

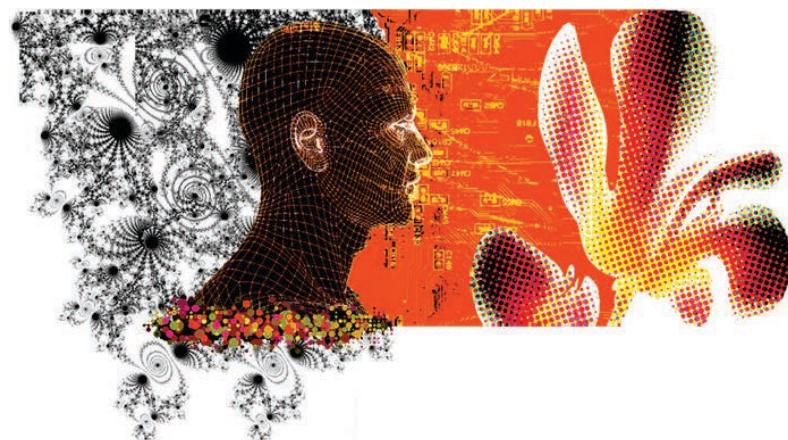
Data Science Training

November 2017

Deep Learning on Graphs

Xavier Bresson

Data Science and AI Research Centre
NTU, Singapore



<http://data-science-optum17.tk>

Outline

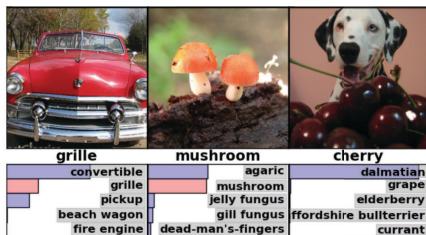
- **Convolutional Neural Networks**
 - *Architecture*
 - *Non-Euclidean Data*
- **Spectral Graph Theory**
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- **Spectral ConvNets**
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- **Spectral ConvNets on Multiple Graphs**
 - *Recommender Systems*
- **Conclusion**

Outline

- **Convolutional Neural Networks**
 - *Architecture*
 - *Non-Euclidean Data*
- **Spectral Graph Theory**
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- **Spectral ConvNets**
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- **Spectral ConvNets on Multiple Graphs**
 - *Recommender Systems*
- **Conclusion**

Convolutional neural networks^[1]

- Data (image, video, sound) are **compositional**, i.e. they are formed of patterns that are:
 - Local
 - Stationary
 - Multi-scale (/hierarchical)
 - ConvNets leverage the compositionality structure: They extract compositional features and feed them to classifier, recommender, etc (end-to-end).



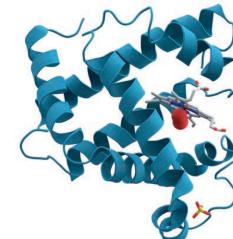
Computer Vision

```

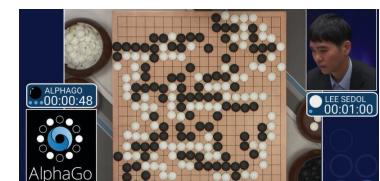
/* Duplicate LSM field information. The lsm_rule is opaque, so
 * we re-initialize it. */
static inline int audit_dupe_lsm_field(struct audit_field *df,
                                      struct audit_field *sf)
{
    int ret = 0;
    /* If same LSM str */
    /* Our own copy of lsm str */
    if (sf->lsm_rule.lsm == df->lsm_rule.lsm) {
        /* Invalid if lsm str is invalid */
        if (!lsm_is_valid(sf->lsm_str))
            return -EINVAL;
        /* If same type */
        if (sf->lsm_rule.type == df->lsm_rule.type)
            /* Our own (refreshed) copy of lsm_rule */
            ret = security_audit_rule_init(df->type, df->op, df->lsm_str,
                                           sf->lsm_rule.advice);
        /* Keep current valid entries around in case they
         * become valid after a policy reload. */
        if (sf->valid_till != ~0U) {
            /* If same timestamp */
            if (sf->valid_till == df->valid_till) {
                pr_warn("audit rule file for LSM '%s' is invalid\n",
                       df->lsm_str);
                ret = 0;
            }
        }
    }
    return ret;
}

```

NLP



Drug discovery

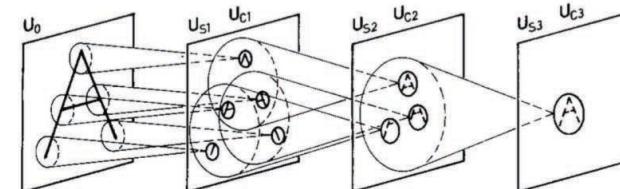
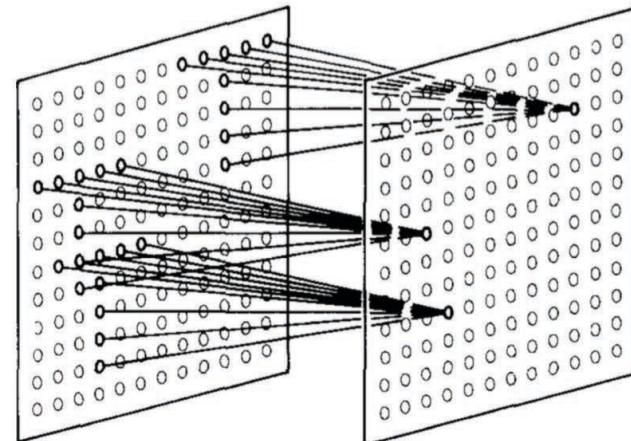
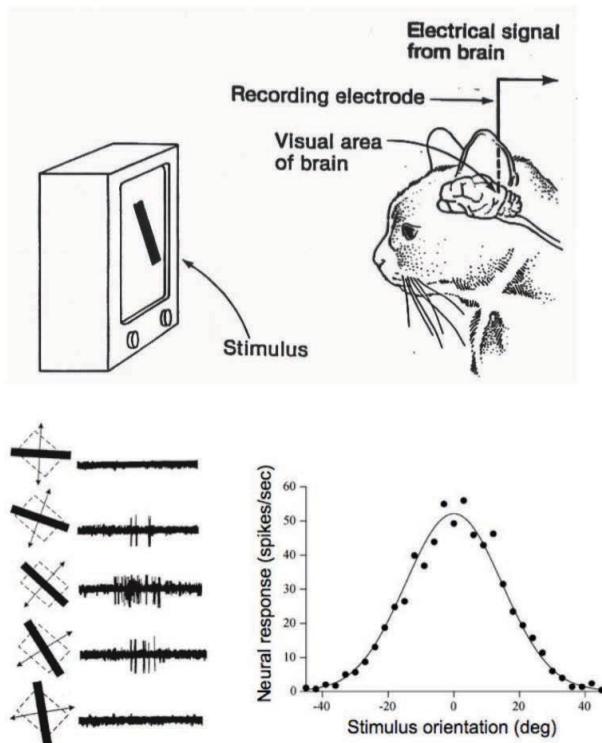


Games

[1] LeCun et al. 1998

Key properties

- **Locality:** Property inspired by visual cortex neurons^[7].
- **Local receptive fields** activate in the presence of (learned) local features.

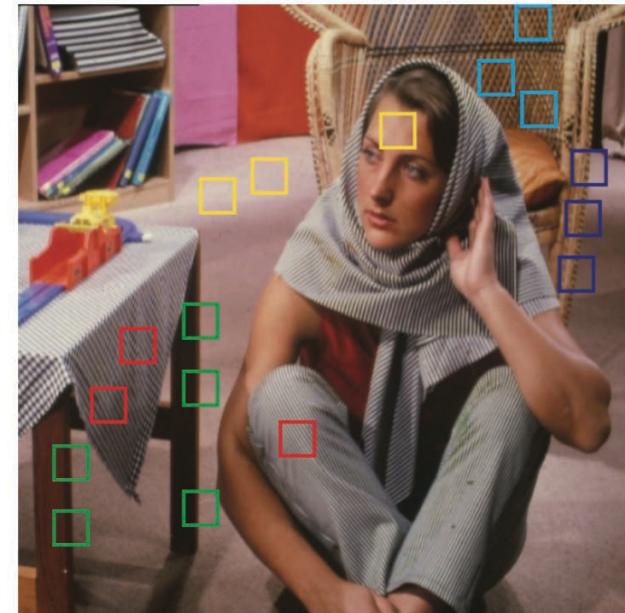
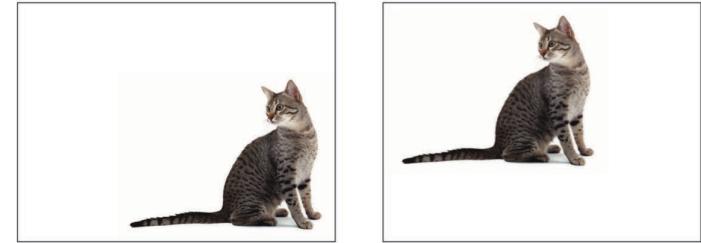


Neocognitron^[8]

[7] Cybenko 1989, Hornik 1991, [8] Fukushima 1980

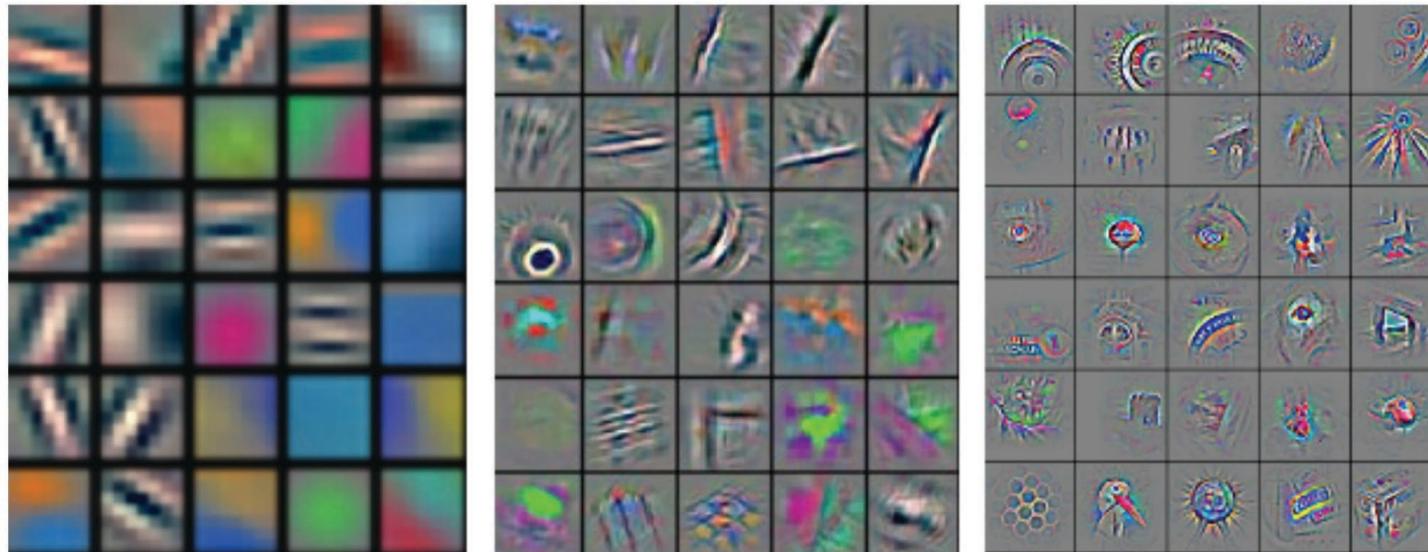
Key properties

- Stationarity \Leftrightarrow Translation invariance
- Local stationarity \Leftrightarrow Similar patches shared across the data domain.



