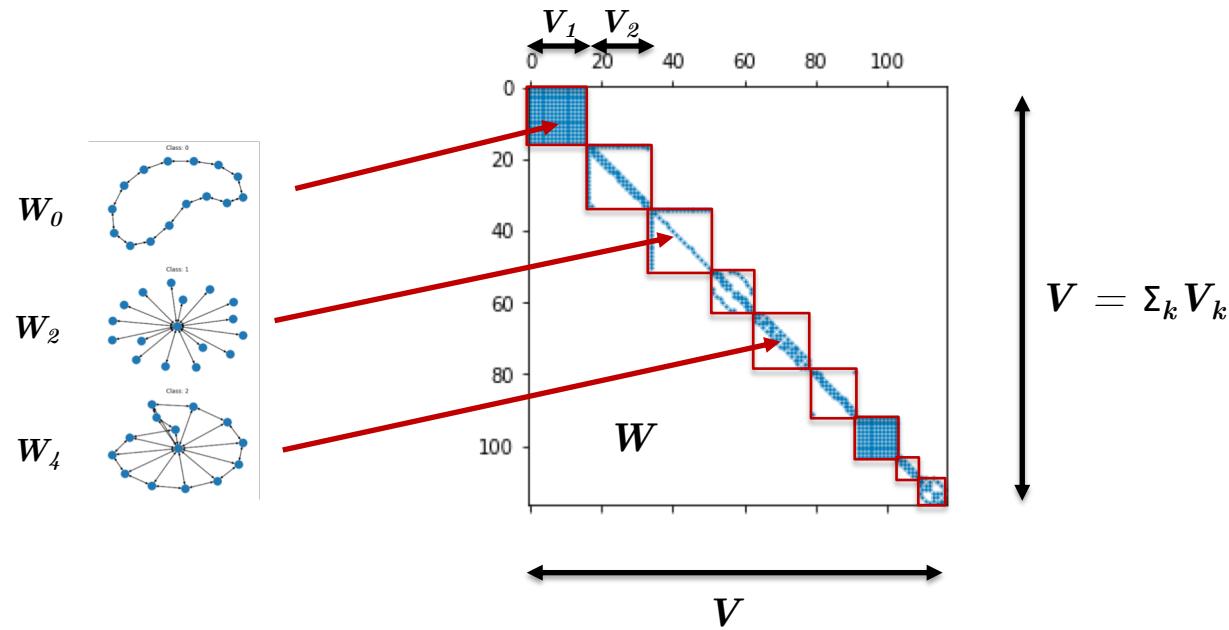


Lab 05

- Understanding DGL :

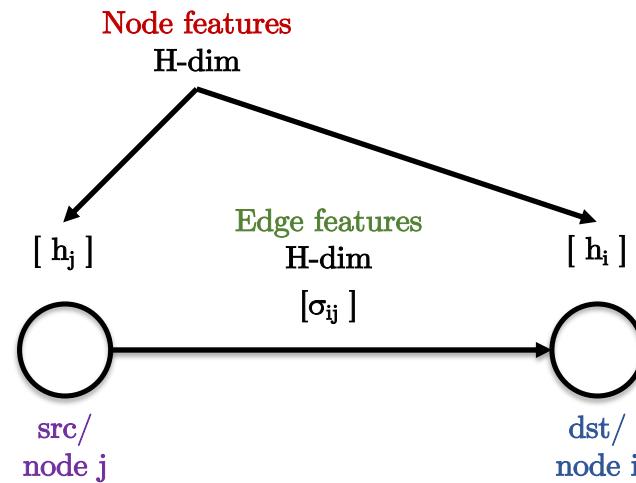
- How to process K graphs of different lengths?

Form a (big) sparse block diagonal matrix W with K adjacency matrices W_k .



