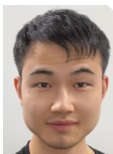# A Large-scale Training Paradigm for Graph Generative Models

Yu Wang[1]　Ryan Rossi[4]　Namyong Park　Huiyuan Chen　Nesreen Ahmed[5]

Puja Trivedi[3]　Franck Dernoncourt[4]　Danai Koutra[3]　Tyler Derr[2]

KIND Lab[1], University of Oregon,
NDS Lab[2], Vanderbilt University
GEMS Lab[3], University of Michigan
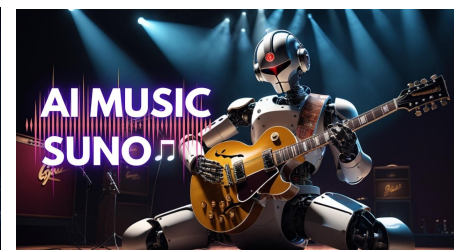Adobe Research[4], Cisco AI Lab[5]

# Background

**Image**

**Language**

**Video**

**Audio**

LLama 3
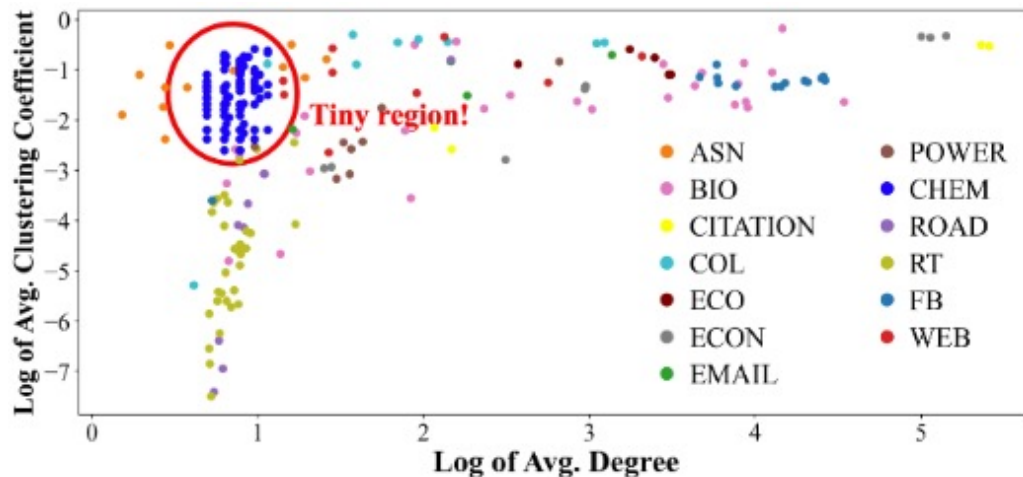
AI MUSIC SUNO

**Previous Small Model (Limited Data)** → **Recent Large Model (Sufficient well-curated Data)**
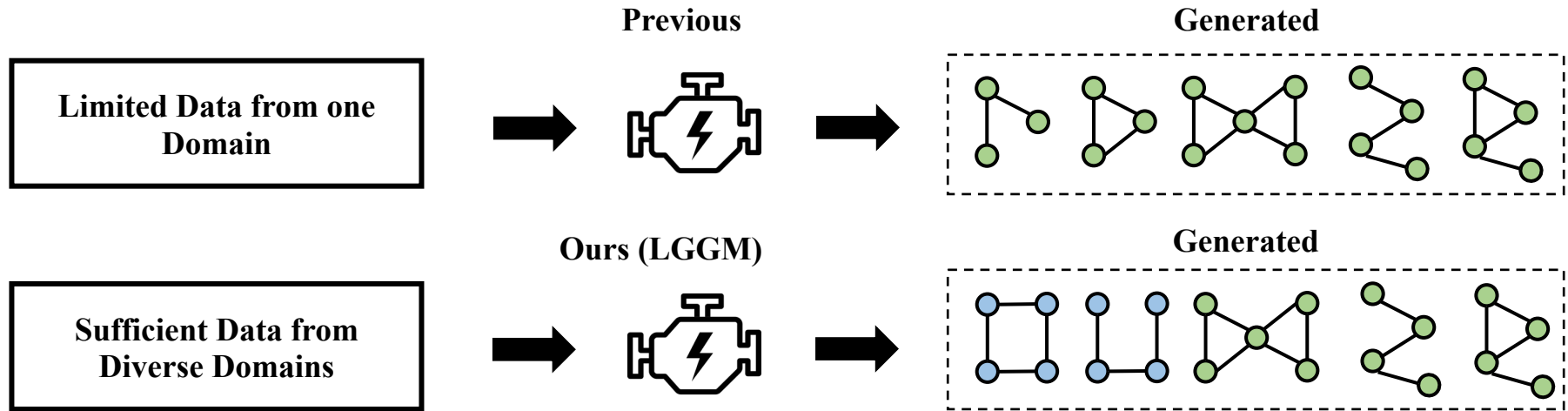
**How about graph generative models?**

