

Graph Representation Learning, Deep Generative Models On Graphs & Multiresolution Machine Learning

Truong Son Hy

Department of Computer Science
The University of Chicago



① Graph representation learning

- Message passing neural networks
- Permutation equivariance
- Covariant compositional networks

② Multiresolution matrix factorization

- Reinforcement Learning & Stiefel manifold optimization
- Graph wavelets
- Wavelet neural networks

③ Deep generative models on graphs

- Variational Autoencoder (VAE)
- Equivariant graph/molecule generation
- Multiresolution graph VAE

Motivation for graph learning

The Anatomy of a Large-Scale Hypertextual Web Search Engine

Sergey Brin and Lawrence Page

Computer Science Department
Stanford University, Stanford, CA 94305, USA
sergey@cs.stanford.edu and page@cs.stanford.edu

Abstract

In this paper, we present Google, a prototype of a large-scale search engine which makes heavy use of the structure present in hypertext. Google is designed to crawl and index the Web efficiently and produce much more satisfying search results than existing systems. The prototype, with a full-text and hypertext database of at least 24 million pages is available at <http://google.stanford.edu/>. To engineer a search engine is a challenging task. Search engines index tens to hundreds of millions of pages every day. Despite the importance of large-scale search engines on the web, very little academic research has been done on them. Furthermore, due to rapid advances in technology and software development, creating a search engine today is very different from three years ago. This paper provides an in-depth description of our large-scale search engine – the first such detailed public description we know of to date. Apart from the problems of scaling, traditional search techniques to date of this magnitude, there are new technical challenges involved with using the additional information present in hypertext to produce better search results. This paper addresses this question of how to build a practical large-scale system which can exploit the additional information present in hypertext. Also we look at the problem of how to effectively deal with uncontrolled hypertext collections where anyone can publish anything they want.

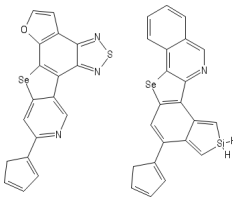
Keywords

World Wide Web, Search Engines, Information Retrieval, PageRank, Google

1. Introduction

(Note: There are two versions of this paper – a longer full version and a shorter printed version. The full version is available on the web and the conference CD-ROM.)

The web creates new challenges for information retrieval. The amount of information on the web is growing rapidly, as well as the number of users who are interested in the web. People are likely to surf the web using a link graph, often entering only high quality human maintained indices such as Yahoo! or with search engines. Human maintained lists cover popular topics effectively but are subjective, expensive to build and maintain, slow to improve, and cannot cover all current topics. Automated search engines that only on keyword matching usually return too many low quality matches. To make matters worse, some advertisers attempt to gain people's attention by linking keywords to unrelated automated search engines. We have built a large-scale search engine which addresses many of the problems of existing systems. It makes especially heavy use of the additional structure present in hypertext to provide much higher quality search results. We chose our system name, Google, because it is a common spelling of googol, or 10^{100} and fits well with our goal of building very large-scale search



(a) Citation network

(b) Molecules

(c) Knowledge graph

Thomas Jefferson <

3rd U.S. President



Thomas Jefferson was an American Founding Father who was the principal author of the Declaration of Independence and later served as the third President of the United States from 1801 to 1809. Previously, he had been elected the second Vice President of the United States, serving under John Adams from 1797 to 1801. [Wikipedia](#)

Born: April 13, 1743, [Shadwell, VA](#)

Died: July 4, 1826, [Monticello, VA](#)

Presidential term: March 4, 1801 – March 4, 1809

Spouse: [Martha Jefferson](#) (m. 1772–1782)

Children: [Martha Jefferson Randolph](#), [Madison Hemings](#). [MORE](#)

Vice presidents: [Aaron Burr](#) (1801–1805), [George Clinton](#) (1805–1809)

People also search for

View 15+ more



[John Adams](#)



[George Washington](#)



[James Madison](#)



[Benjamin Franklin](#)



[Abraham Lincoln](#)

[Feedback](#)

A form of Graph Neural Networks (GNNs)

DFT = Density Functional Theory

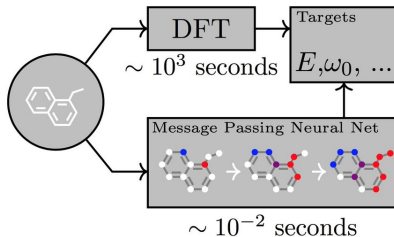


Figure 1. A Message Passing Neural Network predicts quantum properties of an organic molecule by modeling a computationally expensive DFT calculation.

Gilmer et al., *Neural Message Passing for Quantum Chemistry*, ICML 2017

Message Passing Scheme (1)

Given an input graph / network $G = (V, E)$:

- 1 Initially, each vertex v of the graph is associated with a feature representation ℓ_v (label) or f_v^0 . This feature representation can also be called as a *message*.

Message Passing Scheme (1)

Given an input graph / network $G = (V, E)$:

- 1 Initially, each vertex v of the graph is associated with a feature representation ℓ_v (label) or f_v^0 . This feature representation can also be called as a *message*.
- 2 Iteratively, at iteration t , each vertex collects / aggregates all messages of the previous iteration $\{f_{v_1}^{t-1}, \dots, f_{v_k}^{t-1}\}$ from other vertices in its neighborhood $\mathcal{N}(v) = \{v_1, \dots, v_k\}$, and then produces a new message f_v^t via some *hashing function* $\psi(\cdot)$.

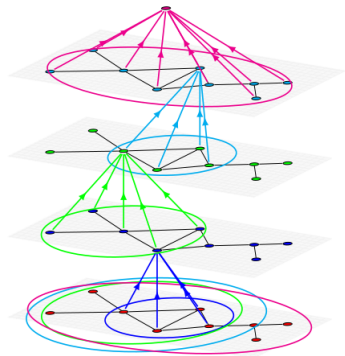
Message Passing Scheme (1)

Given an input graph / network $G = (V, E)$:

- 1 Initially, each vertex v of the graph is associated with a feature representation ℓ_v (label) or f_v^0 . This feature representation can also be called as a *message*.
- 2 Iteratively, at iteration t , each vertex collects / aggregates all messages of the previous iteration $\{f_{v_1}^{t-1}, \dots, f_{v_k}^{t-1}\}$ from other vertices in its neighborhood $\mathcal{N}(v) = \{v_1, \dots, v_k\}$, and then produces a new message f_v^t via some *hashing function* $\psi(\cdot)$.
- 3 The graph representation $\phi(G)$ is obtained by aggregating all messages in the last iteration of every vertex. $\phi(G)$ is then used for downstream application.

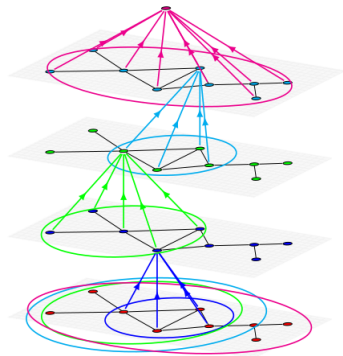
Message Passing Scheme (2)

```
1: for  $v \in V$  do  
2:    $f_v^0 \leftarrow \ell_v$   
3: end for  
4: for  $t = 1 \rightarrow T$  do  
5:   for  $v \in V$  do  
6:      $f_v^t \leftarrow \psi(\{f_i^{t-1}\}_{i \in \mathcal{N}(v)})$   
7:   end for  
8: end for  
9:  $\phi(G) \leftarrow \psi(\{f_v^T\}_{v \in V})$ 
```



Message Passing Scheme (2)

```
1: for  $v \in V$  do  
2:    $f_v^0 \leftarrow \ell_v$   
3: end for  
4: for  $t = 1 \rightarrow T$  do  
5:   for  $v \in V$  do  
6:      $f_v^t \leftarrow \psi(\{f_i^{t-1}\}_{i \in \mathcal{N}(v)})$   
7:   end for  
8: end for  
9:  $\phi(G) \leftarrow \psi(\{f_v^T\}_{v \in V})$ 
```



Note

This procedure is used in Weisfeiler–Lehman **graph isomorphism** test (NP-complete problem).