Lab 05

- Basic structure of an **edge** in DGL :

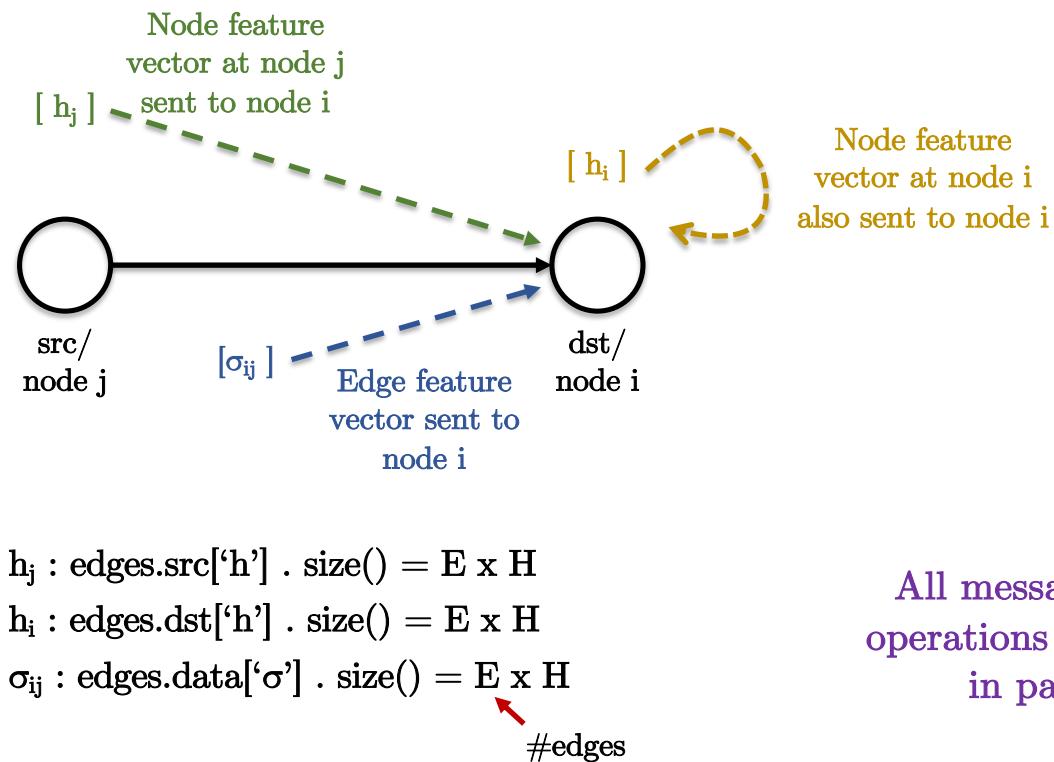


- Goal is to **compute efficiently** expressions of the **form** :

$$f_i = h_i + \sum_{j \rightarrow i} \sigma_{ij} \circ h_j$$

Lab 05

- Step 1 : Message passing function defined on edges
 - Node feature and edge feature are passed along all edges connecting a node.



Lab 05

- Step 2 : Reduce function defined on nodes

- Reduce functions of the form : $f_i = h_i + \sum_{j \rightarrow i} \sigma_{ij} \circ h_j$
- Reduce function collects all messages passed in Step 1. #nodes
- Code : $f = h_i + \text{torch.sum}(h_j \times \sigma_{ij}, \text{dim}=1) . \text{size()} = V \times H$
Sum over neighbors

- GPU acceleration :

- DGL batches the nodes with the same number of neighbors.

$$h_j = \text{nodes.mailbox}[h_j] = \left\{ \begin{array}{l} \# \text{nodes in batch}_1 \\ \text{batch}_1 . \text{size}() = 11 \times 12 \times H \\ \vdots \\ \# \text{neighbors} \\ \text{batch}_{34} . \text{size}() = 14 \times 9 \times H \end{array} \right\}$$

$\sigma_{ij} = \text{nodes.mailbox}[\sigma_{ij}] =$ same structure than h_j

$h_i = \text{nodes.data}[h] . \text{size()} = V \times H$

Lab 05

GatedGCN layer :

$$h_i^{\ell+1} = h_i^\ell + \text{ReLU}\left(\text{BN}\left(A^\ell h_i^\ell + \sum_{j \sim i} \eta(e_{ij}^\ell) \odot B^\ell h_j^\ell\right)\right)/\sqrt{V_k},$$

$$\eta(e_{ij}^\ell) = \frac{\sigma(e_{ij}^\ell)}{\sum_{j' \sim i} \sigma(e_{ij'}^\ell) + \varepsilon},$$