| Type | Model | # Domains | Multi-Domain Training |
|------|-------|-----------|----------------------|
| Auto-Regressive | GraphRNN [69] | 2 | ✗ |
| | EdgeRNN [3] | 3 | ✗ |
| | MolRNN [46] | 2 | ✗ |
| VAE | MDVAE [43] | 1 | ✗ |
| | PCVAE [14, 29] | 3 | ✗ |
| | (DE)CO-VAE [19] | 1 | ✗ |
| | GraphVAE [57] | 1 | ✗ |
| GAN | Mol-CycleGAN [40] | 1 | ✗ |
| | LGGAN [18] | 2 | ✗ |
| Flow | GraphNVP [33] | 1 | ✗ |
| | MoFlow [70] | 1 | ✗ |
| | GraphDF [56] | 2 | ✗ |
| Diffusion | GDSS [24] | 3 | ✗ |
| | DiGress [64] | 2 | ✗ |
| | GraphEBM [23] | 1 | ✗ |

ASN: Animal Social Networks   EMAIL: Email Networks   ROAD: Road Networks
FB: Facebook Networks   WEB: Web Graphs   POWER: Power Networks
BIO: Biological Networks   RT: Retweet Networks   ECO: Ecological Networks
ECON: Economic Networks   COL: Collaboration Networks   CITATION: Citation Networks

Tiny region!

- ASN
- BIO
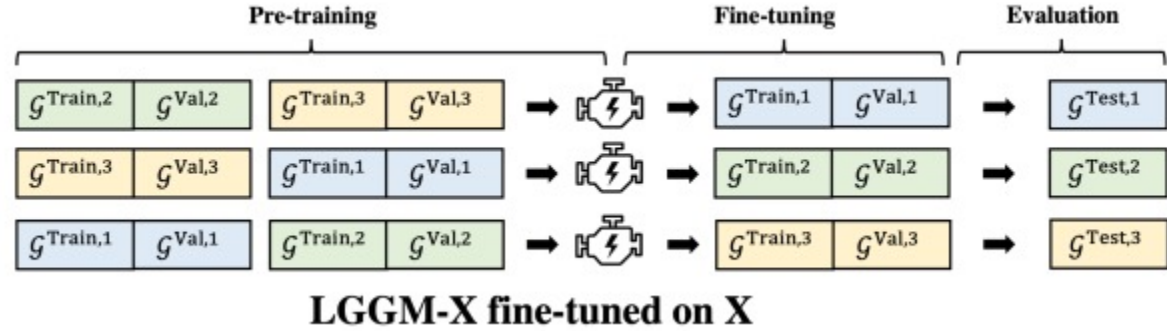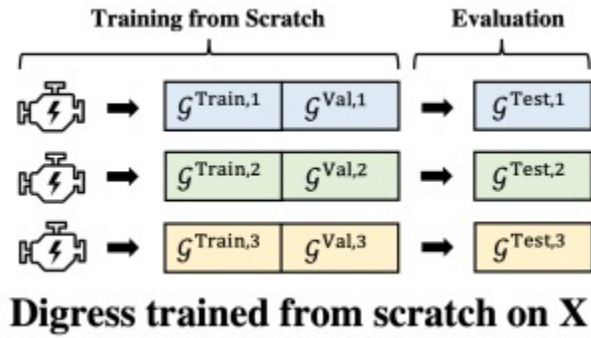- CITATION
- COL
- ECO
- ECON
- EMAIL
- POWER
- CHEM
- ROAD
- RT
- FB
- WEB

Log of Avg. Clustering Coefficient

Log of Avg. Degree

# Our Large-scale Training Paradigm



Previous

Generated

Limited Data from one Domain

Ours (LGGM)