Key properties

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

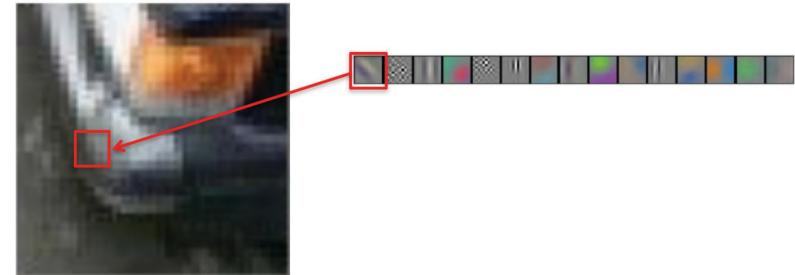


Features learned by a ConvNet become increasingly complex at deeper layers.^[9]

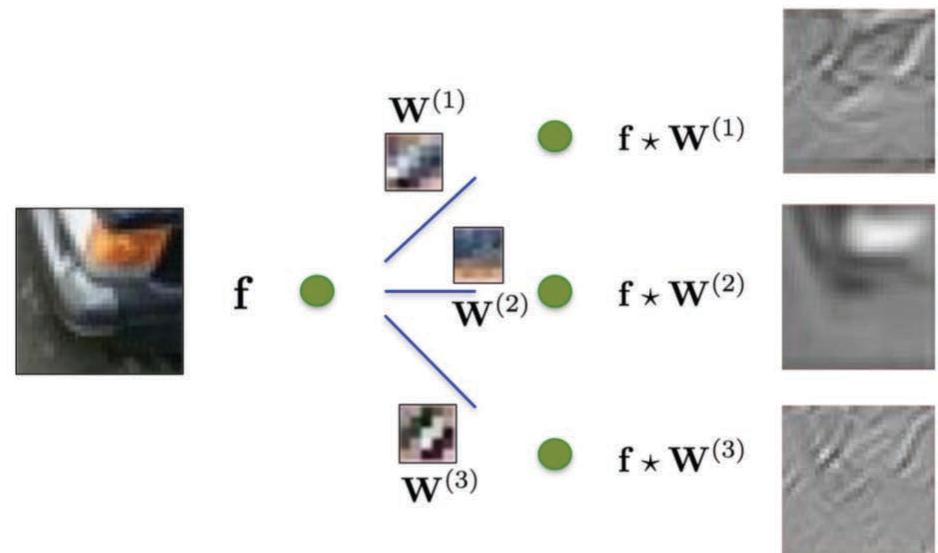
[9] Zeiler, Fergus 2013

Implementation

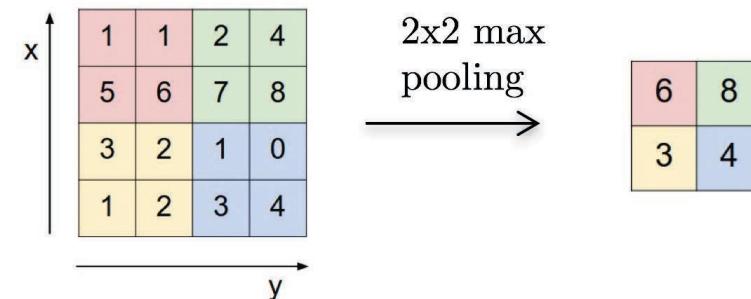
- **Locality:** compact support kernels
 $\Rightarrow O(1)$ parameters per filter.



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

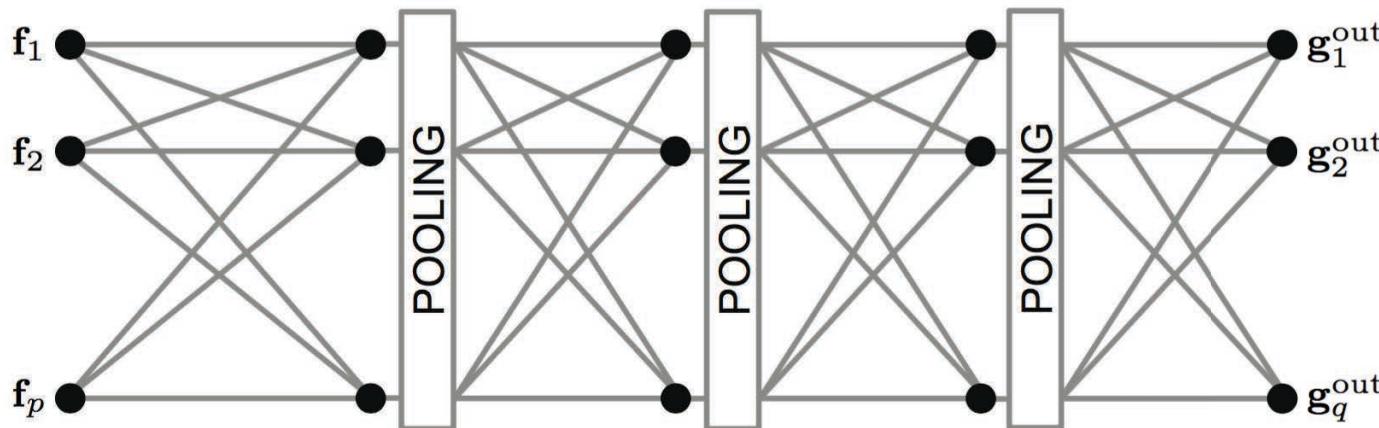


- **Multi-scale:** downsampling + pooling $\Rightarrow O(n)$



Compositional features

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



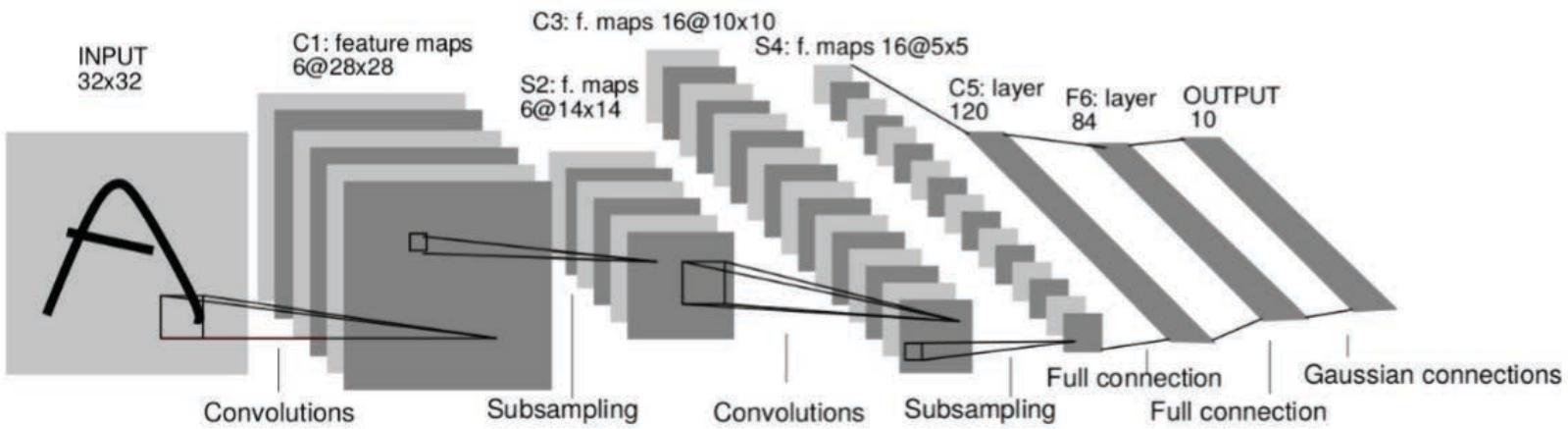
Compositional features consist of multiple convolutional + pooling layers.

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

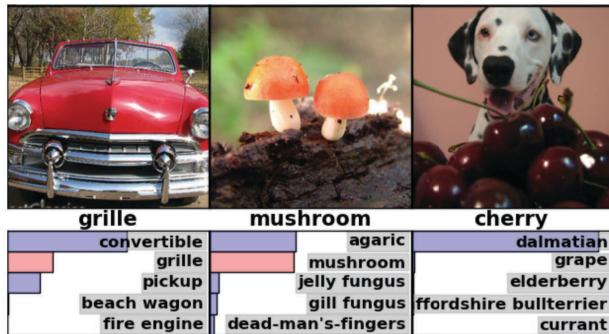
ConvNets^[1]



- ☺ 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)

[1] LeCun et al. 1998

ConvNets in computer vision



Image/object recognition

[Krizhevsky-Sutskever-Hinton'12]



Image inpainting

[Pathak-Efros-etal'16]



Image captioning

[Karpathy-Li'15]



Self-driving cars

⇒ Highly successful - ConvNets is the Swiss knife of Computer Vision!

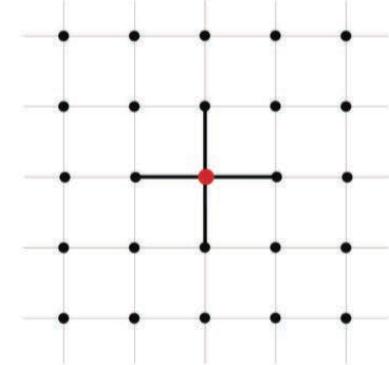
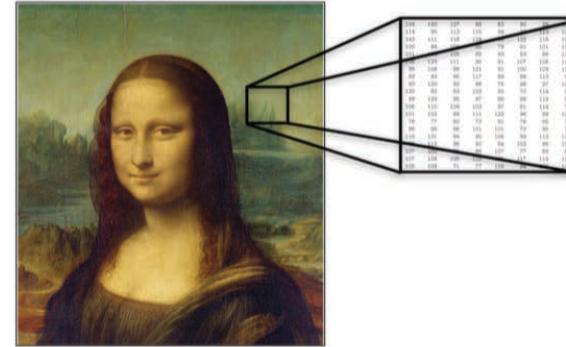


Outline

- **Convolutional Neural Networks**
 - *Architecture*
 - *Non-Euclidean Data*
- **Spectral Graph Theory**
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- **Spectral ConvNets**
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- **Spectral ConvNets on Multiple Graphs**
 - *Recommender Systems*
- **Conclusion**

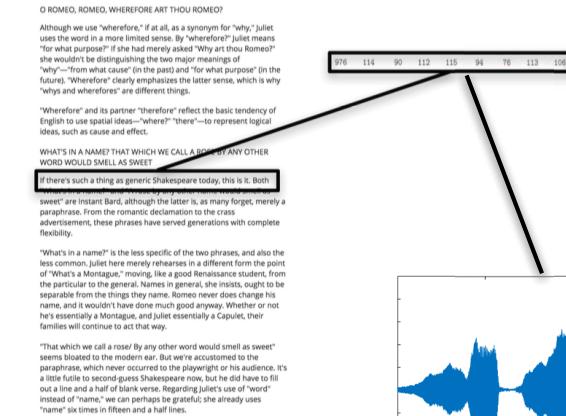
Data domain for ConvNets

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



2D grids

- Sentence, word, sound: 1D Euclidean domain



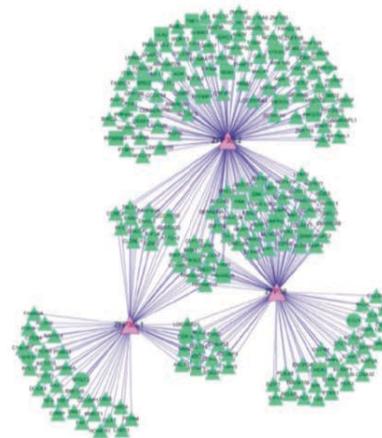
1D grid

- These domains have nice regular spatial structures.
⇒ All CNN operations are math well defined and fast (convolution, pooling).