Message Passing Scheme (3)

With learnable parameters:

```
1: for  $v \in V$  do  
2:    $f_v^0 \leftarrow \ell_v$   
3: end for  
4: for  $t = 1 \rightarrow T$  do  
5:   for  $v \in V$  do  
6:      $f_v^t \leftarrow \psi(\{f_i^{t-1}\}_{i \in \mathcal{N}(v)}; W^t)$   
7:   end for  
8: end for  
9:  $\phi(G) \leftarrow \psi(\{f_v^T\}_{v \in V}; W^{T+1})$ 
```

Message Passing Scheme (3)

With learnable parameters:

```
1: for  $v \in V$  do
2:    $f_v^0 \leftarrow \ell_v$ 
3: end for
4: for  $t = 1 \rightarrow T$  do
5:   for  $v \in V$  do
6:      $f_v^t \leftarrow \psi(\{f_i^{t-1}\}_{i \in \mathcal{N}(v)}; W^t)$ 
7:   end for
8: end for
9:  $\phi(G) \leftarrow \psi(\{f_v^T\}_{v \in V}; W^{T+1})$ 
```

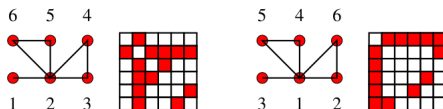
Given a graph properties $y(G) \in \mathbb{R}^d$ to regress, we have the optimization:

$$\min_{\{W^t\}_{t=1}^{T+1}} \|y(G) - \phi(G)\|_2^2$$

The gradient with respect to $\{W^t\}_{t=1}^{T+1}$ can be computed via Back-propagation.

Invariance

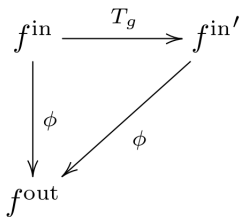
We renumber the vertices by a permutation $\sigma : \{1, 2, \dots, 6\} \mapsto \{1, 2, \dots, 6\}$. The adjacency matrices of G (left) and G' (right) are different, but topologically they represent the same graph:



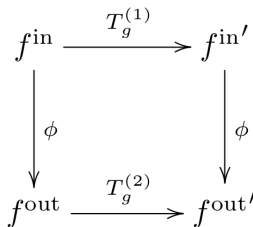
Therefore, ϕ must be **invariant** wrt permutation, i.e. $\phi(G) = \phi(G')$.

Invariance vs. Equivariance

T_g is an action of a group G on the space of inputs and outputs. In case of graphs, G is the symmetry group \mathbb{S}_n .



Invariance: $\phi(T_g(f)) = \phi(f)$



Equivariance: $\phi(T_g^{(1)}(f)) = T_g^{(2)}(\phi(f))$

Message Passing Neural Networks and its limitation (1)

To preserve the **permutation invariance**, the aggregation function ψ of MPNNs basically sums the messages from each node's neighborhood. The algorithm is simply expressed in matrix form as:

$$F^t = \sigma(AF^{t-1}W^t)$$

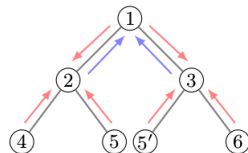
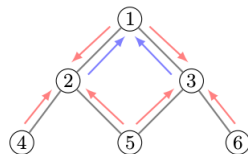
where:

- $A \in \{0, 1\}^{n \times n}$ is the adjacency matrix.
(or graph Laplacian $I_n - D^{-1/2}AD^{-1/2}$)
- $F^t \in \mathbb{R}^{n \times d}$ is the node feature matrix.
- $W^t \in \mathbb{R}^{d \times d'}$ is the weight (channels mixing) matrix, that is learnable.
- σ is the nonlinearity.

Message Passing Neural Networks and its limitation (2)

The summing operator **limits** the representative power of MPNNs such that each node loses their identity after being aggregated. For example:

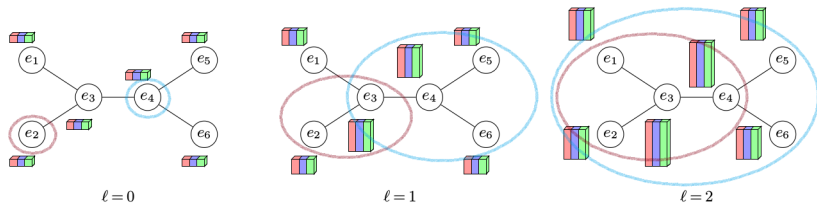
These two graphs are **not** isomorphic, but message passing scheme **fails** to distinguish whether 5 and 5' are the same vertex or not.



Weisfeiler-Lehman isomorphism test fails for highly symmetric structures such as regular graphs.

Covariant Compositional Networks (1)

We propose a new general architecture called **Covariant Compositional Networks** (CCNs) in which the messages are represented by higher order tensors and transform covariantly/equivariantly according to a specific representation of the symmetry group of its receptive field.

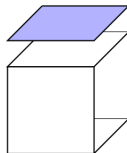


Feature tensors in a **first order** CCN for ethylene (C_2H_4) assuming three channels (red, green, blue).

Covariant Compositional Networks (2)

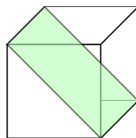
Permutation covariant operators:

1. Projections



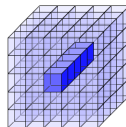
$$C_{i,j} = \sum_a A_{a,i,j}$$

2. Diagonals



$$C_{i,j} = \sum_i A_{i,i,j}$$

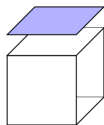
3. Contractions



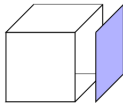
$$C_k = \sum_{i,j} A_{i,j,k}$$

Covariant Compositional Networks (3)

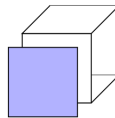
There are six different ways of covariantly reducing (contracting) a third order tensor to a second order tensor:



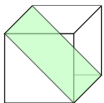
$$C_{i,j} = \sum_a A_{a,i,j}$$



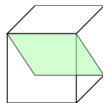
$$C_{i,j} = \sum_j A_{i,a,j}$$



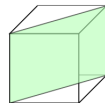
$$C_{i,j} = \sum_k A_{i,j,a}$$



$$C_{i,j} = \sum_i A_{i,i,j}$$



$$C_{i,j} = \sum_i A_{i,j,i}$$

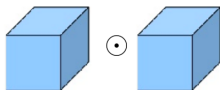


$$C_{i,j} = \sum_i A_{i,j,j}$$

Covariant Compositional Networks (4)

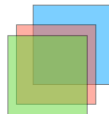
Permutation covariant operators (continued):

4. Hadamard products



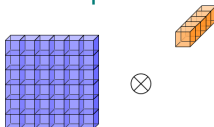
$$C_{i,j,k} = A_{i,j,k} B_{i,j,k}$$

5. Stacking



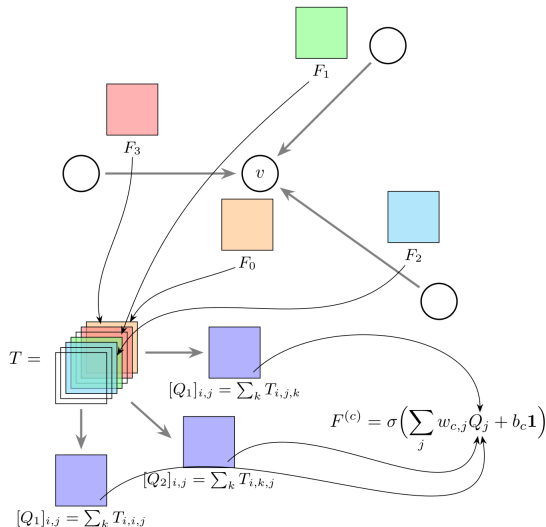
$$C_{i,j,k} = A_{i,j}^{(k)}$$

6. Tensor products



$$C_{i,j,k} = A_{i,j} B_k$$

Covariant Compositional Networks (5)

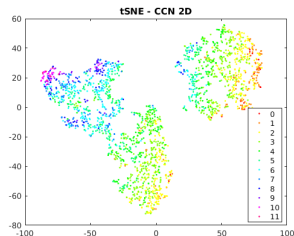
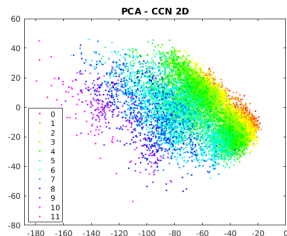