Generated

Sufficient Data from Diverse Domains

| | | | |
|---|---|---|---|
| ANIMAL SOCIAL NETWORKS | 816 | INTERACTION NETWORKS | 29 |
| BIOLOGICAL NETWORKS | 37 | INFRASTRUCTURE NETWORKS | 8 |
| BRAIN NETWORKS | 116 | LABELED NETWORKS | 105 |
| COLLABORATION NETWORKS | 20 | MASSIVE NETWORK DATA | 21 |
| CHEMINFORMATICS | 646 | MISCELLANEOUS NETWORKS | 2669 |
| CITATION NETWORKS | 4 | POWER NETWORKS | 8 |
| ECOLOGY NETWORKS | 6 | PROXIMITY NETWORKS | 13 |
| ECONOMIC NETWORKS | 16 | GENERATED GRAPHS | 221 |
| EMAIL NETWORKS | 6 | RECOMMENDATION NETWORKS | 36 |
| GRAPH 500 | 8 | ROAD NETWORKS | 15 |
| HETEROGENEOUS NETWORKS | 15 | RETWEET NETWORKS | 34 |

| | |
|---|---|
| SCIENTIFIC COMPUTING | 11 |
| SOCIAL NETWORKS | 77 |
| FACEBOOK NETWORKS | 114 |
| TECHNOLOGICAL NETWORKS | 12 |
| WEB GRAPHS | 36 |
| DYNAMIC NETWORKS | 115 |
| TEMPORAL REACHABILITY | 38 |
| BHOSLIB | 36 |
| DIMACS | 78 |
| DIMACS10 | 84 |
| NON-RELATIONAL ML DATA | 211 |

## Degree (DEG), Clustering Coefficient (CC)

$$\mathrm{MMD}(\mathcal{G}_g, \mathcal{G}_r) = \frac{1}{m^2} \sum_{i,j=1}^{m} k(\mathbf{x}_i^r, \mathbf{x}_j^r) + \frac{1}{n^2} \sum_{i,j=1}^{n} k(\mathbf{x}_i^g, \mathbf{x}_j^g) - \frac{2}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} k(\mathbf{x}_i^g, \mathbf{x}_j^r) \qquad k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-d(\mathbf{x}_i, \mathbf{x}_j)/2\sigma^2)$$
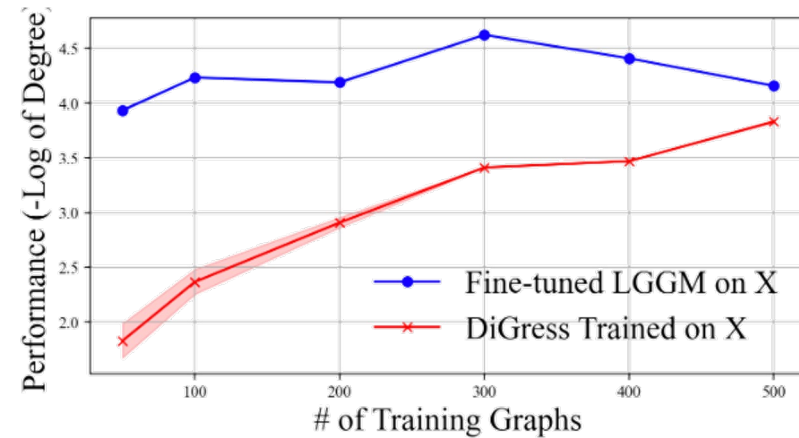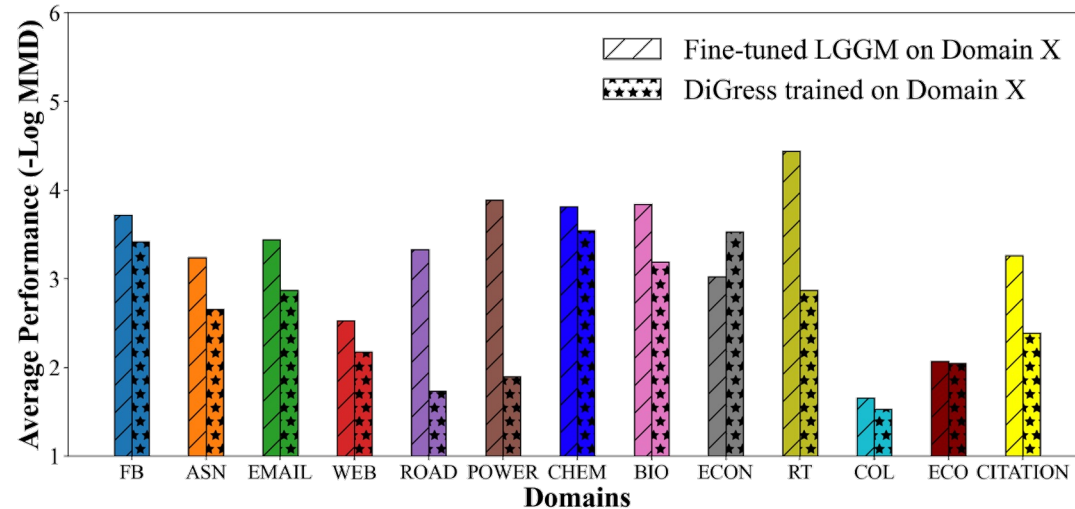
$$\Theta^{\star} = \underset{\Theta}{\operatorname{argmin}} \mathcal{L} = \mathbb{E}_{G \sim P(\mathbb{G})} \mathbb{E}_{t \sim \mathcal{T}} \mathbb{E}_{G^t \sim q(\mathbb{G}^t | \mathbb{G})} \left( - \log p_{\Theta}(G | G^t) \right)$$

