

# Efficient Graph Generation with Graph Recurrent Attention Networks

In **NeurIPS 2019**

Seonguk Seo

2018-20721

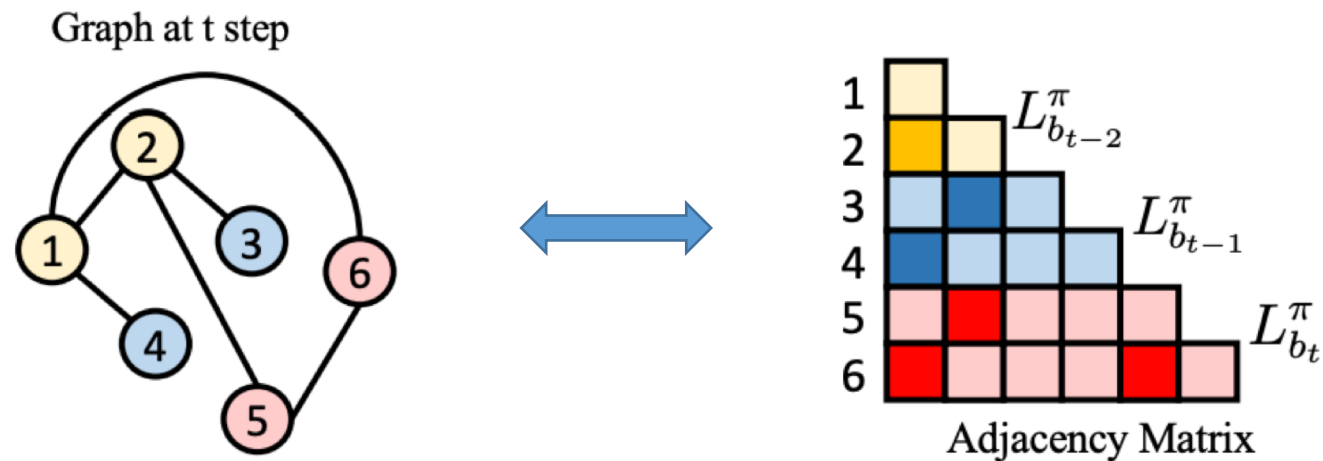
Seoul National University

# Generative Model of Graphs

- Model the distribution of undirected graph  $G = (V, E)$ :

$$P(G) = \sum_{\pi} P(G, \pi) = \sum_{\pi} P(A_{\pi}) = \sum_{\pi} P(L_{\pi})$$

where  $\pi$  denotes a node ordering,  $A_{\pi}$  is an adjacency matrix and  $L_{\pi}$  denotes a lower triangular part of  $A_{\pi}$ .



# Existing Work

- **Graph RNN** [You18]
  - It models graph generation as a sequential process, which accommodate complex dependencies between generated edges.
  - $O(N^2)$  for the best model (not scalable).
  - It has significant bottlenecks in handling long-term dependencies, and the results depend on node orderings.

[You18] Jiaxuan You, Rex Ying, Xiang Ren, William Hamilton, and Jure Leskovec. **Graphrnn: Generating realistic graphs with deep auto-regressive models**. In ICML 2018

# Contributions

- Their model consists of  $O(N)$  auto-regressive generation steps.
- Compared to RNN-based model, they propose an attention-based GNN that better utilizes the topology of the already generated graphs.
- They approximate the likelihood by marginalizing over a family of canonical node orderings.