Non-Euclidean data



Social networks

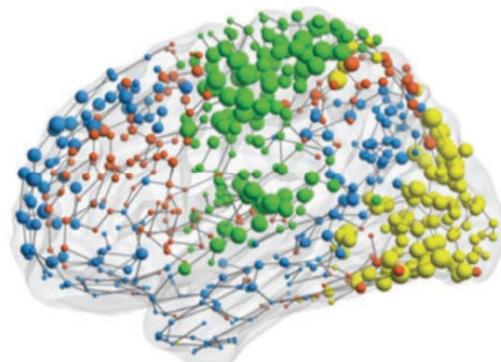


Regulatory networks

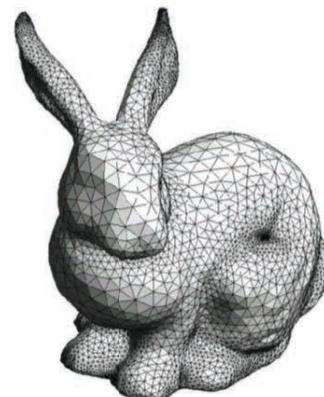
=



Graphs/
Networks



Functional networks



3D shapes

- Also NLP, physics, social science, communication networks, etc.

Challenges of graph deep learning

- Extend neural network techniques to **graph-structured data**.
- **Assumption:** Non-Euclidean data are locally stationary and manifest hierarchical structures.
- How to define **compositionality on graphs**? (convolution and pooling on graphs)
- How to make them **fast**? (linear complexity)

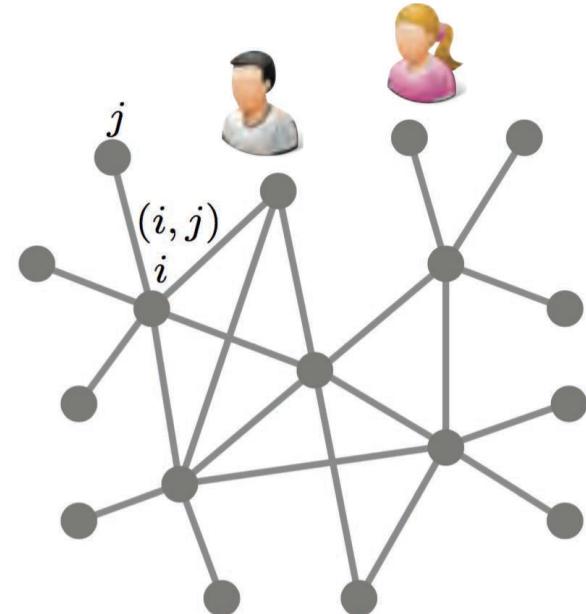
Outline

- Convolutional Neural Networks
 - *Architecture*
 - *Non-Euclidean Data*
- **Spectral Graph Theory**
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- Spectral ConvNets
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- Spectral ConvNets on Multiple Graphs
 - *Recommender Systems*
- Conclusion

Graphs[10]

- **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**

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



[10] Chung 1994

Graph Laplacian

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

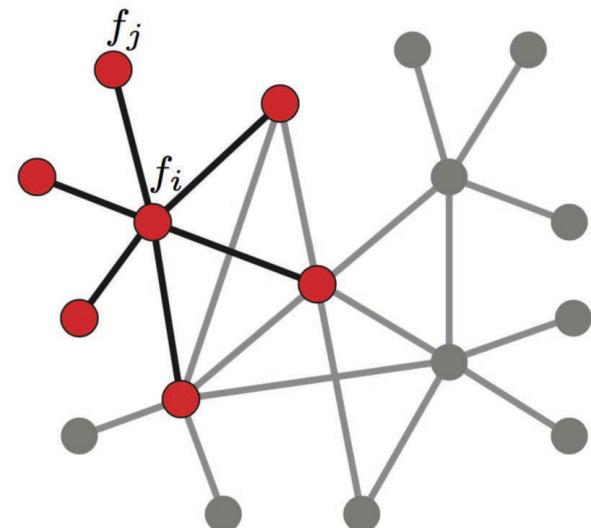
$$(\Delta f)_i = \frac{1}{b_i} \sum_{j:(i,j) \in \mathcal{E}} a_{ij} (f_i - f_j)$$

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

- **Unnormalized 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 $\mathbf{D} = \text{diag}(\sum_{j \neq i} a_{ij})$



Outline

- Convolutional Neural Networks
 - *Architecture*
 - *Non-Euclidean Data*
- **Spectral Graph Theory**
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- Spectral ConvNets
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- Spectral ConvNets on Multiple Graphs
 - *Recommender Systems*
- Conclusion

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 $\langle \phi_k, \phi_{k'} \rangle_{L^2(\mathcal{V})} = \delta_{kk'}$ (due to self-adjointness)
- Eigenvalues are non-negative $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ (due to positive-semidefiniteness)
- Eigendecomposition of a 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

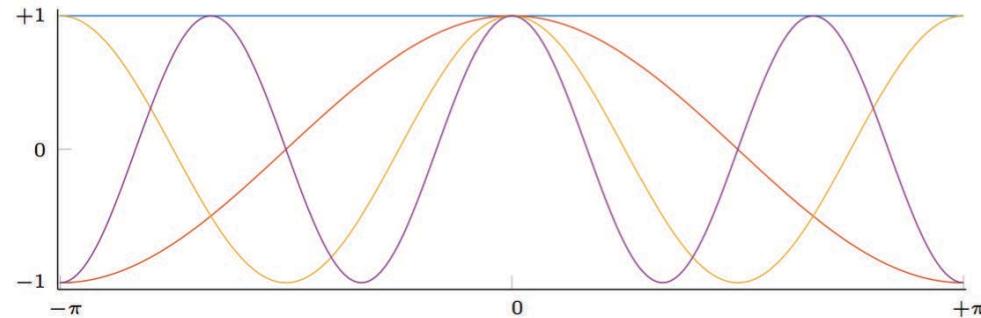
$$E_{\text{Dir}}(\mathbf{f}) = \mathbf{f}^\top \Delta \mathbf{f}$$

- **Solution:** first n Laplacian eigenvectors

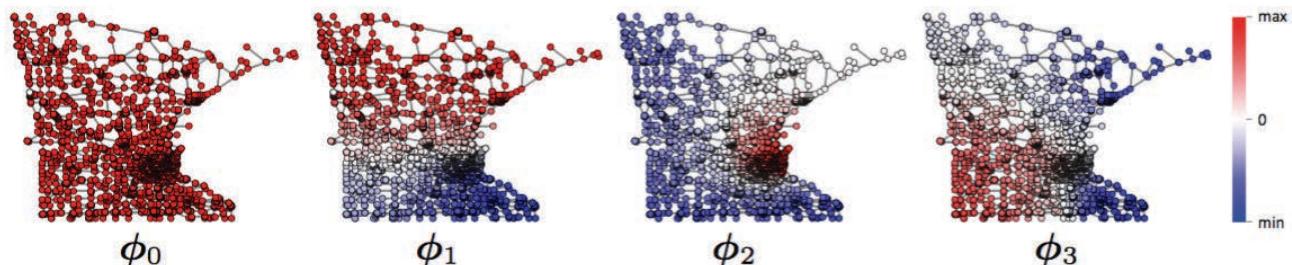
$$\begin{aligned} \min_{\Phi \in \mathbb{R}^{n \times n}} \underbrace{\text{trace}(\Phi^\top \Delta \Phi)}_{\|\Phi\|_{\mathcal{G}} \text{ Dirichlet norm}} \quad & \text{s.t.} \quad \Phi^\top \Phi = \mathbf{I} \end{aligned}$$

Laplacian eigenvectors

- Euclidean domain:



- Graph domain:



First Laplacian eigenvectors of a graph

Lap eigenvectors related to graph geometry
(s.a. communities, hubs, etc), spectral clustering^[10]

[10] Von Luxburg 2007

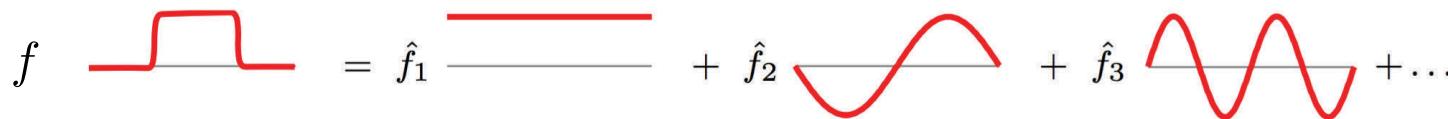
Outline

- Convolutional Neural Networks
 - *Architecture*
 - *Non-Euclidean Data*
- **Spectral Graph Theory**
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- Spectral ConvNets
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- Spectral ConvNets on Multiple Graphs
 - *Recommender Systems*
- Conclusion

Fourier analysis on Euclidean spaces

- 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:

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

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

Fourier analysis on graphs^[11]

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

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

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

$$\hat{f} = \Phi^\top f \quad \text{and} \quad f = \Phi \hat{f}$$

Fourier transform Inverse Fourier transform

- Graph Fourier basis ϕ_k = Laplacian eigenvectors :

$$\Delta \phi_k = \lambda_k \phi_k$$

$$\begin{aligned}\phi_k &= \text{graph Fourier mode} \\ \lambda_k &= \text{(square) frequency}\end{aligned}$$

[11] Hammond, Vandergheynst, Gribonval, 2011

Outline

- Convolutional Neural Networks
 - *Architecture*
 - *Non-Euclidean Data*
- **Spectral Graph Theory**
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- Spectral ConvNets
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- Spectral ConvNets on Multiple Graphs
 - *Recommender Systems*
- Conclusion

Convolution on Euclidean spaces

- 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 operator **commutes with Laplacian:** $(\Delta f) \star g = \Delta(f \star g)$
- Convolution theorem:** Convolution can be computed in the Fourier domain as

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

- Efficient computation** using FFT: $O(n \log n)$

Convolution on discrete Euclidean spaces

- 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)

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

Convolution on graphs^[11]

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

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

inverse Fourier transform

- 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)
- Commutes with Laplacian: $\mathbf{G} \Delta \mathbf{f} = \Delta \mathbf{G} \mathbf{f}$
- Filter coefficients depend on basis ϕ_1, \dots, ϕ_n
- Expensive computation (no FFT): $O(n^2)$

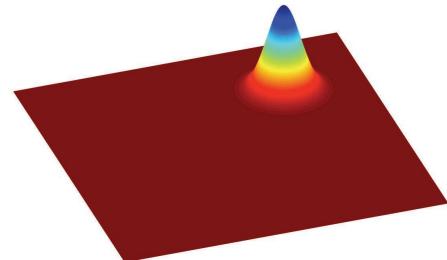
[11] Hammond, Vandergheynst, Gribonval, 2011

No shift invariance on graphs

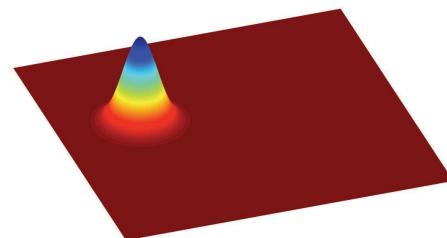
- A signal f on graph can be translated to vertex i :

$$T_i f = f \star \delta_i$$

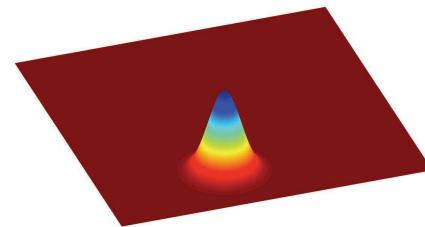
- Euclidean domain



(a) $T_s f$

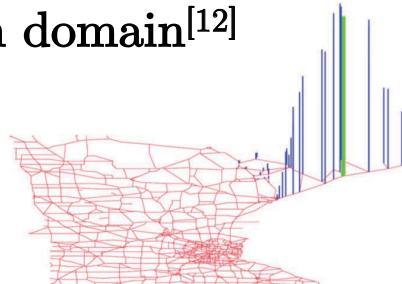


(b) $T_{s'} f$

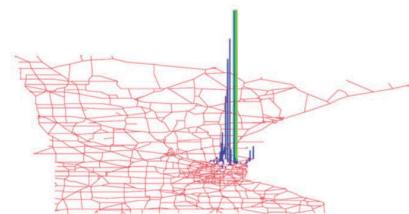


(c) $T_{s''} f$

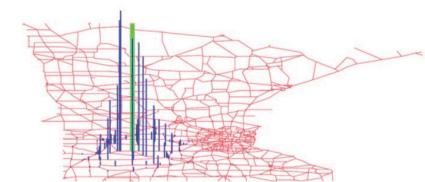
- Graph domain^[12]