$$\Theta^{\star\star} = \underset{\Theta}{\operatorname{argmin}} \mathcal{L} = \mathbb{E}_{\widetilde{G} \sim P(\widetilde{\mathbb{G}})} \mathbb{E}_{t \sim \mathcal{T}} \mathbb{E}_{\widetilde{G}^t \sim q(\widetilde{\mathbb{G}}^t | \widetilde{\mathbb{G}})} \left( - \log p_{\Theta}(\widetilde{G} | \widetilde{G}^t) \right)$$

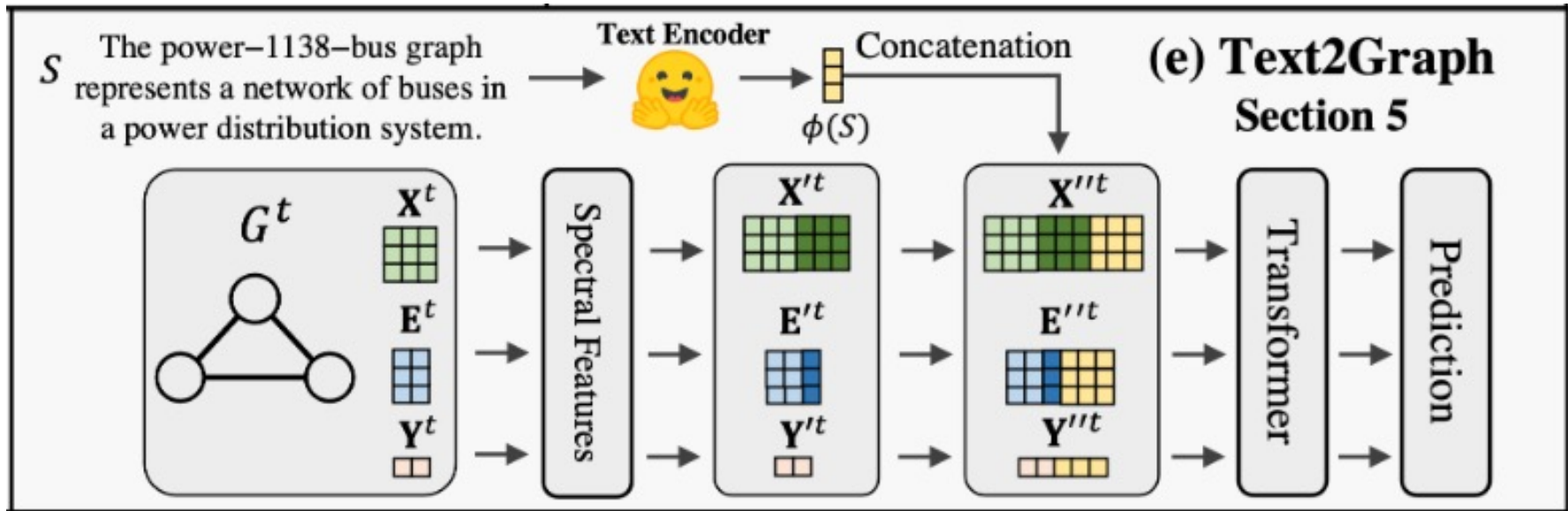Generated Images are of Low-quality due to Online Free Version