$$e_{ij}^{\ell+1} = e_{ij}^\ell + \text{ReLU}\left(\text{BN}\left(C^\ell e_{ij}^\ell + D^\ell h_i^{\ell+1} + E^\ell h_j^{\ell+1}\right)\right)/\sqrt{E_k}.$$

```

class GatedGCN_layer(nn.Module):
    def __init__(self, input_dim, output_dim):
        super(GatedGCN_layer, self).__init__()
        self.A = nn.Linear(input_dim, output_dim, bias=True)
        self.B = nn.Linear(input_dim, output_dim, bias=True)
        self.C = nn.Linear(input_dim, output_dim, bias=True)
        self.D = nn.Linear(input_dim, output_dim, bias=True)
        self.E = nn.Linear(input_dim, output_dim, bias=True)
        self.bn_node_h = nn.BatchNorm1d(output_dim)
        self.bn_node_e = nn.BatchNorm1d(output_dim)

    def message_func(self, edges):
        Bh_j = edges.src['Bh']
        e_ij = edges.data['Ce'] + edges.src['Dh'] + edges.dst['Eh'] # e_ij = Ce_ij + Dhi + Ehj
        edges.data['e'] = e_ij
        return {'Bh_j': Bh_j, 'e_ij': e_ij}

    def reduce_func(self, nodes):
        Ah_i = nodes.data['Ah']
        Bh_j = nodes.mailbox['Bh_j']
        e = nodes.mailbox['e_ij']
        sigma_ij = torch.sigmoid(e) # sigma_ij = sigmoid(e_ij)
        h = Ah_i + torch.sum(sigma_ij * Bh_j, dim=1) / torch.sum(sigma_ij, dim=1) # hi = Ah_i + sum_j eta_ij * Bj
        return {'h': h}

    def forward(self, g, h, e, snorm_n, snorm_e):
        h_in = h # residual connection
        e_in = e # residual connection

        g.ndata['h'] = h
        g.ndata['Ah'] = self.A(h)
        g.ndata['Bh'] = self.B(h)
        g.ndata['Dh'] = self.D(h)
        g.ndata['Eh'] = self.E(h)
        gedata['e'] = e
        g.data['Ce'] = self.C(e)
        g.update_all(self.message_func, self.reduce_func)
        h = g.ndata['h'] # result of graph convolution
        e = g.ndata['e'] # result of graph convolution

        h = h * snorm_n # normalize activation w.r.t. graph node size
        e = e * snorm_e # normalize activation w.r.t. graph edge size

        h = self.bn_node_h(h) # batch normalization
        e = self.bn_node_e(e) # batch normalization

        h = F.relu(h) # non-linear activation
        e = F.relu(e) # non-linear activation

        h = h_in + h # residual connection
        e = e_in + e # residual connection

    return h, e

```

Lab 05

MLP classifier layer

```
class MLP_layer(nn.Module):

    def __init__(self, input_dim, output_dim, L=2): # L = nb of hidden layers
        super(MLP_layer, self).__init__()
        list_FC_layers = [ nn.Linear( input_dim, input_dim, bias=True ) for l in range(L) ]
        list_FC_layers.append(nn.Linear( input_dim, output_dim , bias=True ))
        self.FC_layers = nn.ModuleList(list_FC_layers)
        self.L = L

    def forward(self, x):
        y = x
        for l in range(self.L):
            y = self.FC_layers[l](y)
            y = F.relu(y)
        y = self.FC_layers[self.L](y)
        return y
```

GatedGCN Network

Node input embedding

Edge input embedding

Run graphNN layers

Compute graph vectorial representation by a (simple) average of all node features with DGL.

Use MLP classifier

```
class GatedGCN_Net(nn.Module):

    def __init__(self, net_parameters):
        super(GatedGCN_Net, self).__init__()
        input_dim = net_parameters['input_dim']
        hidden_dim = net_parameters['hidden_dim']
        output_dim = net_parameters['output_dim']
        L = net_parameters['L']
        self.embedding_h = nn.Linear(input_dim, hidden_dim)
        self.embedding_e = nn.Linear(1, hidden_dim)
        self.GatedGCN_layers = nn.ModuleList([ GatedGCN_layer(hidden_dim, hidden_dim) for _ in range(L) ])
        self.MLP_layer = MLP_layer(hidden_dim, output_dim)

    def forward(self, g, h, e, snorm_n, snorm_e):

        # input embedding
        h = self.embedding_h(h)
        e = self.embedding_e(e)

        # graph convnet layers
        for GGCN_layer in self.GatedGCN_layers:
            h,e = GGCN_layer(g,h,e,snorm_n,snorm_e)

        # MLP classifier
        g.ndata['h'] = h
        y = dgl.mean_nodes(g, 'h')
        y = self.MLP_layer(y)

        return y
```

