INTRODUCTION
oo
MOTIVATION
o
DATASET
ooo
METHODOLOGY
ooo
RESULTS
ooo

# HETEROGENEOUS GRAPH GENERATION: A HIERARCHICAL APPROACH USING NODE FEATURE POOLING

Hritaban Ghosh[*]     Changyu Chen[†]     Arunesh Sinha[‡]
Shamik Sural[*]

[*]Indian Institute of Technology Kharagpur, India

[†]Singapore Management University, Singapore

[‡]Rutgers University, USA

July 02, 2025

1 INTRODUCTION

2 MOTIVATION

3 DATASET

4 METHODOLOGY

5 RESULTS

# INTRODUCTION

◄ Heterogeneous graphs can be used to model systems in various domains such as social networks, recommendation systems, and biological networks.

◄ Generating realistic heterogeneous graphs that capture complex interactions among diverse entities is a difficult task.

◄ The generator has to capture both the node-type distribution and the feature distribution for each node type.

◄ In this paper, we address the challenges in the generation of heterogeneous graphs employing a two-phase hierarchical approach called HG2NP (Heterogeneous Graph Generation using Node Feature Pooling).

# HETEROGENEOUS GRAPHS

◄ Heterogeneous graphs are also known as multi-type graphs.

◄ Characterized by the presence of multiple types of nodes and edges representing diverse entities and relationships.

◄ We denote a heterogeneous graph by $G = (V, E)$, where $V$ is the set of nodes/vertices and $E$ is the set of edges.

◄ Given a set of $M$ observed heterogeneous graphs $\{G_i\}_{i=1}^{M}$, the problem of generating heterogeneous graphs is to learn the probability distribution $\mathbb{P}$ of these graphs from which new graphs can be sampled $G \sim \mathbb{P}$.

INTRODUCTION
oo
MOTIVATION
●
DATASET
ooo
METHODOLOGY
ooo
RESULTS
ooo

# MOTIVATION

◀ In order to highlight the importance of a heterogeneous graph, let us take a simplified example of a social network.

◀ In any social network, an individual can have any relationship (including none) with another individual. Individuals can be part of a group. Groups can be part of a community.
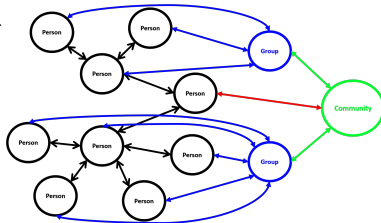


FIGURE 1: Example Heterogeneous Graph

◀ This level of detail and meaningful relationships cannot be achieved by homogeneous graphs without several assumptions and constraints.

◀ The limited availability of real-world data can hinder the development and performance of downstream analytic models, leading to suboptimal results.

◀ Need for methods to generate synthetic heterogeneous graphs that closely resemble their real-world counterparts.

◀ Scarcity of research work on heterogeneous graph generation.

## DATASETS

◄ Digital Bibliography & Library Project (DBLP) and Internet Movie Database (IMDB).

◄ The original DBLP dataset is a heterogeneous graph consisting of authors (4,057 nodes), papers (14,328 nodes), terms (7,723 nodes), and conferences (20 nodes). $\sim$ [26,128 nodes]

◄ The original IMDB dataset is a heterogeneous graph consisting of three types of entities - movies (4,278 nodes), actors (5,257 nodes), and directors (2,081 nodes). $\sim$ [11,616 nodes]

| Node category | Representation | Size of feature vector |
|---|---|---|
| author | Bag-Of-Words of paper keywords | 334 |
| paper | Bag-Of-Words of paper titles | 4231 |
| term | GloVe vector | 50 |

TABLE 1: DBLP graph description

| Node category | Representation | Size of feature vector |
|---|---|---|
| movie | Bag-Of-Words of plot keywords | 3066 |
| actor | Mean of associated movies' features | 3066 |
| director | Mean of associated movies' features | 3066 |

INTRODUCTION
oo

MOTIVATION
o

DATASET
o●o

METHODOLOGY
ooo

RESULTS
ooo

# PROCESSING THE DATASETS - DBLP

◄ Modified the DBLP dataset to include another categorical
feature - the type of the publication.

◄ Used categories and their combinations to segregate the large
graph into smaller components.

- author research area (indicated as author)
- conference
- type of paper publication

◄ Restricted ourselves to work with graphs whose number of
nodes $\leq$ 200.

| Split criteria | Number of graphs with nodes $\leq$ 200 |
|---|---|
| conference | 1 |
| type | 4 |
| author and conference | 25 |
| author and type | 8 |
| conference and type | 29 |
| author, conference and type | 77 |

TABLE 3: DBLP graph sets

# PROCESSING THE DATASETS - IMDB

◄ Modified the IMDB dataset to include three categorical features - year, language and country.