(d) $T_i f$



(e) $T_{i'} f$



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

[12] Shuman et al. 2016

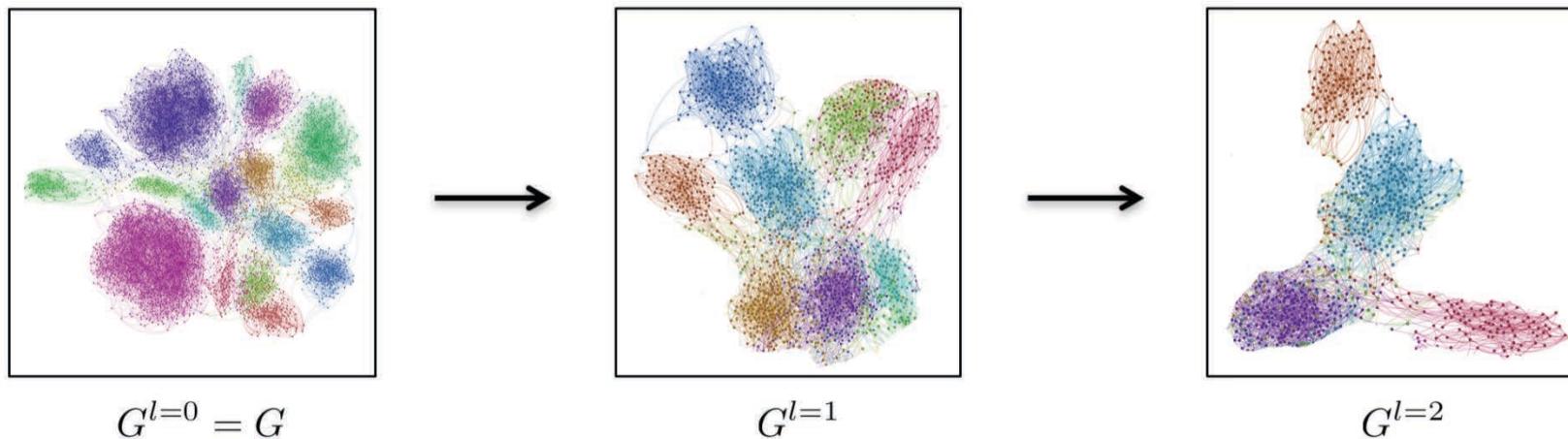
Outline

- Convolutional Neural Networks
 - *Architecture*
 - *Non-Euclidean Data*
- **Spectral Graph Theory**
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- Spectral ConvNets
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- Spectral ConvNets on Multiple Graphs
 - *Recommender Systems*
- Conclusion

Graph coarsening for graph pooling

- Goals:
 - Pool similar local features (**max pooling**).
 - 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**?

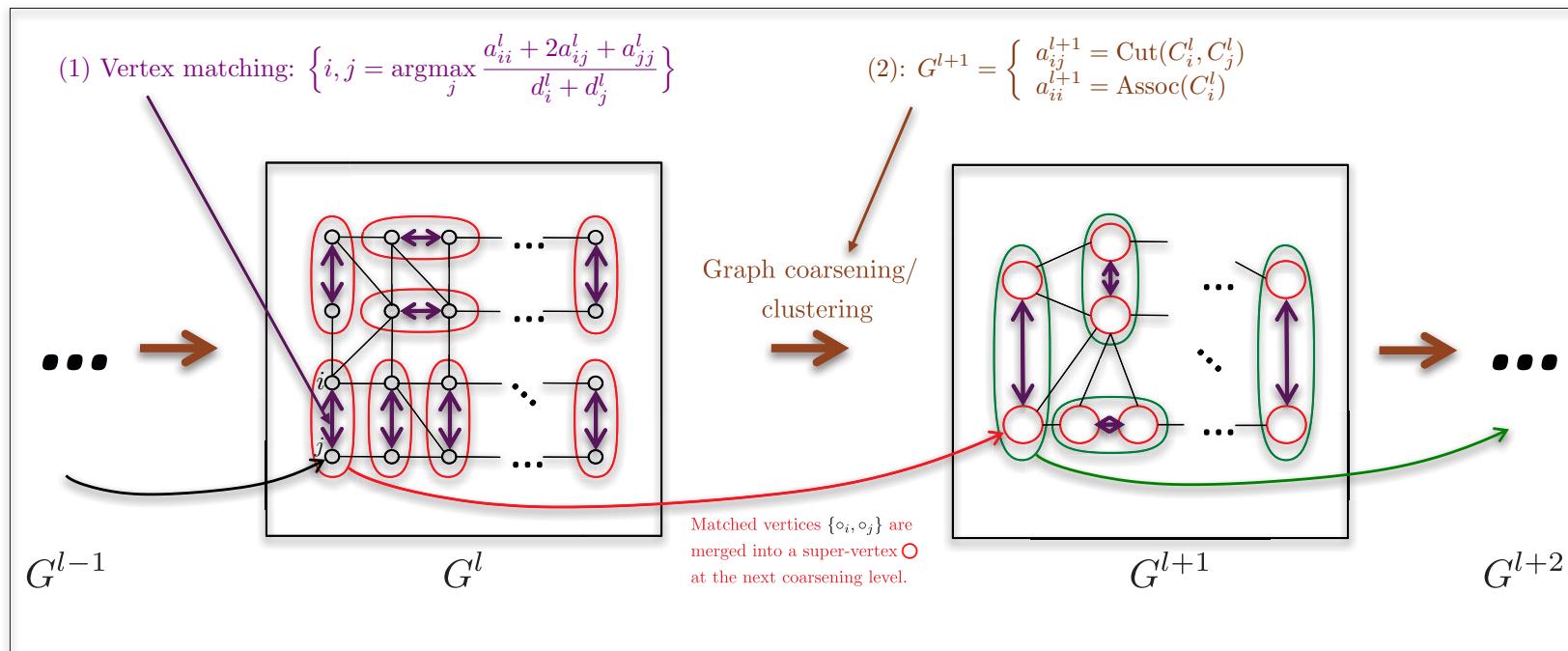
Graph coarsening = graph partitioning



- Graph partitioning decomposes G into smaller meaningful clusters.
- NP-hard \Rightarrow Approximation

Coarsening with Heavy-Edge Matching

- HEM proceeds by two successive steps, vertex matching and graph coarsening (that guarantees a local solution of Norm assoc):

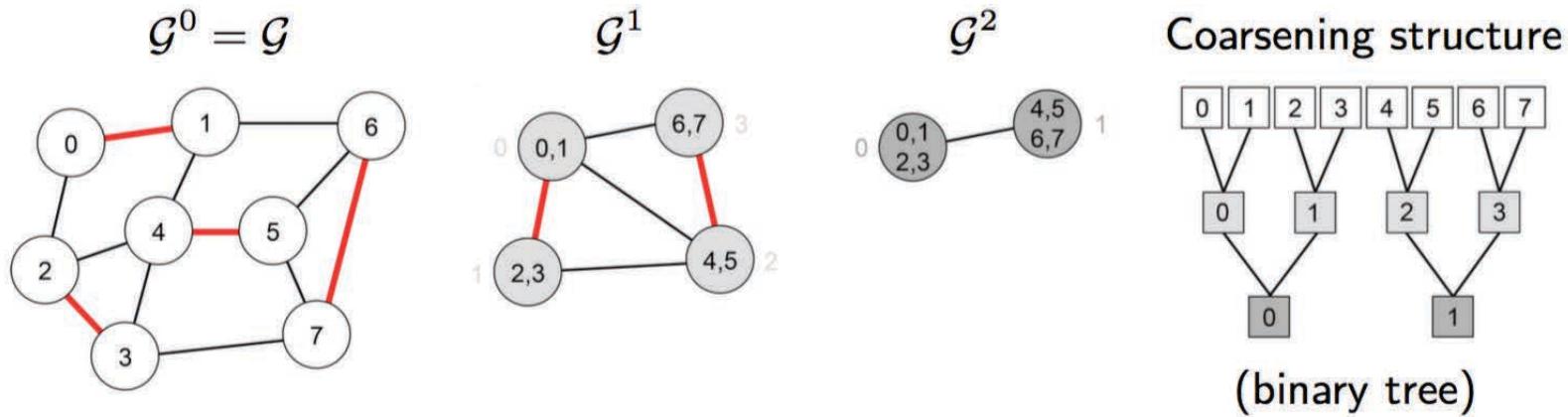


[14] Dhillon, Guan, Kulis 2007

[15] Karypis, Kumar 1995

Fast graph pooling^[18]

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



☺ As efficient as 1D-Euclidean grid pooling.

[18] Defferrard, Bresson, Vandergheynst 2016

Outline

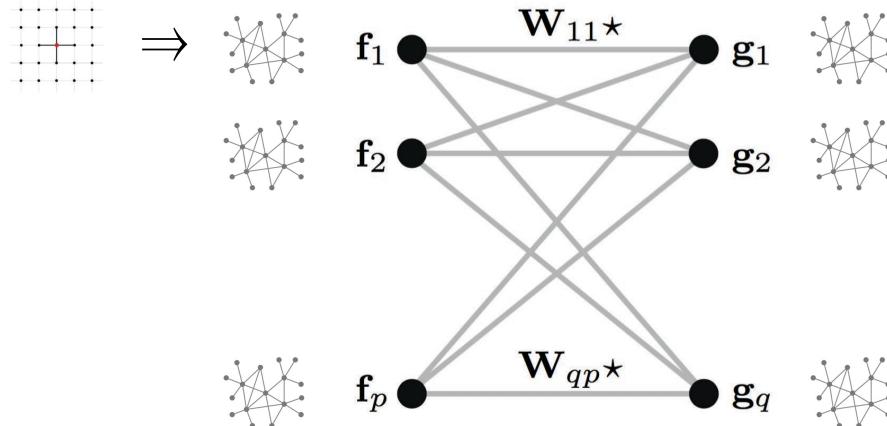
- Convolutional Neural Networks
 - *Architecture*
 - *Non-Euclidean Data*
- Spectral Graph Theory
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- **Spectral ConvNets**
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- Spectral ConvNets on Multiple Graphs
 - *Recommender Systems*
- Conclusion

Vanilla spectral graph ConvNets^[16]

- Graph convolutional layer

\mathbf{f}_l = l -th data feature on graphs, $\dim(\mathbf{f}_l) = n \times 1$

\mathbf{g}_l = l -th feature map, $\dim(\mathbf{g}_l) = n \times 1$



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

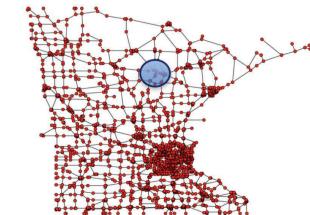
Activation, e.g. $\xi(x) = \max\{x, 0\}$ rectified linear unit (ReLU)

[16] Bruna, Zaremba, Szlam, LeCun 2014

Spectral graph convolution

- Convolutional layer in the spatial domain:

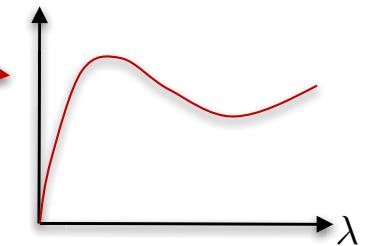
$$\mathbf{g}_l = \xi \left(\sum_{l'=1}^p \mathbf{W}_{l,l'} \star \mathbf{f}_{l'} \right),$$