# Graph Recurrent Attention Networks (GRNN)

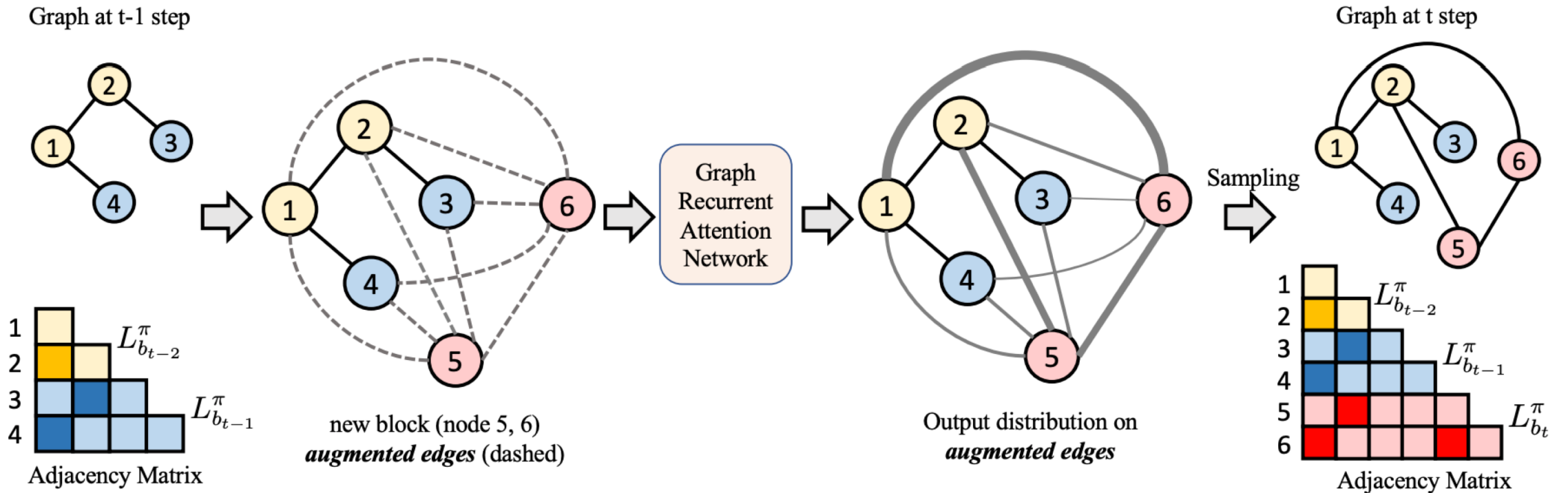
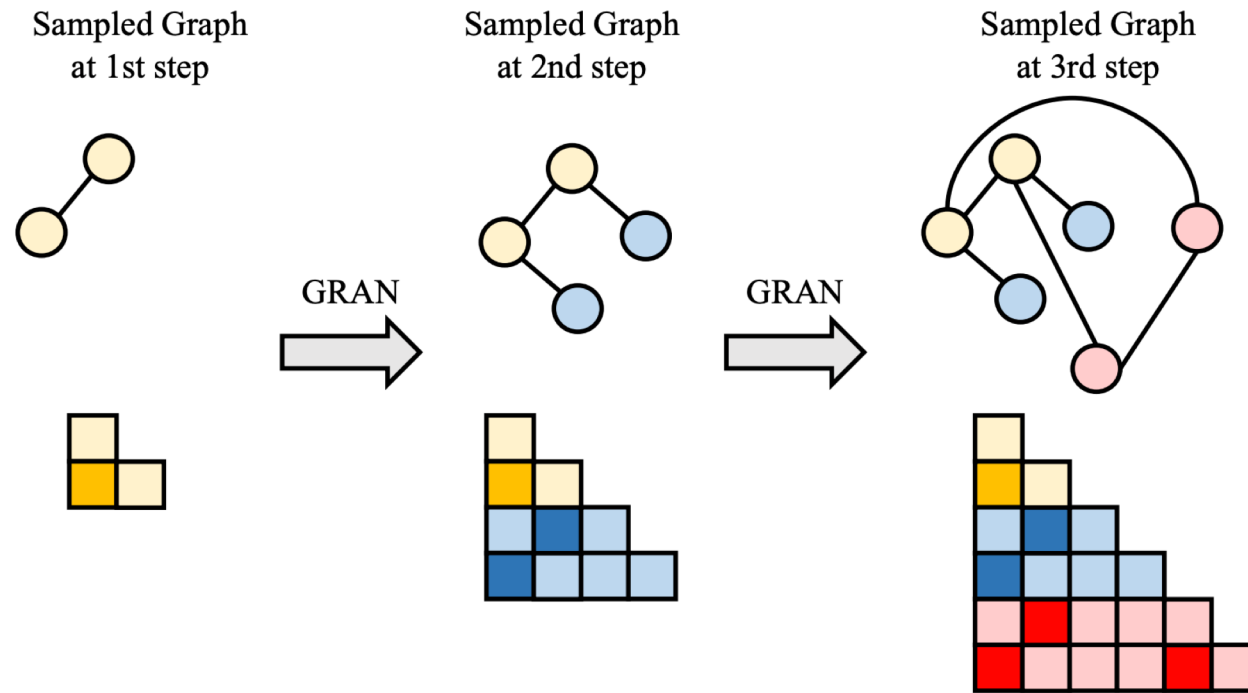