Second-order CCN

Covariant Compositional Networks (6)

Regression results on Havard Clean Energy Project (CEP) dataset

	Test MAE	Test RMSE
Lasso	0.867	1.437
Ridge regression	0.854	1.376
Random forest	1.004	1.799
Gradient boosted trees	0.704	1.005
Weisfeiler-Lehman kernel ^[33]	0.805	1.096
Neural graph fingerprints ^[21]	0.851	1.177
PSCN ($k = 10$) ^[37]	0.718	0.973
Second order CCN (our method)	0.340	0.449



Covariant Compositional Networks (7)

Regression results on QM9 dataset

	CCN	DFT error
α (Bohr ³)	0.22	0.4
C_v (cal/(mol K))	0.07	0.34
G (eV)	0.06	0.1
GAP (eV)	0.12	1.2
H (eV)	0.06	0.1
HOMO (eV)	0.09	2.0
LUMO (eV)	0.09	2.6
μ (Debye)	0.48	0.1
ω_1 (cm ⁻¹)	2.81	28
R_2 (Bohr ²)	4.00	-
U (eV)	0.06	0.1
U_0 (eV)	0.05	0.1
ZPVE (eV)	0.0039	0.0097

(Hy et al., 2018, Kondor et al., 2018)

Software:

- 1 <https://github.com/HyTruongSon/GraphFlow>
(C++/CUDA)
- 2 <https://github.com/HyTruongSon/LibCCNs>
(PyTorch/TensorFlow)

Future works:

- Hypergraph neural networks
 - Simplicial complexes of graphs (algebraic geometry)
 - Hypergraph Fourier Transform
- k-fold atomic interaction in N-Body networks

Multiresolution Matrix Factorization (1)

Not based on the low-rank assumption

$$\Pi \left(\begin{array}{c} \text{[Solid Square]} \\ A \end{array} \right) \Pi^T \approx \left(\begin{array}{c} \text{[Diagonal with Small Square]} \\ U_1^T \end{array} \right) \dots \left(\begin{array}{c} \text{[Diagonal with Small Square]} \\ U_L^T \end{array} \right) \left(\begin{array}{c} \text{[Solid Square]} \\ H \end{array} \right) \left(\begin{array}{c} \text{[Diagonal with Small Square]} \\ U_L \end{array} \right) \dots \left(\begin{array}{c} \text{[Diagonal with Small Square]} \\ U_1 \end{array} \right)$$

MMF of a symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ (Kondor et al., 2014) is:

$$\mathbf{A} = \mathbf{U}_1^T \mathbf{U}_2^T \dots \mathbf{U}_L^T \mathbf{H} \mathbf{U}_L \dots \mathbf{U}_2 \mathbf{U}_1,$$

where:

- Each \mathbf{U}_ℓ is an orthogonal matrix that is a k -point rotation (small k),
- There is a nested sequence of sets $\mathbb{S}_L \subseteq \dots \subseteq \mathbb{S}_1 \subseteq \mathbb{S}_0 = [n]$ such that the coordinates rotated by \mathbf{U}_ℓ are a subset of \mathbb{S}_ℓ ,
- \mathbf{H} is an \mathbb{S}_L -core-diagonal matrix meaning that is diagonal with a an additional small $\mathbb{S}_L \times \mathbb{S}_L$ dimensional “core”.

Multiresolution Matrix Factorization (2)

Finding the best MMF to a symmetric matrix \mathbf{A} involves solving

$$\min_{\substack{\mathbb{S}_L \subseteq \dots \subseteq \mathbb{S}_1 \subseteq \mathbb{S}_0 = [n] \\ \mathbf{H} \in \mathbb{H}_n^{\mathbb{S}_L}; \mathbf{U}_1, \dots, \mathbf{U}_L \in \mathbb{O}}} \|\mathbf{A} - \mathbf{U}_1^T \dots \mathbf{U}_L^T \mathbf{H} \mathbf{U}_L \dots \mathbf{U}_1\|_{\mathcal{F}}.$$

It is equivalent to

$$\min_{\substack{\mathbb{S}_L \subseteq \dots \subseteq \mathbb{S}_1 \subseteq \mathbb{S}_0 = [n] \\ \mathbf{U}_1, \dots, \mathbf{U}_L \in \mathbb{O}}} \|\mathbf{U}_L \dots \mathbf{U}_1 \mathbf{A} \mathbf{U}_1^T \dots \mathbf{U}_L^T\|_{\text{resi}}^2,$$

where $\|\cdot\|_{\text{resi}}^2$ is the squared residual norm

$$\|\mathbf{H}\|_{\text{resi}}^2 = \sum_{i \neq j; (i,j) \notin \mathbb{S}_L \times \mathbb{S}_L} |\mathbf{H}_{i,j}|^2.$$

Multiresolution Matrix Factorization (3)

There are two fundamental difficulties in MMF optimization:

- ① Finding the optimal nested sequence of \mathbb{S}_ℓ is a **combinatorially hard** (e.g., there are $\binom{d_\ell}{k}$ ways to choose k indices out of \mathbb{S}_ℓ);
- ② The solution for \mathbf{U}_ℓ must satisfy the orthogonality constraint such that $\mathbf{U}_\ell^T \mathbf{U}_\ell = \mathbf{I}$.

Multiresolution Matrix Factorization (3)

There are two fundamental difficulties in MMF optimization:

- 1 Finding the optimal nested sequence of \mathbb{S}_ℓ is a **combinatorially hard** (e.g., there are $\binom{d_\ell}{k}$ ways to choose k indices out of \mathbb{S}_ℓ);
- 2 The solution for \mathbf{U}_ℓ must satisfy the orthogonality constraint such that $\mathbf{U}_\ell^T \mathbf{U}_\ell = \mathbf{I}$.

Existing greedy algorithms to approximate MMF (Kondor et al., 2014)

Limitations:

- Greedy heuristics (e.g., clustering) used in selecting k rows/columns for each rotation. **Not globally optimal sequence of \mathbb{S}_ℓ .**
- Iterative algorithm to optimize only 1 rotation at a time. **Not globally optimal sequence of \mathbf{U}_ℓ .**
- Limited to only $k = 2$. **Bigger rotation matrices are better!!**

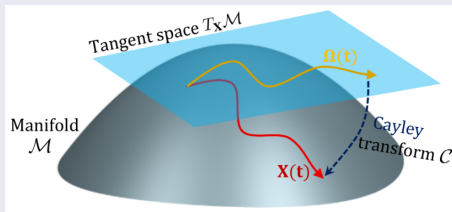
Learnable MMF (1)

The MMF optimization problem is equivalent to

$$\min_{\mathbb{S}_L \subseteq \dots \subseteq \mathbb{S}_1 \subseteq \mathbb{S}_0 = [n]} \min_{\mathbf{U}_1, \dots, \mathbf{U}_L \in \mathbb{O}} \|\mathbf{U}_L \dots \mathbf{U}_1 \mathbf{A} \mathbf{U}_1^T \dots \mathbf{U}_L^T\|_{\text{resi}}^2,$$

Proposal 1 – Stiefel manifold optimization

Given a fixed $\mathbb{S}_L \subseteq \dots \subseteq \mathbb{S}_1 \subseteq \mathbb{S}_0 = [n]$, we use gradient descent algorithm on the Stiefel manifold to optimize all rotations $\{\mathbf{U}_\ell\}_{\ell=1}^L$ *simultaneously*, whilst satisfying the orthogonality constraints.



Learnable MMF (2)

The MMF optimization problem is equivalent to

$$\min_{S_L \subseteq \dots \subseteq S_1 \subseteq S_0 = [n]} \min_{U_1, \dots, U_L \in \mathbb{O}} \|U_L \dots U_1 A U_1^T \dots U_L^T\|_{\text{resi}}^2,$$

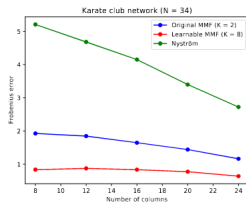
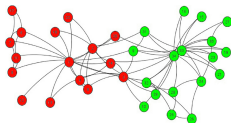
Proposal 2 – Reinforcement Learning

- Formulate the problem of finding the nested sequence of indices as learning a **Markov Decision Process** (MDP).
- **REINFORCE** – Gradient (stochastic) policy method of Reinforcement Learning (RL).
- The RL agent is modeled by **graph neural networks** (GNN).