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 graph ConvNets^[16]

- 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

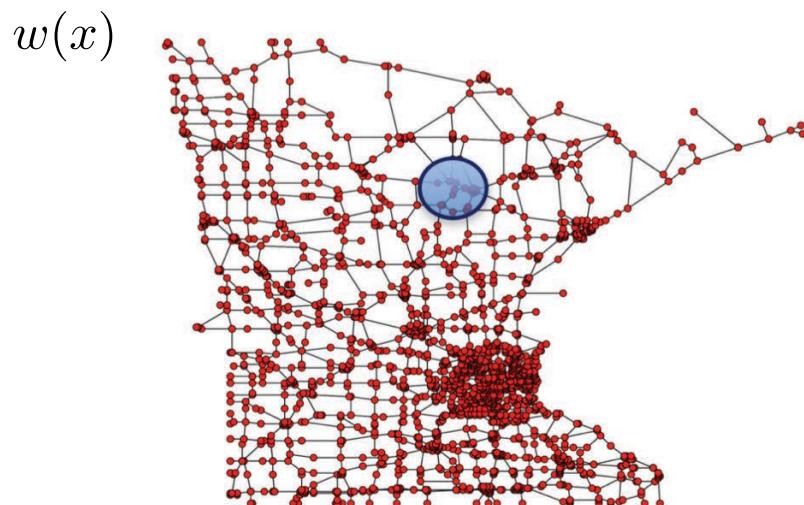
[16] Bruna, Zaremba, Szlam, LeCun 2014

Spatial localization and spectral smoothness^[17]

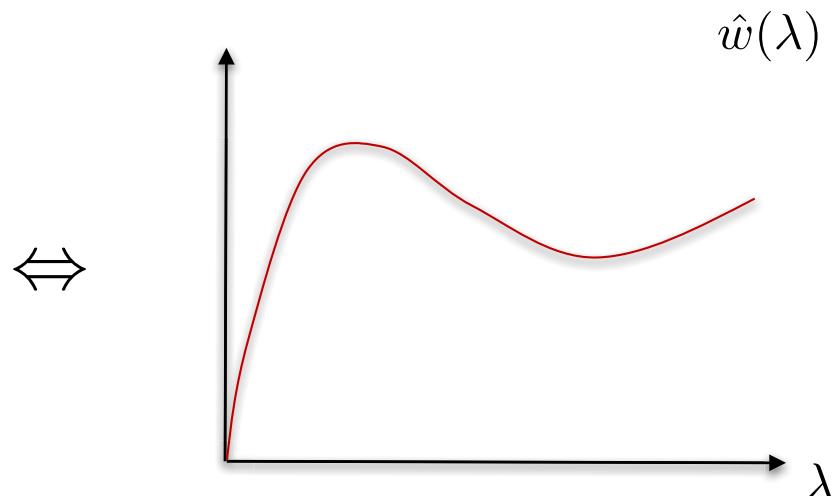
- In the Euclidean setting (by Parseval's identity)

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

⇒ Localization in space = smoothness in frequency domain



Spatial localization



Smooth spectral filter function

[17] Henaff, Bruna, LeCun 2015

Smooth parametric spectral filter

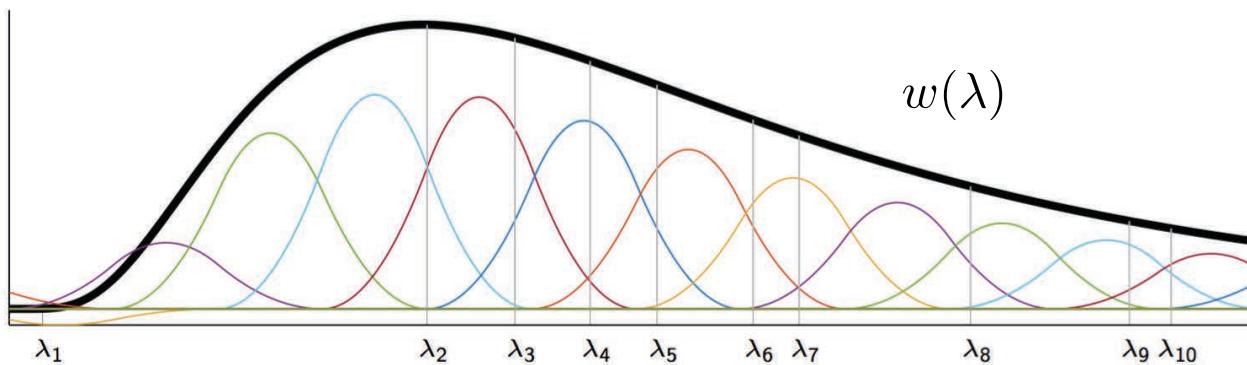
- Parametrize the smooth spectral filter function $w(\lambda)$ with a linear combination of smooth kernel functions $\beta_1(\lambda), \dots, \beta_r(\lambda)$, e.g. splines^[17]

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

\Downarrow

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

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



[17] (Litman, Bronstein, 2014); Henaff, Bruna, LeCun 2015

SplineNets^[17]

- 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 smooth spectral parametric coefficients $\mathbf{W}_{l,l'}^{(k)}$ to be learned 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

[17] Henaff, Bruna, LeCun 2015

Outline

- Convolutional Neural Networks
 - *Architecture*
 - *Non-Euclidean Data*
- Spectral Graph Theory
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- **Spectral ConvNets**
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- Spectral ConvNets on Multiple Graphs
 - *Recommender Systems*
- Conclusion

Spectral polynomial filters[18]

- Represent smooth spectral functions with polynomials of Laplacian eigenvalues

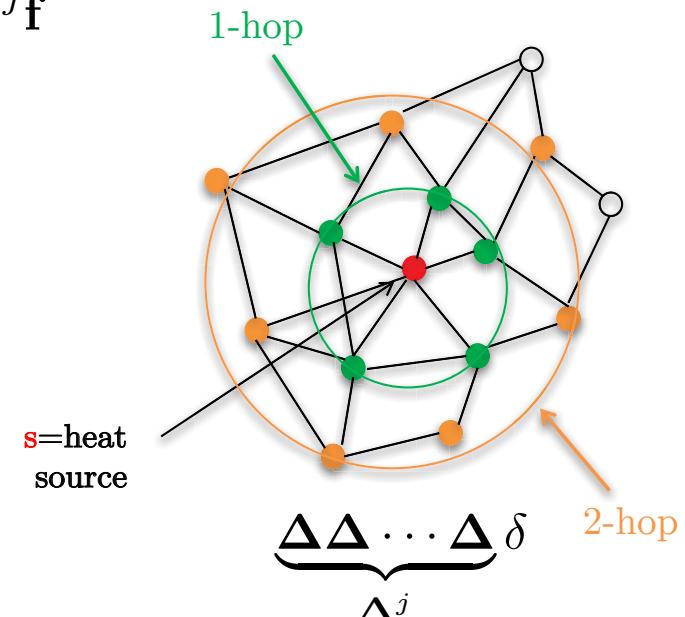
$$w_{\alpha}(\lambda) = \sum_{j=0}^r \alpha_j \lambda^j$$

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.



[18] Defferrard, Bresson, Vandergheynst 2016

Linear complexity

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

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

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

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

- Two important observations:
 1. *No* need to compute the eigendecomposition of the Laplacian (ϕ, Λ) .
 2. Observe that $\{\mathbf{X}_j\}$ are generated by **multiplication of a sparse matrix and a vector** \Rightarrow Complexity is $O(Er) = O(n)$ for sparse graphs.
- Graph convolutional layers are **GPU friendly**.

Spectral graph ConvNets with polynomial filters

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

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

[18] Defferrard, Bresson, Vandergheynst 2016

Chebyshev polynomials

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

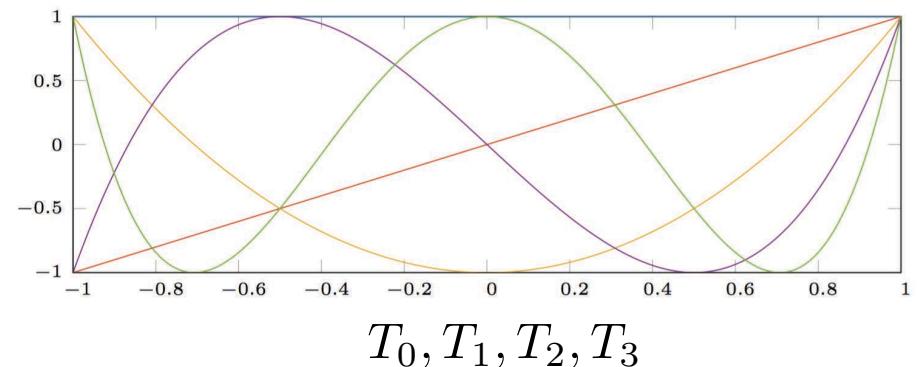
$$w_{\alpha}(\Delta)\mathbf{f} = \sum_{j=0}^r \alpha_j \Delta^j \mathbf{f}$$

- Graph convolution with (orthogonal) **Chebyshev polynomials**

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

- 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**



ChebNets^[18]

- Application of the filter with the scaled Laplacian $\tilde{\Delta} = 2\lambda_n^{-1}\Delta - \mathbf{I}$

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

with

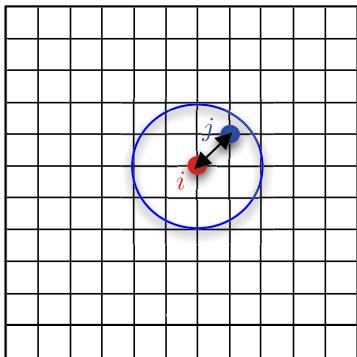
$$\begin{aligned}\mathbf{X}^{(j)} &= T_j(\tilde{\Delta})\mathbf{f} \\ &= 2\tilde{\Delta}\mathbf{X}^{(j-1)} - \mathbf{X}^{(j-2)}, \quad \mathbf{X}^{(0)} = \mathbf{f}, \quad \mathbf{X}^{(1)} = \tilde{\Delta}\mathbf{f}\end{aligned}$$

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

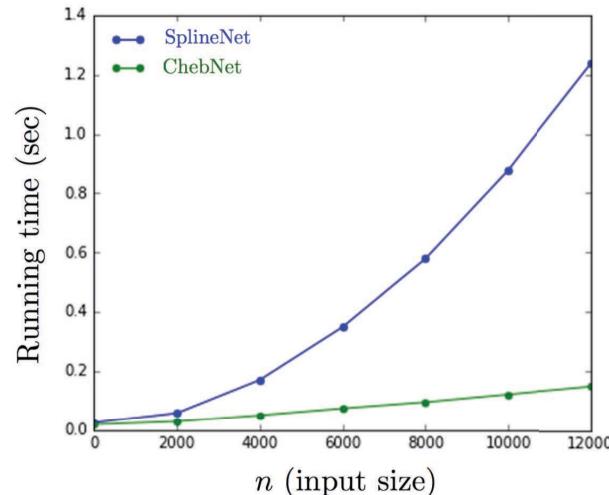
[18] Defferrard, Bresson, Vandergheynst 2016