Figure 1: Overview of our model. Dashed lines are *augmented edges*. Nodes with the same color belong to the same block (block size = 2). In the middle right, for simplicity, we visualize the output distribution as a single Bernoulli where the line width indicates the probability of generating the edge.

# Generation Process

- They generate one block of B rows of  $L_\pi$  for each time step.



- Breaking the dependency between generation steps allows parallel training and generated subgraph.
- Varying the block size permits the efficiency-quality trade-off.

# Graph Recurrent Attention Networks (GRNN)

- Initial Node Representation

$$h_{\mathbf{b}_i}^0 = W L_{\mathbf{b}_i}^\pi + b, \quad \forall i < t.$$

- Message Passing

$$\begin{aligned} m_{ij}^r &= f(h_i^r - h_j^r), & a_{ij}^r &= \text{Sigmoid} \left( g(\tilde{h}_i^r - \tilde{h}_j^r) \right), \\ \tilde{h}_i^r &= [h_i^r, x_i], & h_i^{r+1} &= \text{GRU}(h_i^r, \sum_{j \in \mathcal{N}(i)} a_{ij}^r m_{ij}^r). \end{aligned}$$

- Output Distribution

$$\begin{aligned} p(L_{\mathbf{b}_t}^\pi | L_{\mathbf{b}_1}^\pi, \dots, L_{\mathbf{b}_{t-1}}^\pi) &= \sum_{k=1}^K \alpha_k \prod_{i \in \mathbf{b}_t} \prod_{1 \leq j \leq i} \theta_{k,i,j}, \\ \alpha_1, \dots, \alpha_K &= \text{Softmax} \left( \sum_{i \in \mathbf{b}_t, 1 \leq j \leq i} \text{MLP}_\alpha(h_i^R - h_j^R) \right), \\ \theta_{1,i,j}, \dots, \theta_{K,i,j} &= \text{Sigmoid} (\text{MLP}_\theta(h_i^R - h_j^R)) \end{aligned}$$

# Node Representation

- They first compute the initial node representations of the already-generated graph via a linear mapping

$$h_{\mathbf{b}_i}^0 = W L_{\mathbf{b}_i}^\pi + b, \quad \forall i < t.$$

- For current block,  $h_{\mathbf{b}_t}^0 = \mathbf{0}$ .
- In practice, computing  $h_{\mathbf{b}_{t-1}}^0$  alone at  $t^{\text{th}}$  generation step is enough, because  $\{h_{\mathbf{b}_i}^0 | i < t - 1\}$  can be cached from previous steps, which reduces computation.



# Graph Neural Network with Attentive Message Passing

- From the initial Node Representation, all edges associated with the current block are generated using a GNN.
- The  $r^{\text{th}}$  round of message passing is

$$m_{ij}^r = f(h_i^r - h_j^r),$$

$$\tilde{h}_i^r = [h_i^r, x_i],$$

$$a_{ij}^r = \text{Sigmoid} \left( g(\tilde{h}_i^r - \tilde{h}_j^r) \right),$$

$$h_i^{r+1} = \text{GRU}(h_i^r, \sum_{j \in \mathcal{N}(i)} a_{ij}^r m_{ij}^r).$$

# Output Distribution

- After R round message passing, they obtain the final node representation vector, and then model the probability of generating edges in the current block via a mixture of Bernoulli distributions.