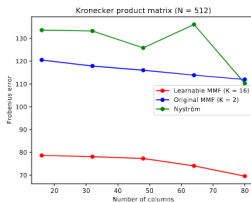
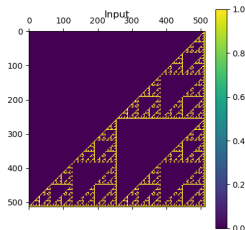
Note: RL has been applied to solve combinatorial, NP-hard, graph problems such as Traveling Saleman (TSP), etc.

Learnable MMF (3)

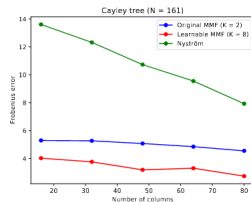
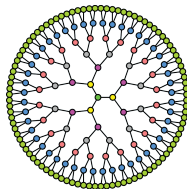
Learnable MMF (red) outperforms the **original MMF** (blue) and the **Nyström method** (green) in matrix approximation:



Karate club network



Kronecker matrix



Cayley tree

Graph Fourier Transform (1)

The eigendecomposition of the normalized graph Laplacian operator $\tilde{\mathbf{L}} = \mathbf{U}^T \mathbf{H} \mathbf{U}$ can be used as the basis of a **Graph Fourier Transform** (GFT) of a graph signal/function $\mathbf{f} : V \rightarrow \mathbb{R}$ (Shuman et al., 2013):

- Forward: $\hat{\mathbf{f}} = \mathbf{U}^T \mathbf{f}$
- Inverse: $\mathbf{f} = \mathbf{U} \hat{\mathbf{f}}$

Graph Fourier Transform (1)

The eigendecomposition of the normalized graph Laplacian operator $\tilde{\mathbf{L}} = \mathbf{U}^T \mathbf{H} \mathbf{U}$ can be used as the basis of a **Graph Fourier Transform** (GFT) of a graph signal/function $\mathbf{f} : V \rightarrow \mathbb{R}$ (Shuman et al., 2013):

- Forward: $\hat{\mathbf{f}} = \mathbf{U}^T \mathbf{f}$
- Inverse: $\mathbf{f} = \mathbf{U} \hat{\mathbf{f}}$
- Graph convolution: $\mathbf{f} *_G \mathbf{g} = \mathbf{U}((\mathbf{U}^T \mathbf{g}) \odot (\mathbf{U}^T \mathbf{f}))$, where \mathbf{g} denotes the convolution kernel, and \odot is the element-wise Hadamard product.
- Based on GFT, (Bruna et al., 2014) and (Defferrard et al., 2016) construct convolutional neural networks (CNNs) **learning on spectral domain** of graphs.

Graph Fourier Transform (2)

Limitations of GFT

- High computational cost:
 - Eigendecomposition of the graph Laplacian has complexity $O(n^3)$
 - “Fourier transform” itself involves multiplying the signal with a **dense** matrix of eigenvectors.
- The graph convolution is **not localized** in the vertex domain, even if the graph itself has well defined local communities.

Graph Fourier Transform (2)

Limitations of GFT

- High computational cost:
 - Eigendecomposition of the graph Laplacian has complexity $O(n^3)$
 - “Fourier transform” itself involves multiplying the signal with a **dense** matrix of eigenvectors.
- The graph convolution is **not localized** in the vertex domain, even if the graph itself has well defined local communities.

Multiresolution Analysis (MRA)

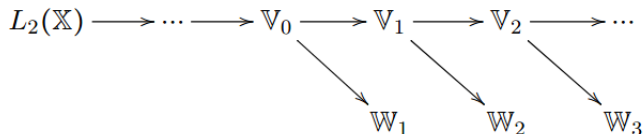
The functional analytic view of wavelets is provided by Multiresolution Analysis (Mallat, 1989), which, similarly to Fourier analysis, is a way of filtering some function space into a sequence of subspaces

$$\dots \subset \mathbb{V}_{-1} \subset \mathbb{V}_0 \subset \mathbb{V}_1 \subset \mathbb{V}_2 \subset \dots$$

Multiresolution analysis (1)

Iteratively, each \mathbb{V}_ℓ is splitted into the orthogonal sum $\mathbb{V}_\ell = \mathbb{V}_{\ell+1} \oplus \mathbb{W}_{\ell+1}$:

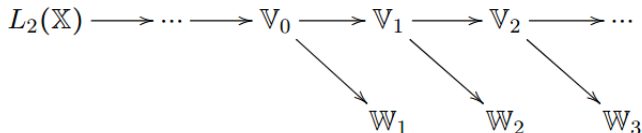
- **Approximation space:** The a smoother part $\mathbb{V}_{\ell+1}$.
- **Detail space:** The rougher part $\mathbb{W}_{\ell+1}$.



Multiresolution analysis (1)

Iteratively, each \mathbb{V}_ℓ is splitted into the orthogonal sum $\mathbb{V}_\ell = \mathbb{V}_{\ell+1} \oplus \mathbb{W}_{\ell+1}$:

- **Approximation space:** The a smoother part $\mathbb{V}_{\ell+1}$.
- **Detail space:** The rougher part $\mathbb{W}_{\ell+1}$.



- Each \mathbb{V}_ℓ has an orthonormal basis $\Phi_\ell \triangleq \{\phi_m^\ell\}_m$ in which each ϕ is called a **father** wavelet.
- Each complementary space \mathbb{W}_ℓ is also spanned by an orthonormal basis $\Psi_\ell \triangleq \{\psi_m^\ell\}_m$ in which each ψ is called a **mother** wavelet.

Multiresolution analysis (2)

Instead of diagonalizing \mathbf{A} in a single step as in PCA, multiresolution analysis will involve a sequence of basis transforms $\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_L$, transforming \mathbf{A} step by step as:

$$\mathbf{A} \rightarrow \mathbf{U}_1 \mathbf{A} \mathbf{U}_1^T \rightarrow \mathbf{U}_2 \mathbf{U}_1 \mathbf{A} \mathbf{U}_1^T \mathbf{U}_2^T \rightarrow \dots \rightarrow \mathbf{U}_L \dots \mathbf{U}_2 \mathbf{U}_1 \mathbf{A} \mathbf{U}_1^T \mathbf{U}_2^T \dots \mathbf{U}_L^T,$$

Multiresolution analysis (2)

Instead of diagonalizing \mathbf{A} in a single step as in PCA, multiresolution analysis will involve a sequence of basis transforms $\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_L$, transforming \mathbf{A} step by step as:

$$\mathbf{A} \rightarrow \mathbf{U}_1 \mathbf{A} \mathbf{U}_1^T \rightarrow \mathbf{U}_2 \mathbf{U}_1 \mathbf{A} \mathbf{U}_1^T \mathbf{U}_2^T \rightarrow \dots \rightarrow \mathbf{U}_L \dots \mathbf{U}_2 \mathbf{U}_1 \mathbf{A} \mathbf{U}_1^T \mathbf{U}_2^T \dots \mathbf{U}_L^T,$$

Each individual rotation $\mathbf{U}_\ell : \mathbb{V}_{\ell-1} \rightarrow \mathbb{V}_\ell \oplus \mathbb{W}_\ell$ is a sparse basis transform that expresses $\Phi_\ell \cup \Psi_\ell$ in the previous basis $\Phi_{\ell-1}$ such that:

$$\phi_m^\ell = \sum_{i=1}^{\dim(\mathbb{V}_{\ell-1})} [\mathbf{U}_\ell]_{m,i} \phi_i^{\ell-1},$$

$$\psi_m^\ell = \sum_{i=1}^{\dim(\mathbb{V}_{\ell-1})} [\mathbf{U}_\ell]_{m+\dim(\mathbb{V}_{\ell-1}),i} \phi_i^{\ell-1}.$$

Multiresolution analysis (3)

In the case \mathbf{A} is the normalized graph Laplacian of a graph $\mathcal{G} = (V, E)$, the wavelet transform (up to level L) expresses a graph signal (function over the vertex domain) $f : V \rightarrow \mathbb{R}$, without loss of generality $f \in \mathbb{V}_0$, as:

$$f(v) = \sum_{\ell=1}^L \sum_m \alpha_m^\ell \psi_m^\ell(v) + \sum_m \beta_m \phi_m^L(v), \quad \text{for each } v \in V,$$

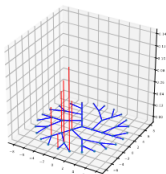
where $\alpha_m^\ell = \langle f, \psi_m^\ell \rangle$ and $\beta_m = \langle f, \phi_m^L \rangle$ are the wavelet coefficients.