**Images have their semantic-meaning for us to specify**

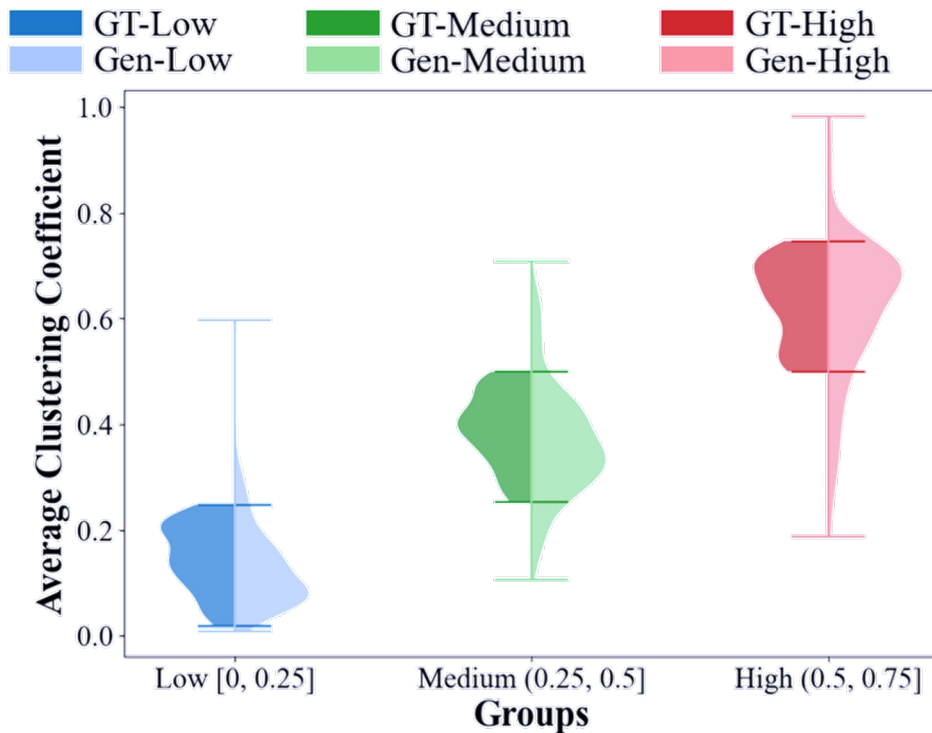**Graphs have their own meta-data: Degree/Density/Clustering Coefficient/Domain…**

**Pretrained encoder to obtain text embeddings**

**Fuse the text embeddings into the latent diffusion**

**Control the Clustering Coefficient (somewhat # of triangles) in the graph**



**GT – Ground-truth ones**

**Gen – Generated ones**
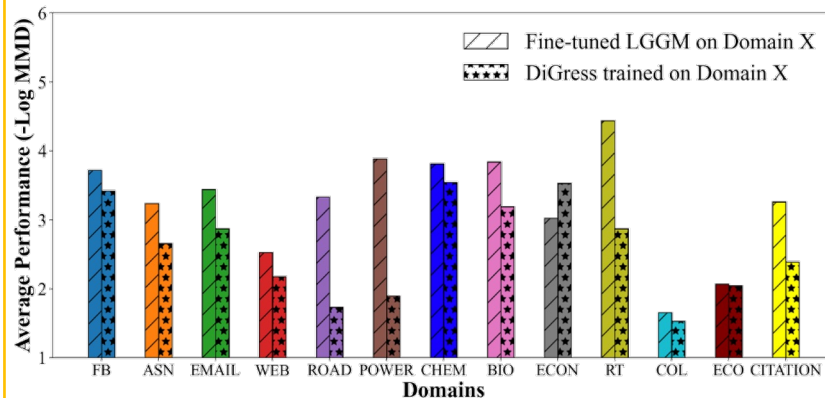
# Summary of LGGM

## Method - LGGM



(a) Graph Universe
- Road
- Power
- Facebook
- Economic
- Collaboration
- Ecological
- Web
- Animal
- Chemical
- Citation
- Biological
- Retweet
- Email

(c) Strategy of our LGGM

(e) Text2Graph — Section 5

$S$ — The power−1138−bus graph represents a network of buses in a power distribution system.

https://kindlab-fly.github.io/

**Knowledge Intelligence for Discovery and Decision-making (KIND) Lab**

Adobe

NSF

## Better Generative Performance



- Fine-tuned LGGM on Domain X
- DiGress trained on Domain X

Domains: FB, ASN, EMAIL, WEB, ROAD, POWER, CHEM, BIO, ECON, RT, COL, ECO, CITATION

## Text-to-Graph Generation Control



- GT-Low / Gen-Low
- GT-Medium / Gen-Medium
- GT-High / Gen-High

Groups: Low [0, 0.25], Medium (0.25, 0.5], High (0.5, 0.75]

**Low-CC**
(20, 0.2250)
(50, 0.0691)
(100, 0.0088)

**Medium-CC**
(20, 0.3645)
(50, 0.4350)
(100, 0.2649)

**High-CC**
(20, 0.5982)
(50, 0.6278)
(100, 0.7372)