$$p(L_{\mathbf{b}_t}^\pi | L_{\mathbf{b}_1}^\pi, \dots, L_{\mathbf{b}_{t-1}}^\pi) = \sum_{k=1}^K \alpha_k \prod_{i \in \mathbf{b}_t} \prod_{1 \leq j \leq i} \theta_{k,i,j},$$

$$\alpha_1, \dots, \alpha_K = \text{Softmax} \left( \sum_{i \in \mathbf{b}_t, 1 \leq j \leq i} \text{MLP}_\alpha(h_i^R - h_j^R) \right),$$

$$\theta_{1,i,j}, \dots, \theta_{K,i,j} = \text{Sigmoid} (\text{MLP}_\theta(h_i^R - h_j^R))$$

# Approximated Likelihood

- They aim to maximize a lower bound:

$$\log p(G) \geq \log \sum_{\pi \in \mathcal{Q}} p(G, \pi)$$

where the size  $Q$  achieve a tradeoff between tightness of the bound (usually correlated with better model quality) and computational cost.

- They adopts DFS, BFS, k-core and degree descent ordering.

# Experiments

- Dataset
  - Grid : 100 standard 2D grid graphs with  $100 < |V| < 400$
  - Protein : 918 protein graphs with  $100 < |V| < 500$
  - Point Cloud : 41 simulated 3D point clouds with  $|V|_{\text{avg}} \geq 1\text{k}$
- 64% train split, 16% valid split, 20% test split

# Evaluation Metrics

- Compare the distributions of graph statistics **between the generated and ground-truth graphs**, by computing the maximum mean discrepancy (**MMD**) over the following 4 statistics.
  1. Degree distributions
  2. Clustering coefficient distributions
  3. The number of occurrences of all orbits with 4 nodes
  4. Spectra of the graphs (the eigenvalues of normalized graph Laplacian)

# Visualization

- Sample graphs generated by GRAN with its comparison.