Multiresolution analysis (4)

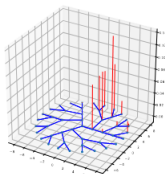
MMF gives us a total of N wavelets:

- L mother wavelets $\overline{\psi} = \{\psi^1, \dots, \psi^L\}$,
- $N - L$ father wavelets $\overline{\phi} = \{\phi_m^L = \mathbf{H}_{m,:}\}_{m \in \mathbb{S}_L}$.

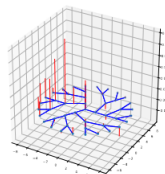
Visualization of some of the wavelets on the Cayley tree of 46 vertices. The low index wavelets (low ℓ) are **highly localized**, whereas the high index ones are smoother and spread out over large parts of the graph.



$\ell = 1$



$\ell = 20$



$\ell = 39$

Wavelet Neural Networks (1)

Analogous to the convolution based on GFT (Bruna et al., 2014), each convolution layer $k = 1, \dots, K$ of our wavelet network transforms an input vector $\mathbf{f}^{(k-1)}$ of size $|V| \times F_{k-1}$ into an output $\mathbf{f}^{(k)}$ of size $|V| \times F_k$ as

$$\mathbf{f}_{:,j}^{(k)} = \sigma \left(\mathbf{W} \sum_{i=1}^{F_{k-1}} \mathbf{g}_{i,j}^{(k)} \mathbf{W}^T \mathbf{f}_{:,i}^{(k-1)} \right) \quad \text{for } j = 1, \dots, F_k,$$

where $\mathbf{W} = [\bar{\phi}, \bar{\psi}]$ is our wavelet basis matrix, $\mathbf{g}_{i,j}^{(k)}$ is a parameter/filter in the form of a diagonal matrix, and σ is an element-wise linearity.

Fast Wavelet Transform

Since the wavelet basis is **sparse**, the wavelet transform can be implemented efficiently by sparse matrix multiplication.

Wavelet Neural Networks (2)

Only **4.69%** and **15.25%** of elements of the wavelet basis are non-zero in Citeseer and Cora, while the GFT basis is completely dense.

Method	Cora	Citeseer
MLP	55.1%	46.5%
ManiReg (Belkin et al., 2006)	59.5%	60.1%
SemiEmb (Weston et al., 2008)	59.0%	59.6%
LP (Zhu et al., 2003)	68.0%	45.3%
DeepWalk (Perozzi et al., 2014)	67.2%	43.2%
ICA (Getoor, 2005)	75.1%	69.1%
Planetoid (Yang et al., 2016)	75.7%	64.7%
Spectral CNN (Bruna et al., 2014)	73.3%	58.9%
ChebyNet (Defferrard et al., 2016)	81.2%	69.8%
GCN (Kipf and Welling, 2017)	81.5%	70.3%
MoNet (Monti et al., 2017)	81.7%	N/A
GWNN (Xu et al., 2019)	82.8%	71.7%
MMF₁	84.35%	68.07%
MMF₂	84.55%	72.76%
MMF₃	87.59%	72.90%

Node classification on citation graphs

(Hy & Kondor, 2021)

Wavelet Neural Networks (3)

Average percentages of non-zero elements of the wavelet basis:

19.23% (MUTAG), 18.18% (PTC), 2.26% (PROTEINS), 11.43% (NCI1)

Method	MUTAG	PTC	PROTEINS	NCI1
DGCNN (Zhang et al., 2018)	85.83 \pm 1.7	58.59 \pm 2.5	75.54 \pm 0.9	74.44 \pm 0.5
PSCN (Niepert et al., 2016)	88.95 \pm 4.4	62.29 \pm 5.7	75 \pm 2.5	76.34 \pm 1.7
DCNN (Atwood and Towsley, 2016)	N/A	N/A	61.29 \pm 1.6	56.61 \pm 1.0
CCN (Kondor et al., 2018)	91.64 \pm 7.2	70.62 \pm 7.0	N/A	76.27 \pm 4.1
GK (Shervashidze et al., 2009)	81.39 \pm 1.7	55.65 \pm 0.5	71.39 \pm 0.3	62.49 \pm 0.3
RW (Vishwanathan et al., 2010)	79.17 \pm 2.1	55.91 \pm 0.3	59.57 \pm 0.1	N/A
PK (Neumann et al., 2015)	76 \pm 2.7	59.5 \pm 2.4	73.68 \pm 0.7	82.54 \pm 0.5
WL (Shervashidze et al., 2011)	84.11 \pm 1.9	57.97 \pm 2.5	74.68 \pm 0.5	84.46 \pm 0.5
IEGN (Maron et al., 2019)	84.61 \pm 10	59.47 \pm 7.3	75.19 \pm 4.3	73.71 \pm 2.6
MMF	86.31 \pm 9.47	67.99 \pm 8.55	78.72 \pm 2.53	71.04 \pm 1.53

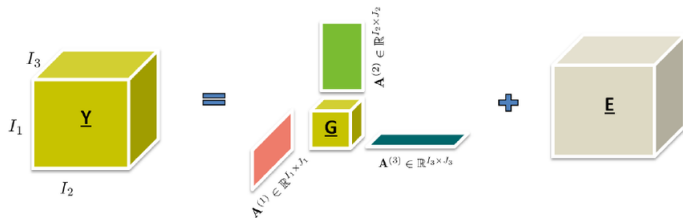
Our WNNs outperform 7/8, 7/8, 8/8, and 2/8 competing methods on molecular graph classification datasets, respectively.

(Hy & Kondor, 2021)

Software:

https://github.com/risilab/Learnable_MMF (PyTorch)

Multiresolution Tensor Factorization:

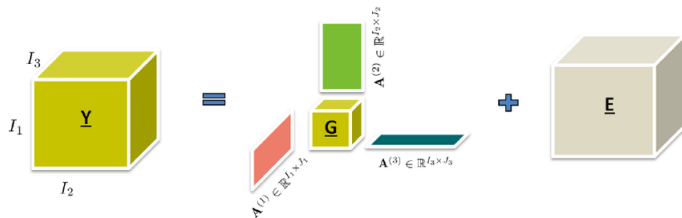


Tensor factorization. Figure taken from (Zhou & Cichocki, 2012).

Software:

https://github.com/risilab/Learnable_MMF (PyTorch)

Multiresolution Tensor Factorization:



Tensor factorization. Figure taken from (Zhou & Cichocki, 2012).

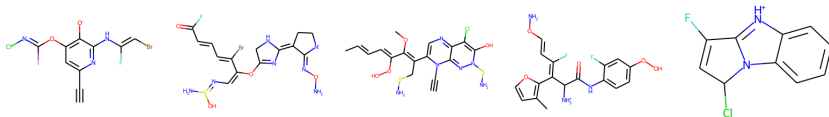
Next question

How to generate new/unseen graphs in an **equivariant** and **multiresolution** manner? Why equivariance is important for graph generation?

Graph generation (1)

Goal

Design models that can observe a set of graphs $\{\mathcal{G}_1, \dots, \mathcal{G}_n\}$ and learn to generate graphs with similar characteristics as this training set.



Look-alike molecules generated from Multiresolution VAE trained on ZINC dataset

Graph generation (2)

Methods:

- ① Traditional graph generation approaches:
 - Erdős-Rényi (ER) Model
 - Stochastic Block Models (SBM)

Graph generation (2)

Methods:

- ① Traditional graph generation approaches:
 - Erdős-Rényi (ER) Model
 - Stochastic Block Models (SBM)

Limitation:

- Rely on a fixed, hand-crafted generation process.
- Lack the ability to learn a generative model from data.