Numerical experiments

Graph: a 8-NN graph of the Euclidean grid



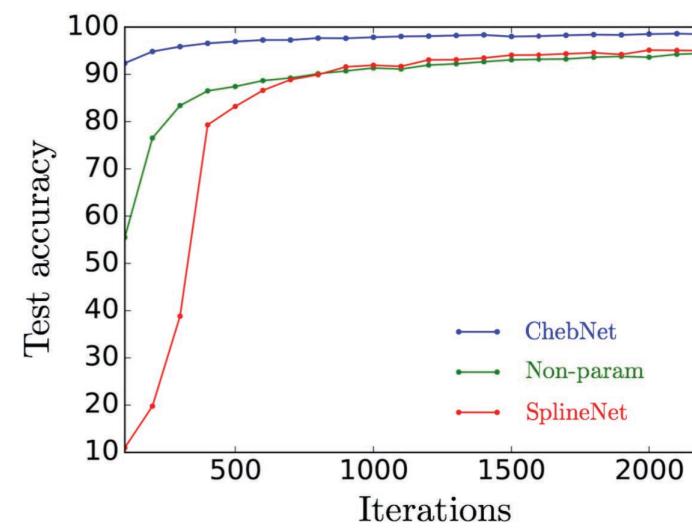
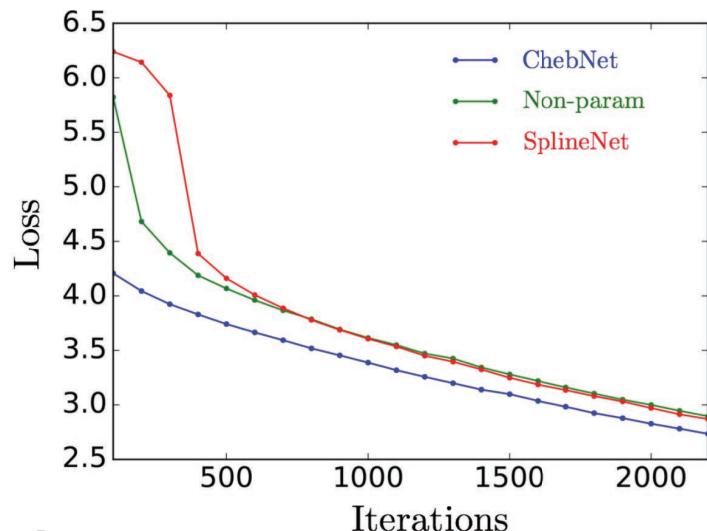
- Running time



- Accuracy

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

- Optimization



Outline

- Convolutional Neural Networks
 - *Architecture*
 - *Non-Euclidean Data*
- Spectral Graph Theory
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- **Spectral ConvNets**
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- Spectral ConvNets on Multiple Graphs
 - *Recommender Systems*
- Conclusion

Graph convolutional nets^[19]: simplified ChebNets

- Use Chebychev polynomials of degree $r=2$ and assume $\lambda_n \approx 2$

$$\begin{aligned} w_{\alpha}(\Delta)\mathbf{f} &= \alpha_0\mathbf{f} + \alpha_1(\Delta - \mathbf{I})\mathbf{f} \\ &= \alpha_0\mathbf{f} - \alpha_1\mathbf{D}^{-1/2}\mathbf{W}\mathbf{D}^{-1/2}\mathbf{f} \end{aligned}$$

- Further constrain $\alpha = \alpha_0 = -\alpha_1$ to obtain a single-parameter filter

$$w_{\alpha}(\Delta)\mathbf{f} = \alpha(\mathbf{I} + \mathbf{D}^{-1/2}\mathbf{W}\mathbf{D}^{-1/2})\mathbf{f}$$

- Caveat: The eigenvalues of $\mathbf{I} + \mathbf{D}^{-1/2}\mathbf{W}\mathbf{D}^{-1/2}$ are now in $[0,2]$
⇒ repeated application of the filter results in **numerical instability**
- Fix: Apply a **renormalization**

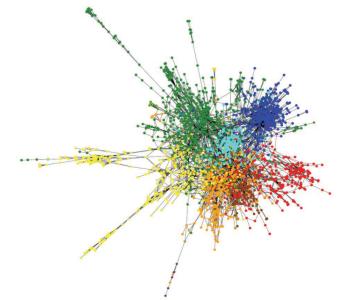
$$w_{\alpha}(\Delta)\mathbf{f} = \alpha\tilde{\mathbf{D}}^{-1/2}\tilde{\mathbf{W}}\tilde{\mathbf{D}}^{-1/2}\mathbf{f}$$

with $\tilde{\mathbf{W}} = \mathbf{W} + \mathbf{I}$ and $\tilde{\mathbf{D}} = \text{diag}(\sum_{j \neq i} \tilde{w}_{ij})$

[19] Kipf, Welling 2016

Example: citation networks

Method	Cora¹	PubMed²
Manifold Regularization ³	59.5%	70.7%
Semidefinite Embedding ⁴	59.0%	71.1%
Label Propagation ⁵	68.0%	63.0%
DeepWalk ⁶	67.2%	65.3%
Planetoid ⁷	75.7%	77.2%
Graph Convolutional Net⁸	81.59%	78.72%



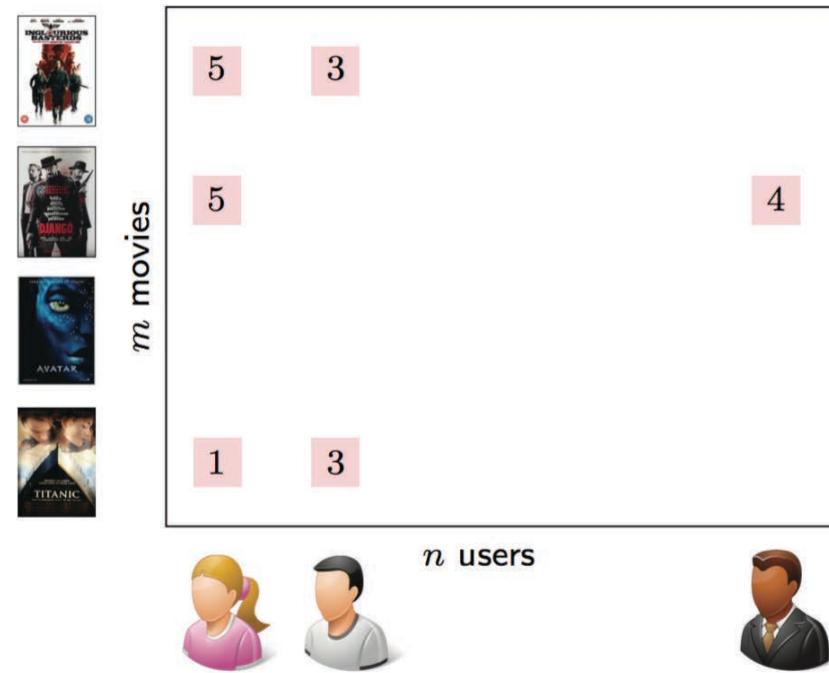
Classification accuracy of different methods on citation network datasets

Monti et al. 2016; data: ^{1,2}Sen et al. 2008; methods: ³Belkin et al. 2006; ⁴Weston et al. 2012; ⁵Zhu et al. 2003; ⁶Perozzi et al. 2014; ⁷Yang et al. 2016; ⁸Kipf, Welling 2016

Outline

- Convolutional Neural Networks
 - *Architecture*
 - *Non-Euclidean Data*
- Spectral Graph Theory
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- Spectral ConvNets
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- **Spectral ConvNets on Multiple Graphs**
 - *Recommender Systems*
- Conclusion

Matrix completion - Netflix challenge



$$\min_{\mathbf{X} \in \mathbb{R}^{m \times n}} \quad \text{rank}(\mathbf{X}) \quad \text{s.t.} \quad x_{ij} = a_{ij} \quad \forall ij \in \Omega$$

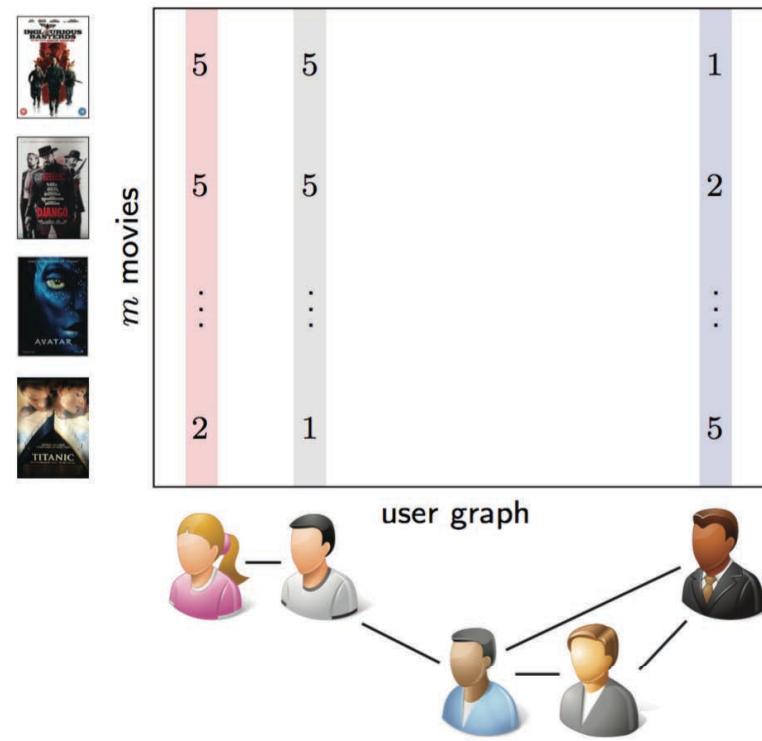
\Downarrow Convex relaxation^[21]

$$\min_{\mathbf{X} \in \mathbb{R}^{m \times n}} \quad \|\mathbf{X}\|_* \quad \text{s.t.} \quad x_{ij} = a_{ij} \quad \forall ij \in \Omega$$

$$\min_{\mathbf{X} \in \mathbb{R}^{m \times n}} \quad \|\mathbf{X}\|_* + \mu \|\mathbf{\Omega} \circ (\mathbf{X} - \mathbf{A})\|_F^2$$

[21] Candes et al. 2008

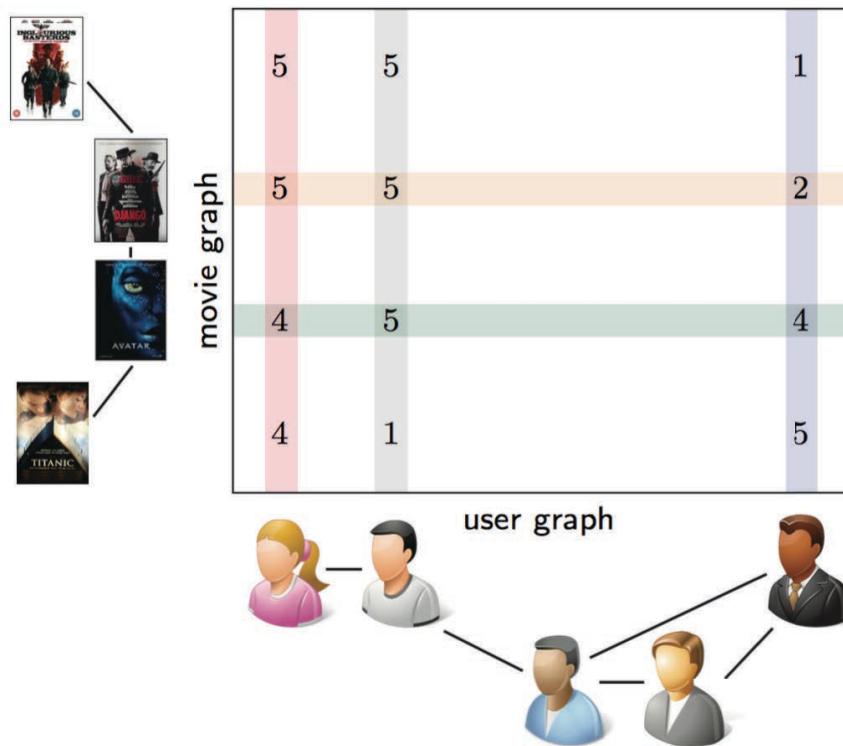
Matrix completion on single graph^[22]



$$\min_{\mathbf{X} \in \mathbb{R}^{m \times n}} \quad \|\mathbf{X}\|_* + \mu \|\boldsymbol{\Omega} \circ (\mathbf{X} - \mathbf{A})\|_F^2 + \mu_c \underbrace{\text{tr}(\mathbf{X} \boldsymbol{\Delta}_c \mathbf{X}^\top)}_{\|\mathbf{X}\|_{\mathcal{G}_c}^2}$$