Lab 05

```
def train_one_epoch(net, data_loader):
    """
    train one epoch
    """
    net.train()
    epoch_loss = 0
    epoch_train_acc = 0
    nb_data = 0
    gpu_mem = 0
    for iter, (batch_graphs, batch_labels, batch_snorm_n, batch_snorm_e) in enumerate(data_loader):
        batch_x = batch_graphs.ndata['feat'].to(device)
        batch_e = batch_graphs.edata['feat'].to(device)
        batch_snorm_n = batch_snorm_n.to(device)
        batch_snorm_e = batch_snorm_e.to(device)
        batch_labels = batch_labels.to(device)
        batch_scores = net.forward(batch_graphs, batch_x, batch_e, batch_snorm_n, batch_snorm_e)
        gpu_mem = net.gpu_memory(gpu_mem)
        loss = net.loss(batch_scores, batch_labels)
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
        epoch_loss += loss.detach().item()
        epoch_train_acc += net.accuracy(batch_scores, batch_labels)
        nb_data += batch_labels.size(0)
    epoch_loss /= (iter + 1)
    epoch_train_acc /= nb_data
    return epoch_loss, epoch_train_acc, gpu_mem
```

Train function

```
# datasets
train_loader = DataLoader(trainset, batch_size=50, shuffle=True, collate_fn=collate)
test_loader = DataLoader(testset, batch_size=50, shuffle=False, collate_fn=collate)
val_loader = DataLoader(valset, batch_size=50, shuffle=False, drop_last=False, collate_fn=collate)

# Create model
net_parameters = {}
net_parameters['input_dim'] = 1
net_parameters['hidden_dim'] = 100
net_parameters['output_dim'] = 8 # nb of classes
net_parameters['L'] = 4
net = GatedGCN_Net(net_parameters)
net = net.to(device)

optimizer = torch.optim.Adam(net.parameters(), lr=0.0005)

epoch_train_losses = []
epoch_test_losses = []
epoch_val_losses = []
epoch_train_accs = []
epoch_test_accs = []
epoch_val_accs = []
for epoch in range(50):

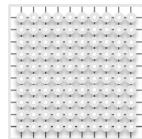
    start = time.time()
    epoch_train_loss, epoch_train_acc, gpu_mem = train_one_epoch(net, train_loader)
    epoch_test_loss, epoch_test_acc = evaluate_network(net, test_loader)
    epoch_val_loss, epoch_val_acc = evaluate_network(net, val_loader)

    print('Epoch {}, time {:.4f}, train_loss: {:.4f}, test_loss: {:.4f}, val_loss: {:.4f} \n
```