Graph generation (2)

Methods:

① Traditional graph generation approaches:

- Erdős-Rényi (ER) Model
- Stochastic Block Models (SBM)

Limitation:

- Rely on a fixed, hand-crafted generation process.
- Lack the ability to learn a generative model from data.

② Deep generative models:

- **All-at-once:**
 - Generate the whole adjacency matrix with all node (atomic) features
 - VAEs, GANs

Graph generation (2)

Methods:

① Traditional graph generation approaches:

- Erdős-Rényi (ER) Model
- Stochastic Block Models (SBM)

Limitation:

- Rely on a fixed, hand-crafted generation process.
- Lack the ability to learn a generative model from data.

② Deep generative models:

• **All-at-once:**

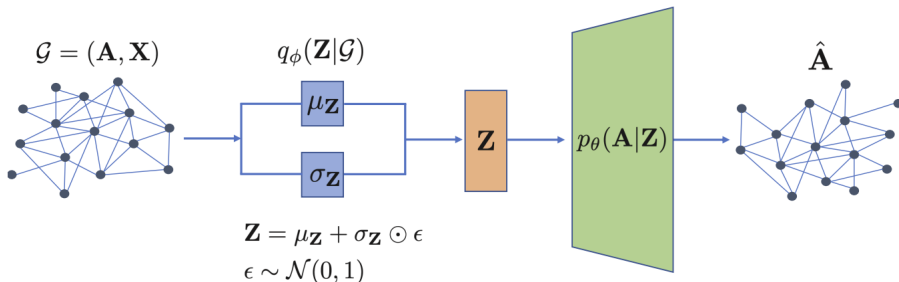
- Generate the whole adjacency matrix with all node (atomic) features
- VAEs, GANs

• **Autoregressive:**

- Generate a graph *incrementally* by adding one node/edge at a time.
- Reinforcement Learning, LSTM-based language models, GRNN, etc.

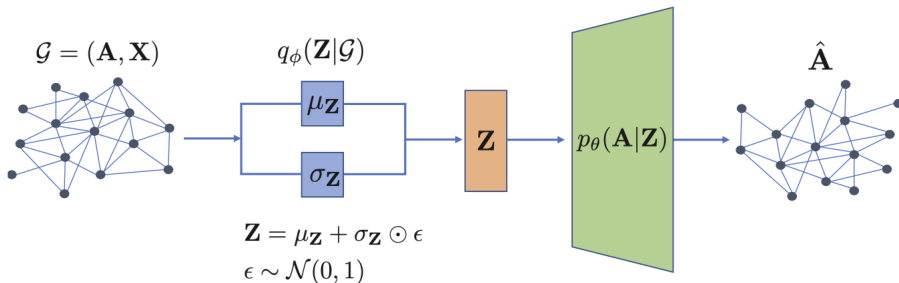
Graph Variational Autoencoder (1)

- **Probabilistic encoder** $q_{\theta}(\mathbf{Z}|\mathcal{G})$ defines a distribution over *latent representations*. We specify the latent conditional distribution as $\mathbf{Z} \sim \mathcal{N}(\mu_{\theta}(\mathcal{G}), \sigma_{\theta}(\mathcal{G}))$, where μ_{θ} and σ_{θ} are neural networks that output the mean and variance parameters for a normal distribution.



Graph Variational Autoencoder (2)

- **Probabilistic decoder** $p_{\theta}(\mathbf{A}|\mathbf{Z})$, from which we can sample realistic graphs (i.e. adjacency matrices) $\hat{\mathbf{A}} \sim p_{\theta}(\mathbf{A}|\mathbf{Z})$ by conditioning on a latent variable \mathbf{Z} .
- **Prior distribution** $p(\mathbf{Z})$ over the latent space. Assume $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$.



Graph Variational Autoencoder (3)

Given a set of training graphs $\{\mathcal{G}_1, \dots, \mathcal{G}_n\}$, we can train a VAE model by maximizing the evidence likelihood lower bound (ELBO):

$$\mathcal{L} = \sum_{\mathcal{G}_i \in \{\mathcal{G}_1, \dots, \mathcal{G}_n\}} \mathbb{E}_{q_{\theta}(\mathbf{Z}|\mathcal{G}_i)}[p_{\theta}(\mathcal{G}_i|\mathbf{Z})] - \text{KL}(q_{\theta}(\mathbf{Z}|\mathcal{G}_i)||p(\mathbf{Z}))$$

Basic idea:

- $\mathbb{E}_{q_{\theta}(\mathbf{Z}|\mathcal{G}_i)}[p_{\theta}(\mathcal{G}_i|\mathbf{Z})]$: Maximize the reconstruction ability of our decoder.
- $\text{KL}(q_{\theta}(\mathbf{Z}|\mathcal{G}_i)||p(\mathbf{Z}))$: Minimize the KL-divergence between our posterior latent distribution $q_{\theta}(\mathbf{Z}|\mathcal{G}_i)$ and the prior $p(\mathbf{Z})$.

Why do we need equivariant generation? (1)

Suppose that we want to map the latent vector $\mathbf{z}_{\mathcal{G}}$ to a matrix $\hat{\mathbf{A}} \in [0, 1]^{|\mathcal{V}| \times |\mathcal{V}|}$ of edge probabilities. The posterior distribution:

$$p_{\theta}(\mathcal{G}|\mathbf{z}_{\mathcal{G}}) = \prod_{(u,v) \in \mathcal{V} \times \mathcal{V}} \hat{\mathbf{A}}_{u,v} \mathbf{A}_{u,v} + (1 - \hat{\mathbf{A}}_{u,v})(1 - \mathbf{A}_{u,v})$$

where \mathbf{A} denotes the true adjacency, and $\hat{\mathbf{A}}$ denotes the predicted edge probabilities.

Problem

But we do **not** know the correct ordering of nodes.

Why do we need equivariant generation? (2)

Graph matching problem (**NP-hard**, quadratic assignment problem):

$$p_{\theta}(\mathcal{G}|\mathbf{z}_{\mathcal{G}}) = \max_{\pi \in \Pi} \prod_{(u,v) \in \mathcal{V} \times \mathcal{V}} \hat{\mathbf{A}}_{u,v}^{\pi} \mathbf{A}_{u,v} + (1 - \hat{\mathbf{A}}_{u,v}^{\pi})(1 - \mathbf{A}_{u,v})$$

Approximate solution: Specify a set of particular orderings $\{\pi_1, \dots, \pi_n\}$
(this is also how autoregressive methods work in practice)

$$p_{\theta}(\mathcal{G}|\mathbf{z}_{\mathcal{G}}) \approx \sum_{\pi_i \in \{\pi_1, \dots, \pi_n\}} \prod_{(u,v) \in \mathcal{V} \times \mathcal{V}} \hat{\mathbf{A}}_{u,v}^{\pi_i} \mathbf{A}_{u,v} + (1 - \hat{\mathbf{A}}_{u,v}^{\pi_i})(1 - \mathbf{A}_{u,v})$$

Why do we need equivariant generation? (2)

Graph matching problem (**NP-hard**, quadratic assignment problem):

$$p_{\theta}(\mathcal{G}|\mathbf{z}_{\mathcal{G}}) = \max_{\pi \in \Pi} \prod_{(u,v) \in \mathcal{V} \times \mathcal{V}} \hat{\mathbf{A}}_{u,v}^{\pi} \mathbf{A}_{u,v} + (1 - \hat{\mathbf{A}}_{u,v}^{\pi})(1 - \mathbf{A}_{u,v})$$

Approximate solution: Specify a set of particular orderings $\{\pi_1, \dots, \pi_n\}$
(this is also how autoregressive methods work in practice)

$$p_{\theta}(\mathcal{G}|\mathbf{z}_{\mathcal{G}}) \approx \sum_{\pi_i \in \{\pi_1, \dots, \pi_n\}} \prod_{(u,v) \in \mathcal{V} \times \mathcal{V}} \hat{\mathbf{A}}_{u,v}^{\pi_i} \mathbf{A}_{u,v} + (1 - \hat{\mathbf{A}}_{u,v}^{\pi_i})(1 - \mathbf{A}_{u,v})$$

Almost perfect solution

- Equivariant (higher-order) latent, encoder, and decoder
- Thiede, Hy & Kondor, 2020

Markov Random Fields (1)

Problem with the prior

- $\mathcal{N}(0, 1)$ is not a good prior for graph generation, because each node (atom)'s latent is sampled **independently**.
- We want a new prior $\mathcal{N}(\mu, \Sigma)$ that is both **learnable** and **equivariant**.
- That also requires a new **reparameterization trick**.
- Hy & Kondor, 2021

Markov network

In general, k -th order graph encoders encode an undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ into a k -th order latent $\mathbf{z} \in \mathbb{R}^{n^k \times d_z}$, with learnable parameters θ , can be represented as a parameterized Markov Random Field (MRF) or Markov network.