[22] Kalofolias, Bresson, Bronstein, Vandergheynst, 2014

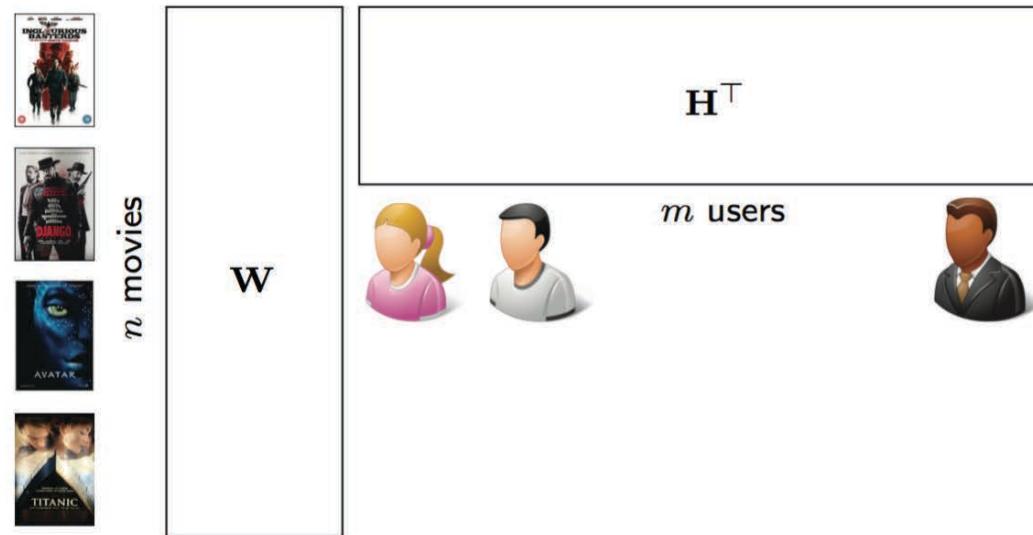
Matrix completion on multiple graphs^[22]



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

[22] Kalofolias, Bresson, Bronstein, Vandergheynst, 2014

Factorized matrix completion models^[23]

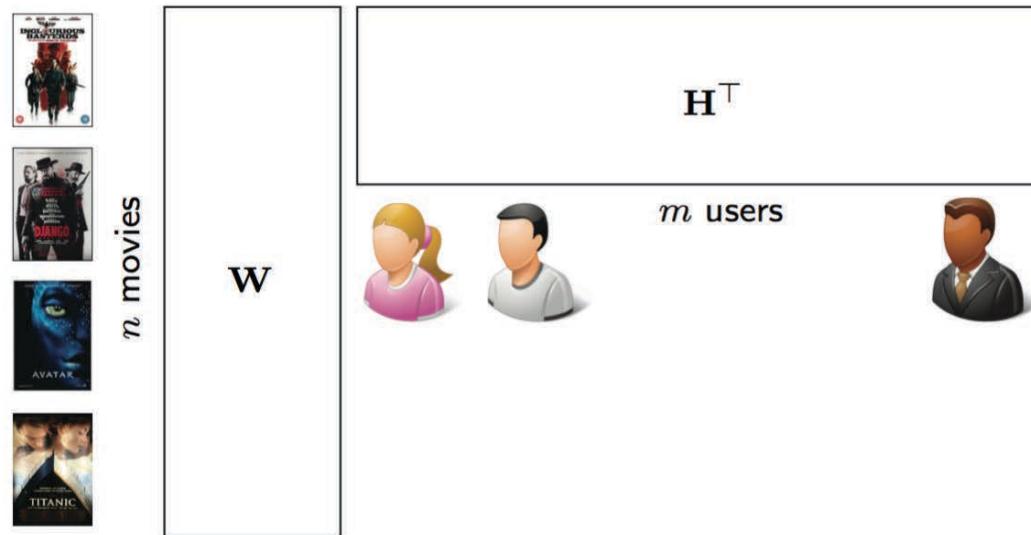


$$\min_{\substack{\mathbf{W} \in \mathbb{R}^{m \times s} \\ \mathbf{H} \in \mathbb{R}^{n \times s}}} \mu \|\Omega \circ (\mathbf{X} - \mathbf{A})\|_F^2 + \mu_c \|\mathbf{W}\|_F^2 + \mu_r \|\mathbf{H}\|_F^2$$

- ☺ $O(n+m)$ variables instead of $O(nm)$
- ☺ $\text{rank}(\mathbf{X}) = \text{rank}(\mathbf{W}\mathbf{H}^T) \leq s$ by construction

[23] Srebro, Rennie, Jaakkola 2004

Factorized matrix completion on graphs^[24]



$$\min_{\substack{\mathbf{W} \in \mathbb{R}^{m \times s} \\ \mathbf{H} \in \mathbb{R}^{n \times s}}} \mu \|\boldsymbol{\Omega} \circ (\mathbf{X} - \mathbf{A})\|_F^2 + \mu_c \text{tr}(\mathbf{H}^\top \boldsymbol{\Delta}_c \mathbf{H}) + \mu_r \text{tr}(\mathbf{W}^\top \boldsymbol{\Delta}_r \mathbf{W})$$

- ☺ $O(n+m)$ variables instead of $O(nm)$
- ☺ $\text{rank}(\mathbf{X}) = \text{rank}(\mathbf{W}\mathbf{H}^\top) \leq s$ by construction

[24] Rao et al. 2015

1D Fourier transform

$$\begin{matrix} \text{Heatmap} \\ (\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_n) \end{matrix} = \begin{matrix} \Phi^\top \\ \text{Matrix} \end{matrix} \times \begin{matrix} \text{Mona Lisa} \\ \mathbf{x} \end{matrix}$$

Column-wise transform

2D Fourier transform

$$\hat{\mathbf{X}} = \Phi^T \times \mathbf{X} \times \Phi$$

The diagram illustrates the 2D Fourier transform process. On the left is a heatmap labeled $\hat{\mathbf{X}}$, representing the Fourier transform of the image. To its right is an equals sign. Following the equals sign are three components: a vertical matrix labeled Φ^T , a grayscale image of the Mona Lisa labeled \mathbf{X} , and a horizontal matrix labeled Φ . Between the Φ^T matrix and the image \mathbf{X} is a multiplication symbol (\times). Between the image \mathbf{X} and the Φ matrix is another multiplication symbol (\times). This visualizes the formula $\hat{\mathbf{X}} = \Phi^T \times \mathbf{X} \times \Phi$.

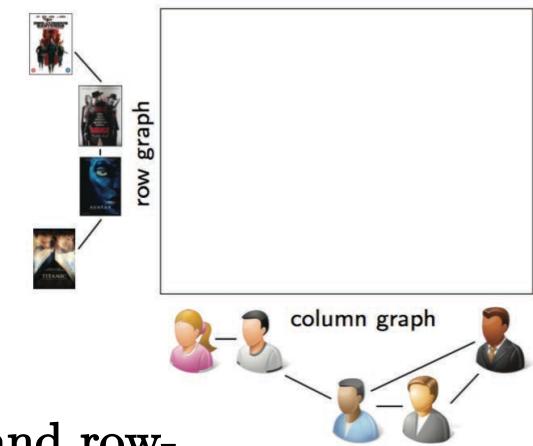
Column-wise transform + Row-wise transform = 2D transform

From 1D signal processing to 2D image processing

Multi-graph Fourier transform and convolution^[25]

- Multi-graph Fourier transform

$$\hat{\mathbf{X}} = \Phi_r^\top \mathbf{X} \Phi_c$$



where Φ_c , Φ_r are the eigenvectors of the column- and row-graph Laplacians Δ_c , Δ_r , respectively.

- Multi-graph spectral convolution

$$\begin{aligned}\mathbf{X} \star \mathbf{Y} &= \Phi_r ((\Phi_r^\top \mathbf{X} \Phi_c) \circ (\Phi_r^\top \mathbf{Y} \Phi_c)) \Phi_c^\top \\ &= \Phi_r (\hat{\mathbf{X}} \circ \hat{\mathbf{Y}}) \Phi_c^\top\end{aligned}$$

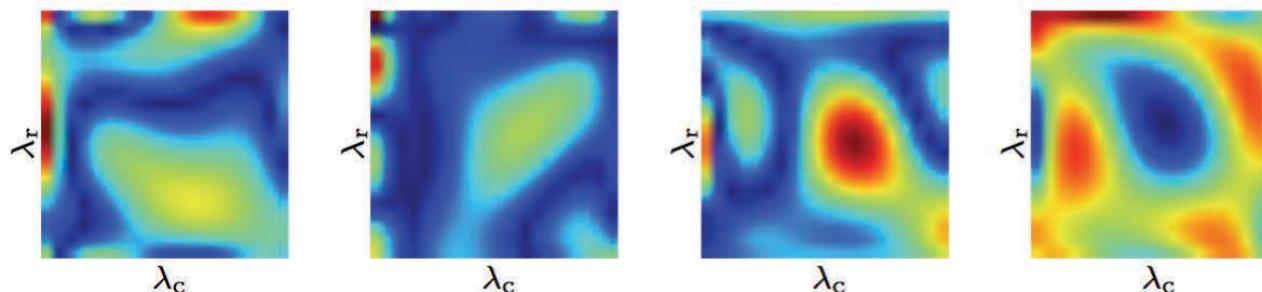
[25] Monti, Bresson, Bronstein 2017

Multi-graph spectral filters

- Parametrize the multi-graph filter as a smooth function of column- and row frequencies w/ e.g. bivariate Chebyshev polynomial

$$w_{\Theta}(\tilde{\lambda}_c, \tilde{\lambda}_r) = \sum_{j,j'=1}^r \theta_{jj'} T_j(\tilde{\lambda}_c) T_{j'}(\tilde{\lambda}_r)$$

where $\Theta = (\theta_{jj'})$ is matrix of coefficients and $-1 \leq \tilde{\lambda}_c, \tilde{\lambda}_r \leq 1$



Examples of multi-graph spectral filters (shown is $|\tau(\tilde{\lambda}_c, \tilde{\lambda}_r)|$)

- Multi-graph convolutional layer: Apply filter to matrix X

$$\mathbf{Y} = \sum_{j,j'=1}^r \theta_{jj'} T_j(\tilde{\Delta}_c) \mathbf{X} T_{j'}(\tilde{\Delta}_r)$$

where the scaled Laplacians $\tilde{\Delta}_c = 2\lambda_{c,n}^{-1} \Delta_c - \mathbf{I}$ and $\tilde{\Delta}_r = 2\lambda_{r,m}^{-1} \Delta_r - \mathbf{I}$

Multi-graph ConvNets^[25]

- Multi-graph spectral filters

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

applied to p input channels ($m \times n$ matrices $\mathbf{Y}_1^{(k-1)}, \dots, \mathbf{Y}_p^{(k-1)}$) and producing q output channels $\mathbf{Y}_1^{(k)}, \dots, \mathbf{Y}_q^{(k)}$

- ☺ Filters have guaranteed r -hops support on both graphs
- ☺ $O(1)$ parameters per layer
- ☹ $O(nm)$ computational complexity (assuming sparse graphs)