Main function

Outline

- Part 1: Traditional ConvNets



- Part 2: Spectral Graph ConvNets



- Part 3: Spatial Graph ConvNets



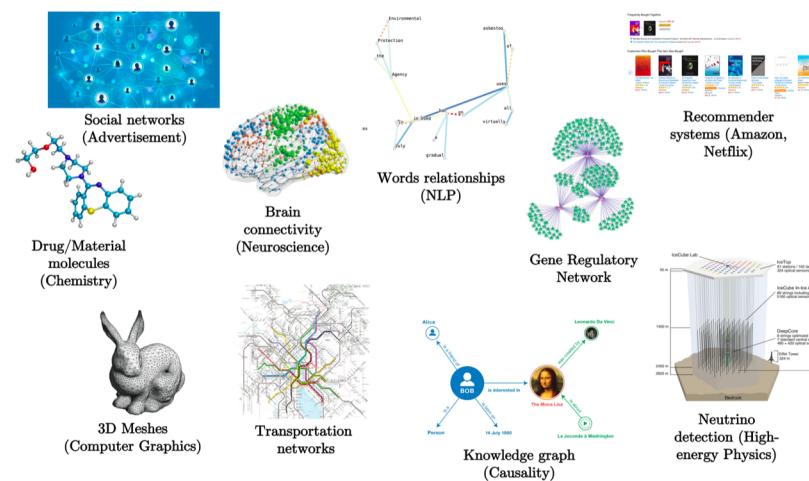
- Conclusion

Conclusion

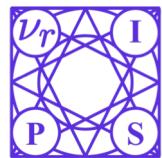
- **Contributions :**

- Generalization of ConvNets to data on graphs
- Re-design compositional property on graphs
- Linear complexity for sparse graphs
- GPU implementation (not yet optimized for sparse matrix-matrix multiplications)
- Universal learning capacity
- Multiple and dynamic graphs

- **Applications :**



Tutorials on Graph Deep Learning



NeurIPS'17

1,000-2,000 participants



CVPR'17

500-1,000 participants

- Organizers :

- Y. LeCun (NYU, Chief AI Scientist at Facebook)
- M. Bronstein (Imperial, Twitter)
- J. Bruna (NYU)
- A. Szlam (Facebook AI Research)
- X. Bresson (NTU, Singapore)

Annual Meeting



July 9-13, 2018
Portland, Oregon, USA
Oregon Convention Center (OCC)



SIAM Annual Meeting'18

2018 IPAM-UCLA Workshop

Organizing committee:

- Y. LeCun (NYU, Chief AI Scientist at Facebook)
- M. Bronstein (Imperial, Twitter)
- J. Bruna (NYU)
- A. Szlam (Facebook AI Research)
- S. Osher (UCLA)
- X. Bresson (NTU, Singapore)

The screenshot shows the homepage of the 'Workshops' section for the 'New Deep Learning Techniques' event. At the top, there's a dark banner with the word 'Workshops' and a decorative graphic of concentric circles and dots. Below it, a white header bar includes links for 'Programs', 'Workshops', and 'New Deep Learning Techniques'. The main content area has a light gray background. It features the title 'New Deep Learning Techniques' and the date 'FEBRUARY 5 - 9, 2018'. Below this are four navigation tabs: 'OVERVIEW' (highlighted in blue), 'SPEAKER LIST', 'SCHEDULE', and 'APPLY/REGISTER'. The 'OVERVIEW' tab contains a detailed paragraph about the challenges of applying deep learning to irregular domains and a small image of a colorful, abstract network graph. The 'SPEAKER LIST' tab is visible but not expanded.

Video talks : <https://lnkd.in/fY2-9UU>

Speaker List

- Alán Aspuru-Guzik (Harvard University)
Samuel Bowman (New York University)
Xavier Bresson (Nanyang Technological University, Singapore)
Michael Bronstein (USI Lugano, Switzerland, / Tel Aviv University, Israel / Intel Perceptual
Joan Bruna (New York University)
Pratik Chaudhari (University of California, Los Angeles (UCLA))
Kyle Cranmer (New York University)
Michael Elad (Technion - Israel Institute of Technology, Computer Science Department)
Sanja Fidler (University of Toronto)
Emily Fox (University of Washington)
Tom Goldstein (University of Maryland)
Leonidas Guibas (Stanford University, Computer Science)
Yann LeCun (New York University, Canadian Institute for Advanced Research)
Jure Leskovec (Stanford University)
Stéphane Mallat (École Normale Supérieure)
Federico Monti (Università della Svizzera Italiana)
Stanley Osher (University of California, Los Angeles (UCLA), Mathematics)
Ellie Pavlick (University of Pennsylvania)
Daniel Rueckert (Imperial College)
Ruslan Salakhutdinov (Carnegie-Mellon University)
Zuwei Shen (National University of Singapore, mathematics)
Stefano Soatto (University of California, Los Angeles (UCLA))
Sainbayar Sukhbaatar (New York University)
Arthur Szlam (Facebook)
Raquel Urtasun (University of Toronto)
Wei Zhu (Duke University, Mathematics)

2019 IPAM-UCLA Workshop

Organizing committee:

- Y. LeCun (NYU, Chief AI Scientist at Facebook)
- Rene Vidal (Johns Hopkins, Director Data Science Institute)
- Rebecca Willett (Chicago U.)
- S. Osher (UCLA)
- X. Bresson (NTU, Singapore)