Markov Random Fields (2)

Hammersley–Clifford theorem

A positive distribution $p(\mathbf{z}) > 0$ satisfies the conditional independent properties of an undirected graph \mathcal{G} iff p can be represented as a product of potential functions ψ , one per *maximal clique*, i.e.,

$$p(\mathbf{z}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{c \in \mathcal{C}} \psi_c(z_c|\theta_c)$$

where \mathcal{C} is the set of all the (maximal) cliques of \mathcal{G} , and $Z(\boldsymbol{\theta})$ is the *partition function* to ensure the overall distribution sums to 1, and given by

$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{z}} \prod_{c \in \mathcal{C}} \psi_c(z_c|\theta_c)$$

Markov Random Fields (3)

Our second order encoder inherits Gaussian MRF as *pairwise* MRF of the following form

$$p(\mathbf{z}|\boldsymbol{\theta}) \propto \prod_{s \sim t} \psi_{st}(z_s, z_t) \prod_t \psi_t(z_t)$$

where

$$\psi_{st}(z_s, z_t) = \exp \left(-\frac{1}{2} z_s \Lambda_{st} z_t \right),$$

$$\psi_t(z_t) = \exp \left(-\frac{1}{2} \Lambda_{tt} z_t^2 + \eta_t z_t \right)$$

are the edge and vertex potentials. The joint distribution can be written in the *information form* of a multivariate Gaussian in which

$$\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}, \quad \boldsymbol{\eta} = \boldsymbol{\Lambda} \boldsymbol{\mu}$$

$$p(\mathbf{z}|\boldsymbol{\theta}) \propto \exp \left(\boldsymbol{\eta}^T \mathbf{z} - \frac{1}{2} \mathbf{z}^T \boldsymbol{\Lambda} \mathbf{z} \right)$$

Markov Random Fields (4)

Our second order encoder produces tensor \mathbf{L} as the second order activation, and set $\Sigma = \mathbf{L}\mathbf{L}^T$ to ensure invertibility. The reparameterization trick is changed to

$$\mathbf{z} = \boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, 1)$$

Markov Random Fields (4)

Our second order encoder produces tensor \mathbf{L} as the second order activation, and set $\mathbf{\Sigma} = \mathbf{L}\mathbf{L}^T$ to ensure invertibility. The reparameterization trick is changed to

$$\mathbf{z} = \boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Hungarian matching

We allow the prior $\mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\mathbf{\Sigma}})$ to be **learnable** in which $\hat{\boldsymbol{\mu}}$ and $\hat{\mathbf{\Sigma}}$ are parameters optimized by back propagation in a data driven manner. We want to solve the following convex optimization problem (new equivariant loss):

$$\min_{\sigma \in \mathbb{S}_n} \mathcal{D}_{\text{KL}}(\mathcal{N}(\mathbf{P}_{\sigma}\boldsymbol{\mu}, \mathbf{P}_{\sigma}\mathbf{\Sigma}\mathbf{P}_{\sigma}^T) || \mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\mathbf{\Sigma}}))$$

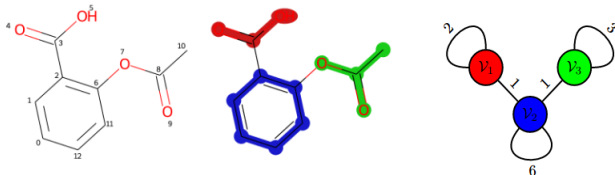
that can be approximated by matching on a bipartite graph.

Multiresolution Graph Network (1)

What we want more?

Learning and then generating graphs in multiple levels of granularity.

The backbone of this coarse-graining architecture, Multiresolution Graph Network (MGN), is the **Learning to cluster** algorithm. The hard clustering can be differentiable (for back-propagation) by the Gumbel-softmax trick.



Aspirin $C_9H_8O_4$, its 3-cluster partition and the corresponding coarsen graph.

Multiresolution Graph Network (2)

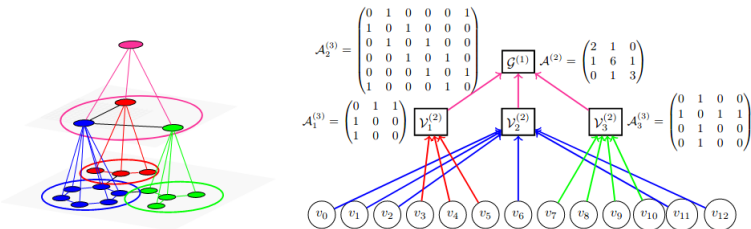
It is desirable to have a *balanced* K-cluster partition in which clusters $\mathcal{V}_1^{(\ell)}, \dots, \mathcal{V}_K^{(\ell)}$ (at the ℓ -th resolution level) have similar sizes that are close to $|\mathcal{V}^{(\ell)}|/K$.

Multiresolution Graph Network (2)

It is desirable to have a *balanced* K-cluster partition in which clusters $\mathcal{V}_1^{(\ell)}, \dots, \mathcal{V}_K^{(\ell)}$ (at the ℓ -th resolution level) have similar sizes that are close to $|\mathcal{V}^{(\ell)}|/K$.

We enforce the clustering procedure to produce a balanced cut by minimizing the following Kullback–Leibler divergence:

$$\mathcal{D}_{KL}(P||Q) = \sum_{k=1}^K P(k) \log \frac{P(k)}{Q(k)}, \quad P = \left(\frac{|\mathcal{V}_1^{(\ell)}|}{|\mathcal{V}^{(\ell)}|}, \dots, \frac{|\mathcal{V}_K^{(\ell)}|}{|\mathcal{V}^{(\ell)}|} \right), \quad Q = \left(\frac{1}{K}, \dots, \frac{1}{K} \right)$$



Multiresolution Equivariant Graph VAE (1)

Based on the construction of multiresolution graph network, the latent hierarchy is partitioned into disjoint groups, $\mathcal{Z}_i = \{\mathcal{Z}_i^{(1)}, \mathcal{Z}_i^{(2)}, \dots, \mathcal{Z}_i^{(L)}\}$ where $\mathcal{Z}_i^{(\ell)}$ is the set of latents at the ℓ -th resolution level. We employ the use of **hierarchical VAEs**.

We write our multiresolution variational lower bound $\mathcal{L}_{\text{MGVAE}}(\phi, \theta)$ on $\log p(\mathcal{G})$ compactly as

$$\mathcal{L}_{\text{MGVAE}}(\phi, \theta) = \sum_i \sum_{\ell} \left[\mathbb{E}_{q_{\phi}(\mathcal{Z}_i^{(\ell)} | \mathcal{G}_i^{(\ell)})} [\log p_{\theta}(\mathcal{G}_i^{(\ell)} | \mathcal{Z}_i^{(\ell)})] - \mathcal{D}_{\text{KL}}(q_{\phi}(\mathcal{Z}_i^{(\ell)} | \mathcal{G}_i^{(\ell)}) || p_0(\mathcal{Z}_i^{(\ell)})) \right]$$

Multiresolution Equivariant Graph VAE (2)

In general, the overall optimization is given as follows:

$$\min_{\phi, \theta, \{\hat{\mu}^{(\ell)}, \hat{\Sigma}^{(\ell)}\}_{\ell}} \mathcal{L}_{\text{MGVAE}}(\phi, \theta; \{\hat{\mu}^{(\ell)}, \hat{\Sigma}^{(\ell)}\}_{\ell}) + \sum_{i, \ell} \lambda^{(\ell)} \mathcal{D}_{\text{KL}}(P_i^{(\ell)} || Q_i^{(\ell)}),$$

where

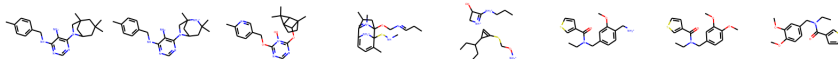
- ϕ denotes all learnable parameters of the encoders,
- θ denotes all learnable parameters of the decoders,
- $\mathcal{D}_{\text{KL}}(P_i^{(\ell)} || Q_i^{(\ell)})$ is the balanced-cut loss for graph \mathcal{G}_i at level ℓ ,
- $\hat{\mu}^{(\ell)}$ and $\hat{\Sigma}^{(\ell)}$ are learnable parameters of the prior in an equivariant manner.