[25] Monti, Bresson, Bronstein 2017

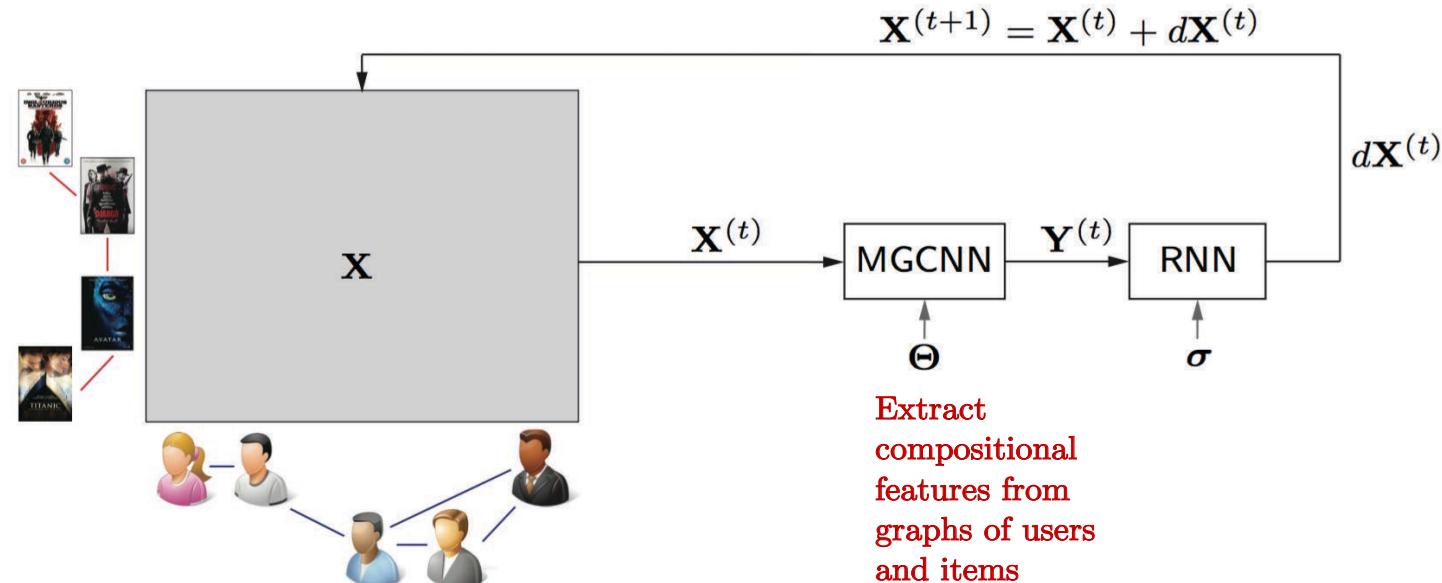
Factorized multi-graph ConvNets

- Two spectral convolutional layers applied to each of the factors \mathbf{W} , \mathbf{H}

$$\begin{aligned}\mathbf{u}_l &= \xi \left(\sum_{l'=1}^p \sum_{j=0}^r \theta_{\text{r},jll'} T_j(\tilde{\Delta}_{\text{r}}) \mathbf{w}_{l'} \right) & l &= 1, \dots, q \\ \mathbf{v}_l &= \xi \left(\sum_{l'=1}^p \sum_{j'=0}^r \theta_{\text{c},j' ll'} T_{j'}(\tilde{\Delta}_{\text{c}}) \mathbf{h}_{l'} \right) & l' &= 1, \dots, p\end{aligned}$$

- ☺ Filters have guaranteed r -hops support on both graphs
- ☺ $O(1)$ parameters per layer
- ☺ $O(n+m)$ computational complexity (assuming sparse graphs)
- ☹ May loose some coupling information

Matrix completion with Recurrent Multi-Graph ConvNets^[25]



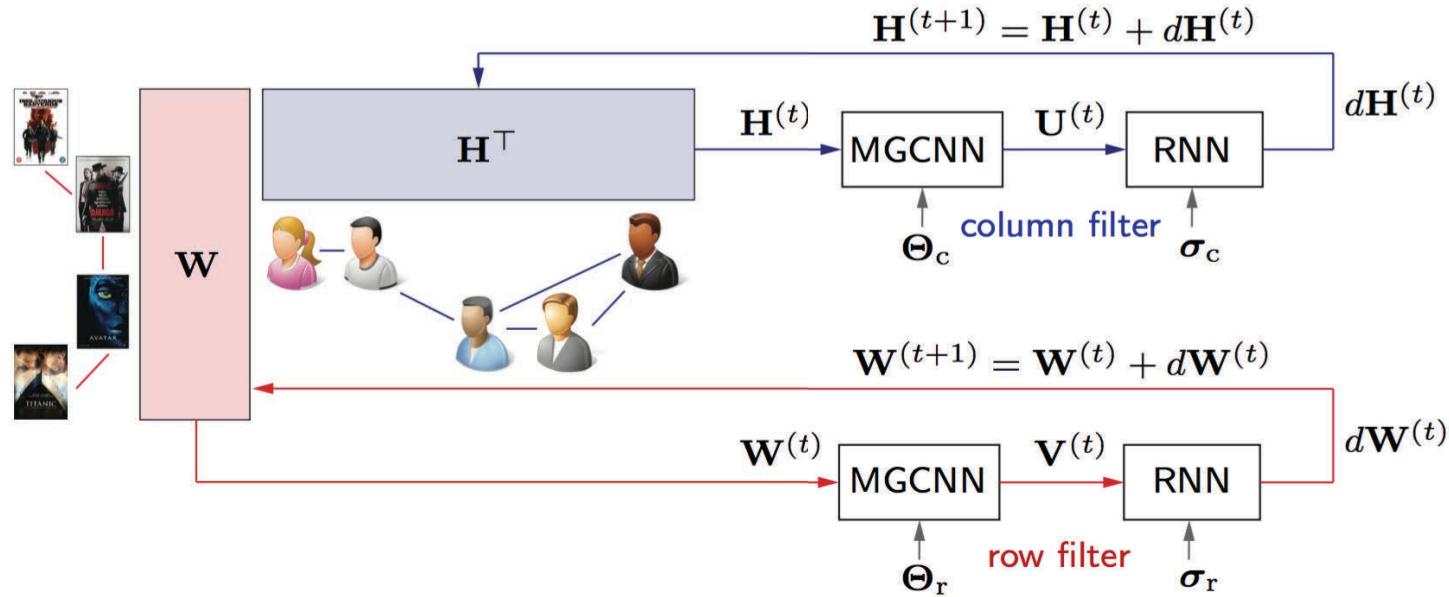
Recurrent multigraph CNN (RMCNN) architecture
for matrix completion

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

Complexity: $O(nm)$

[25] Monti, Bresson, Bronstein 2017

Factorized version



Separable recurrent multigraph CNN (sRMCNN) architecture
for matrix completion in factorized form

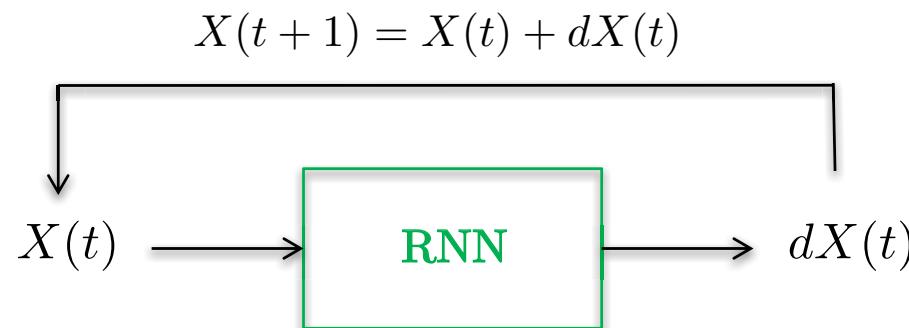
$$\min_{\theta_r, \theta_r, \sigma_r, \sigma_c} \|\mathbf{W}_{\theta_r, \sigma_r}^{(T)}\|_{\mathcal{G}_r}^2 + \|\mathbf{H}_{\theta_c, \sigma_c}^{(T)}\|_{\mathcal{G}_c}^2 + \frac{\mu}{2} \|\Omega \circ (\mathbf{W}_{\theta_r, \sigma_r}^{(T)} (\mathbf{H}_{\theta_c, \sigma_c}^{(T)})^\top - \mathbf{Y})\|_F^2$$

Complexity: $O(n+m)$

[25] Monti, Bresson, Bronstein 2017

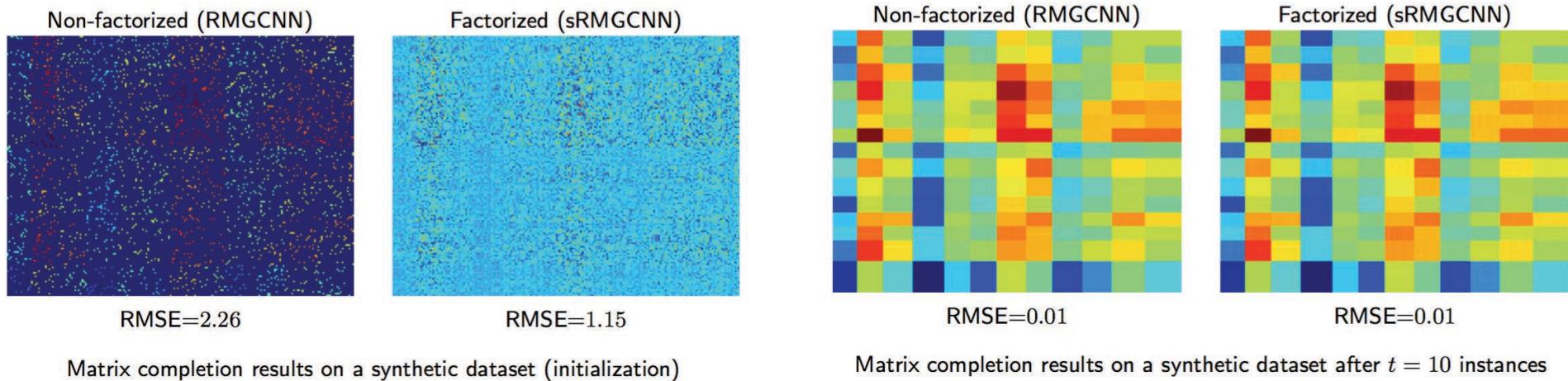
Incremental updates with RNNs

- The RNN/LSTM^[26] cast the completion problem as **non-linear diffusion** process.
 - Replace multi-layer CNNs by simple RNN + diffusion
⇒ **much less parameters to learn** (reduce overfitting).
⇒ **more favorable for highly sparse entries** (better info diffusion).
 - Fix #parameters of the model **independently of entry sparsity**.
 - Allows to **learn temporal dynamics of entries** if such information is available.



[26] Hochreiter, Schmidhuber 1997

Matrix completion methods comparison: synthetic



Method	#Params	Complexity	RMSE
GMC ¹	mn	mn	0.3693
GRALS ²	$m + n$	$m + n$	0.0114
sRMGCNN³	1	$m + n$	0.0106
RMGCNN³	1	mn	0.0053

Comparison of geometric matrix completion methods on synthetic data
using both user+movie graphs

Methods: ¹Kalofolias et al. 2014; ²Rao et al. 2015; ³Monti, Bresson, Bronstein 2017

Matrix completion methods comparison: real

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	—	—	—

Accuracy (RMS error) of matrix completion methods on real data

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

Outline

- Convolutional Neural Networks
 - *Architecture*
 - *Non-Euclidean Data*
- Spectral Graph Theory
 - *Graphs and operators*
 - *Spectral decomposition*
 - *Fourier analysis*
 - *Convolution*
 - *Coarsening*
- Spectral ConvNets
 - *SplineNets*
 - *ChebNets*
 - *GraphConvNets*
- Spectral ConvNets on Multiple Graphs
 - *Recommender Systems*
- Conclusion

Conclusion

- Contributions:
 - Generalization of ConvNets to data on (fixed) graphs
 - Linear complexity for sparse graphs
 - GPU implementation
 - Multiple graphs
- Several potential applications in network sciences.



Questions?