The screenshot shows the workshop's website interface. At the top, there's a navigation bar with 'Programs', 'Workshops', and 'Workshop IV: Deep Geometric Learning of Big Data and Applications'. Below the navigation is a blue header image featuring a complex network of nodes and connections. The main content area has a dark background with white text. It displays the title 'Workshop IV: Deep Geometric Learning of Big Data and Applications', the subtitle 'Part of the Long Program Geometry and Learning from Data in 3D and Beyond', and the date 'MAY 20 - 24, 2019'. There are four tabs at the top of this section: 'OVERVIEW' (selected), 'SPEAKER LIST', 'SCHEDULE', and 'APPLY/REGISTER'. Below these tabs is a section titled 'Overview' with a detailed description of the workshop's goals and applications. To the right of the text is a small image of a colorful, abstract geometric visualization.

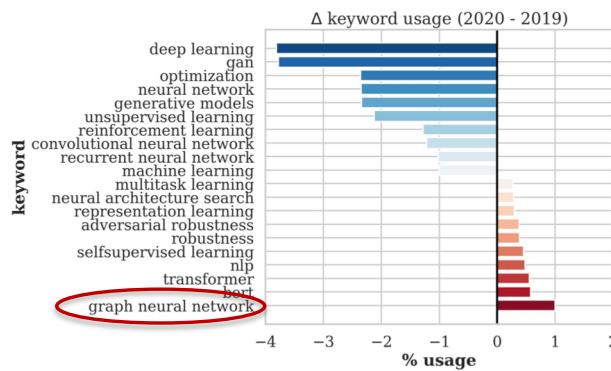
Video talks : <https://bit.ly/2w8EtLV>

Speaker List

Mikhail Belkin (Ohio State University)
Xavier Bresson (Nanyang Technological University, Singapore)
Soumith Chintala (Facebook AI Research)
Taco Cohen (Qualcomm AI Research)
Kostas Daniilidis (University of Pennsylvania)
Tom Goldstein (University of Maryland)
Bahram Jalali (University of California, Los Angeles (UCLA))
Thomas Kipf (Universiteit van Amsterdam)
Roy Lederman (Yale University, Applied Mathematics)
Jure Leskovec (Stanford University)
Federico Monti (Università della Svizzera Italiana)
Mathias Niepert (NEC Laboratories Europe)
Stanley Osher (University of California, Los Angeles (UCLA), Mathematics)
Hamed Pirsiavash (University of Maryland Baltimore County)
Marc Pollefeys (ETH Zurich)
Srikumar Ramalingam (University of Utah)
Thiago Serra (Mitsubishi Electric Research Laboratories (Merl))
Jeremias Sulam (Johns Hopkins University)
Arthur Szlam (Facebook)
Jian Tang (HEC Montréal)
Luc Van Gool (ETH Zurich)
Rene Vidal (Johns Hopkins University)
Ersin Yumer (Uber ATG)
Hongyang Zhang (Carnegie Mellon University)

Graph Deep Learning

- Top AI/DL conferences organize workshops/tutorials on “Graph Neural Networks” :
 - NeurIPS’19, ICML’19, ICLR’19
 - CVPR’19, ICCV’19
 - Etc



ICLR’20

Xavier Bresson @xbresson · Mar 28
Congratulations to visionaries @ylecun, Y. Bengio and @geoffreyhinton. This is so much deserved! Deep learning is changing our view of science. It opens many new opportunities for everyone in research, technology and philosophy. Wonderful news!

Yann LeCun @ylecun · Mar 28
I am extremely honored to be the recipient of the 2018 ACM A.M. Turing Award, and absolutely delighted to be sharing it with my friends and colleagues Geoffrey Hinton and Yoshua Bengio.
NYT... nytimes.com/2019/03/27/tec...

Yann LeCun @ylecun · Mar 28
Merci Xavier!

Michael Bronstein @mmbronstein · Mar 28
I am honored to now have a paper coauthored with a Turing awardee :-) - very well deserved Yann, congratulations!

Yann LeCun @ylecun
Replying to @mmbronstein @xbresson and @geoffreyhinton
Thank you, Michael. Someone, at some point, is going to win a Turing for geometric deep learning.



Questions?