Multiresolution Equivariant Graph VAE (3)

Molecular graph generation results on QM9 & ZINC datasets:

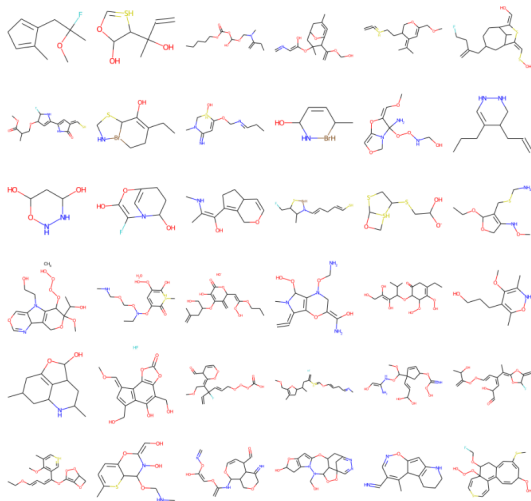
Dataset	Method	Training size	Input features	Validity	Novelty	Uniqueness
QM9	GraphVAE	~ 100K	Graph	61.00%	85.00%	40.90%
	CGVAE			100%	94.35%	98.57%
	MolGAN			98.1%	94.2%	10.4%
	Autoregressive MGN	10K		100%	95.01%	97.44%
	All-at-once MGVAE			100%	100%	95.16%
ZINC	GraphVAE	~ 200K	Graph	14.00%	100%	31.60%
	CGVAE			100%	100%	99.82%
	JT-VAE			100%	-	-
	Autoregressive MGN	1K		100%	99.89%	99.69%
	All-at-once MGVAE	10K	Chemical	99.92%	100%	99.34%

Interpolation on the latent:



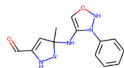
(Hy & Kondor, 2021)

Multiresolution Equivariant Graph VAE (4)

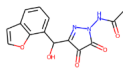


Some generated examples on ZINC by the all-at-once MGVAE with second order \mathbb{S}_n -equivariant decoders.

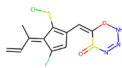
Multiresolution Equivariant Graph VAE (5)



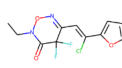
QED = 0.710



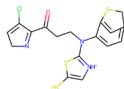
QED = 0.790



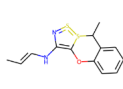
QED = 0.850



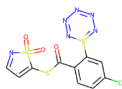
QED = 0.859



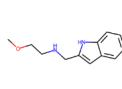
QED = 0.730



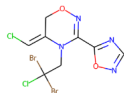
QED = 0.901



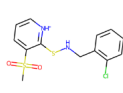
QED = 0.786



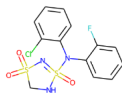
QED = 0.729



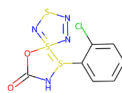
QED = 0.703



QED = 0.855



QED = 0.895



QED = 0.809

Some generated molecules on ZINC by the autoregressive MGN with high QED (drug-likeness score).

Multiresolution Equivariant Graph VAE (6)

Citation graph link prediction results (AUC & AP):

Dataset	Cora		Citeseer	
Method	AUC (ROC)	AP	AUC (ROC)	AP
SC	84.6 ± 0.01	88.5 ± 0.00	80.5 ± 0.01	85.0 ± 0.01
DW	83.1 ± 0.01	85.0 ± 0.00	80.5 ± 0.02	83.6 ± 0.01
VGAE	90.97 ± 0.77	91.88 ± 0.83	89.63 ± 1.04	91.10 ± 1.02
MGVAE (Spectral)	91.19 ± 0.76	92.27 ± 0.73	90.55 ± 1.17	91.89 ± 1.27
MGVAE (K-Means)	93.07 ± 5.61	92.49 ± 5.77	90.81 ± 1.19	91.98 ± 1.02
MGVAE	95.67 ± 3.11	95.02 ± 3.36	93.93 ± 5.87	93.06 ± 6.33

Multiresolution Equivariant Graph VAE (6)

Citaton graph link prediction results (AUC & AP):

Dataset	Cora		Citeseer	
Method	AUC (ROC)	AP	AUC (ROC)	AP
SC	84.6 ± 0.01	88.5 ± 0.00	80.5 ± 0.01	85.0 ± 0.01
DW	83.1 ± 0.01	85.0 ± 0.00	80.5 ± 0.02	83.6 ± 0.01
VGAE	90.97 ± 0.77	91.88 ± 0.83	89.63 ± 1.04	91.10 ± 1.02
MGVAE (Spectral)	91.19 ± 0.76	92.27 ± 0.73	90.55 ± 1.17	91.89 ± 1.27
MGVAE (K-Means)	93.07 ± 5.61	92.49 ± 5.77	90.81 ± 1.19	91.98 ± 1.02
MGVAE	95.67 ± 3.11	95.02 ± 3.36	93.93 ± 5.87	93.06 ± 6.33

Supervised multiresolution graph nets to predict solubility on ZINC:

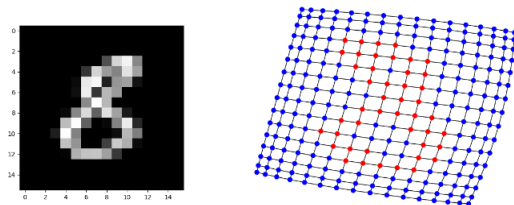
Method	MLP	GCN	GAT	MoNet	DiscenGCN	FactorGCN	GatedGCN _E	MGN
MAE	0.667	0.503	0.479	0.407	0.538	0.366	0.363	0.290

(Hy & Kondor, 2021)

Multiresolution Equivariant Graph VAE (7)

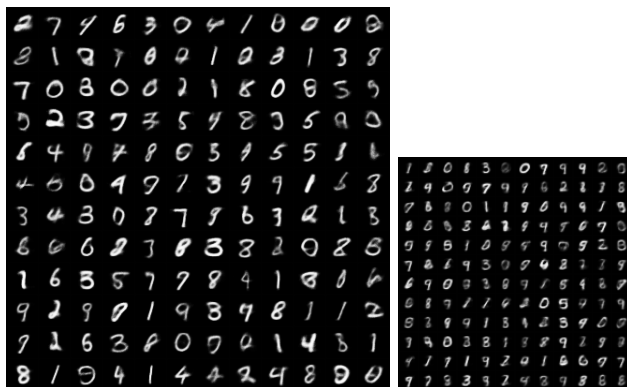
Graph-based image generation at multiple resolutions on MNIST:

Method	FID _↓ (32 × 32)	FID _↓ (16 × 16)	FID _↓ (8 × 8)
DCGAN	113.129	N/A	N/A
VEEGAN	68.749		
PACGAN	58.535		
PresGAN	42.019		
MGVAE	39.474	64.289	39.038



(Hy & Kondor, 2021)

Multiresolution Equivariant Graph VAE (8)



Generated digit images at 32×32 (left) and 16×16 (right) resolutions.

- ① **Drug discovery:** Application of deep generative models on graphs into the lead optimization process that enhances the most promising compounds to improve effectiveness, safety and tolerability.
- ② **Material science:** Constrained generative models to generate stable crystal structures by optimizing the formation energy in the latent space.
- ③ **Proteins:** Multiscale modeling of proteins for the purpose of function prediction and protein design.

Summary:

- Equivariant networks learning graphs
- Multiresolution matrix factorization & graph wavelets
- Equivariant & multiresolution generative models generating graphs

Thank you for your attention!