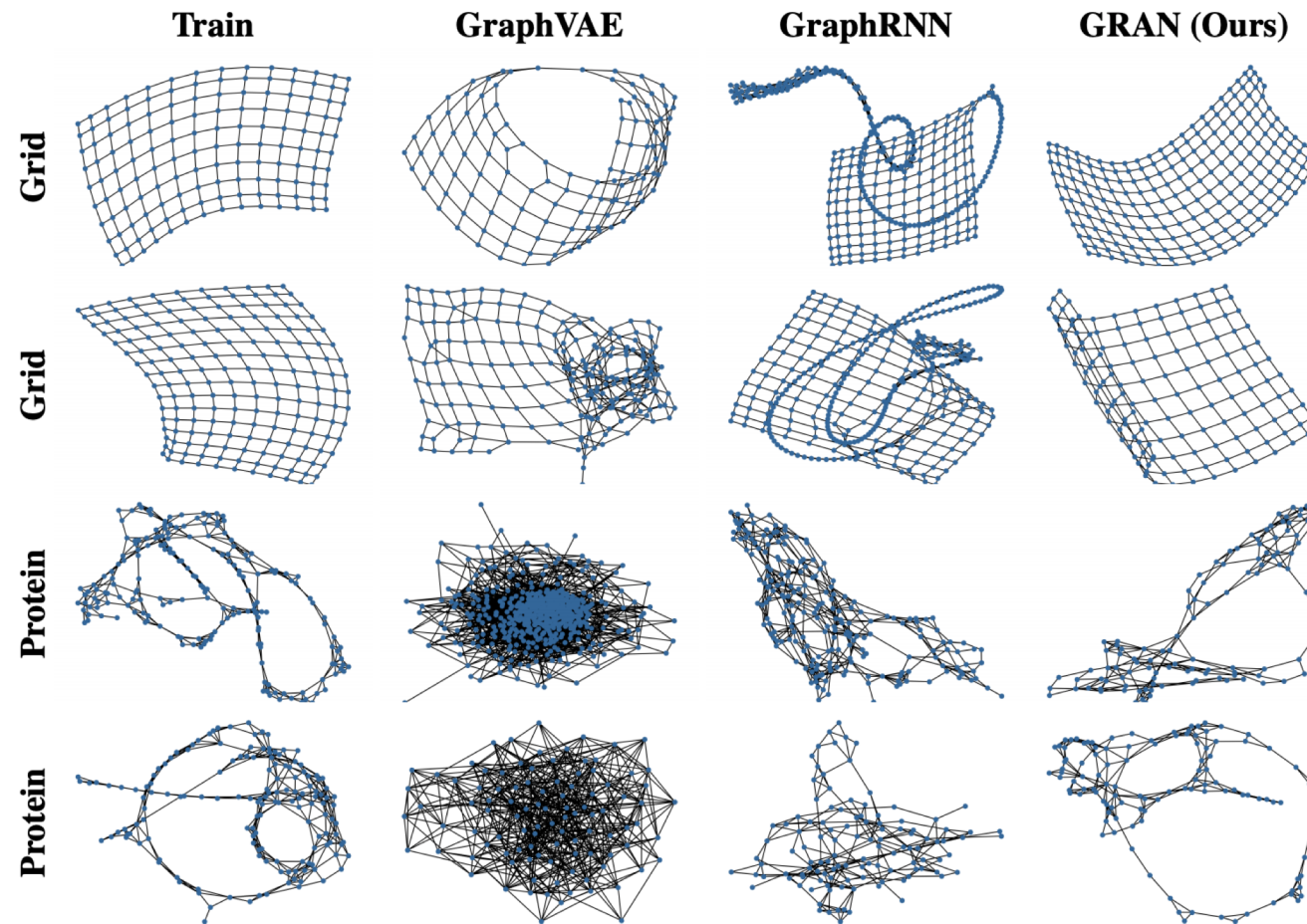


Figure 2: Visualization of sample graphs generated by different models.

# Results

- For all metrics, lower value is preferred.

	Grid				Protein				3D Point Cloud			
	$ V _{\max} = 361,  E _{\max} = 684$ $ V _{\text{avg}} \approx 210,  E _{\text{avg}} \approx 392$				$ V _{\max} = 500,  E _{\max} = 1575$ $ V _{\text{avg}} \approx 258,  E _{\text{avg}} \approx 646$				$ V _{\max} = 5037,  E _{\max} = 10886$ $ V _{\text{avg}} \approx 1377,  E _{\text{avg}} \approx 3074$			
	Deg.	Clus.	Orbit	Spec.	Deg.	Clus.	Orbit	Spec.	Deg.	Clus.	Orbit	Spec.
Erdos-Renyi	0.79	2.00	1.08	0.68	$5.64e^{-2}$	1.00	1.54	$9.13e^{-2}$	0.31	1.22	1.27	$4.26e^{-2}$
GraphVAE*	$7.07e^{-2}$	$7.33e^{-2}$	0.12	$1.44e^{-2}$	0.48	$7.14e^{-2}$	0.74	0.11	-	-	-	-
GraphRNN-S	0.13	$3.73e^{-2}$	0.18	0.19	$4.02e^{-2}$	<b><math>4.79e^{-2}</math></b>	0.23	0.21	-	-	-	-
GraphRNN	$1.12e^{-2}$	<b><math>7.73e^{-5}</math></b>	<b><math>1.03e^{-3}</math></b>	<b><math>1.18e^{-2}</math></b>	$1.06e^{-2}$	0.14	0.88	$1.88e^{-2}$	-	-	-	-
GRAN	<b><math>8.23e^{-4}</math></b>	$3.79e^{-3}$	$1.59e^{-3}$	$1.62e^{-2}$	<b><math>1.98e^{-3}</math></b>	$4.86e^{-2}$	<b>0.13</b>	<b><math>5.13e^{-3}</math></b>	<b><math>1.75e^{-2}</math></b>	<b>0.51</b>	<b>0.21</b>	<b><math>7.45e^{-3}</math></b>

Table 1: Comparison with other graph generative models. For all MMD metrics, the smaller the better. \*: our own implementation, -: not applicable due to memory issue, Deg.: degree distribution, Clus.: clustering coefficients, Orbit: the number of 4-node orbits, Spec.: spectrum of graph Laplacian.

# Efficiency vs Sample Quality

- The main limiting factor for graph generation speed is the number of generation steps  $T$ , which is related to the block size  $B$ .
- If  $B$  grows, which can improve speed, the model quality may suffer.
- They propose “stride sampling”, where neighboring blocks have an overlap