◄ Used categories and their combinations to segregate the large graph into smaller components.
  - movie classes (indicated as movie)
  - year
  - language
  - country

◄ Restricted ourselves to work with graphs whose number of nodes $\leq 200$.

| Split criteria | Number of graphs with nodes $\leq 200$ |
|---|---|
| year | 64 |
| language | 43 |
| country | 56 |
| year and language | 246 |
| year and country | 441 |
| language and country | 111 |
| movie, language and country | 179 |
| year, language and country | 523 |
| movie, year, language and country | 802 |

TABLE 4: IMDB graph sets

# METHODOLOGY

HG2NPis a two-phase scheme, in which the first phase leverages an existing state-of-the-art homogeneous graph generation framework to produce a *skeleton graph* with nodes and edges and node type. Then, in the second phase, we assign feature vectors to the nodes. The second phase is setup as a generative adversarial network.
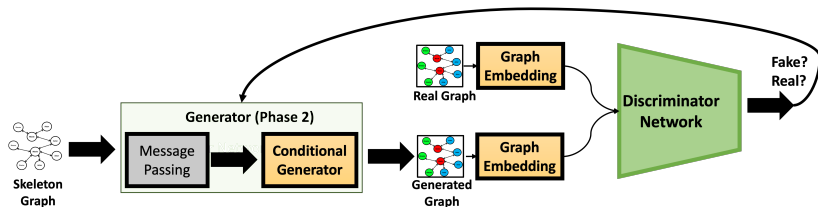


FIGURE 2: Overall generation process for HG2NP. The skeleton graph is the output of Phase 1. Phase 2 uses a GAN architecture for training.

# METHODOLOGY

**Phase 1: Skeleton Graph Generation**

◄ We leverage DiGress, which is designed for generating homogeneous graphs.

◄ It is a diffusion based generative model which is utilized to learn the distribution patterns of node type underlying the heterogeneous graphs in our datasets.

**Phase 2: Heterogeneous Feature Vector Assignment**

◄ Given a skeleton graph from Phase 1, in the second phase, we first utilize a message passing process for updating each node's type vector to obtain a graph with updated node type vectors
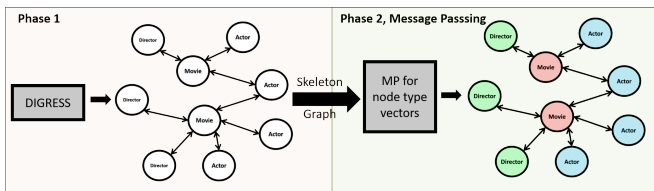


FIGURE 3: In Phase 1, we use DiGress to output a graph with node types only. Subsequently, the first part in Phase 2 is to perform message passing to embed neighbor node type information in each node type vector.

## METHODOLOGY

**Phase 2: Heterogeneous Feature Vector Assignment**

◄ Conditioned on the updated node type vector, we sample a feature vector from a node type specific pool of feature vectors.

◄ In this way, all the nodes of the graph are assigned feature vectors. This yields a generated heterogeneous graph.
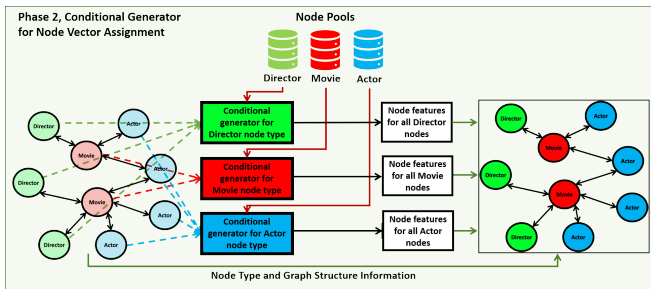


FIGURE 4: The second part is to perform node vector assignment using conditional generators. The node vectors are picked from the feature vector pool based on the probability distribution generated by the conditional generator of that node type.

# MULTIPLE GRAPHS AS TRAIN DATA

## IMDB

| Metrics | Year | | | Country, Language and Movie | | |
|---|---|---|---|---|---|---|
| | VGAE | VGAE-H | HG2NP | VGAE | VGAE-H | HG2NP |
| Degree Dist. | $0.35 \pm 0.004$ | $0.34 \pm 0.008$ | $\mathbf{0.18 \pm 0.001}$ | $0.30 \pm 0.002$ | $0.30 \pm 0.003$ | $\mathbf{0.15 \pm 0.001}$ |
| Clust. Coeff. | $0.92 \pm 0.024$ | $0.88 \pm 0.015$ | $\mathbf{5.2e\text{-}5 \pm 1.4e\text{-}5}$ | $- \pm -$ | $0.75 \pm 0.012$ | $\mathbf{7.8e\text{-}6 \pm 4.6e\text{-}6}$ |
| Spect. Dens. | $0.38 \pm 0.013$ | $0.36 \pm 0.011$ | $\mathbf{0.06 \pm 0.003}$ | $0.23 \pm 0.003$ | $0.24 \pm 0.004$ | $\mathbf{0.04 \pm 0.002}$ |
| Type Degree | $0.28 \pm 0.007$ | $0.26 \pm 0.007$ | $\mathbf{0.05 \pm 3.5e\text{-}4}$ | $0.18 \pm 0.002$ | $0.17 \pm 0.001$ | $\mathbf{0.03 \pm 2.3e\text{-}4}$ |
| W-Dist. | $43.85e\text{-}6 \pm 5.47e\text{-}6$ | $49.64e\text{-}6 \pm 7.95e\text{-}6$ | $\mathbf{2.93e\text{-}6 \pm 2.01e\text{-}6}$ | $39.07e\text{-}6 \pm 3.20e\text{-}6$ | $35.27e\text{-}6 \pm 3.41e\text{-}6$ | $\mathbf{8.63e\text{-}6 \pm 2.35e\text{-}6}$ |
| 1-NN Accuracy | $0.85 \pm 0.112$ | $0.79 \pm 0.019$ | $\mathbf{0.38 \pm 0.196}$ | $0.47 \pm 0.011$ | $0.48 \pm 0.009$ | $\mathbf{0.48 \pm 0.010}$ |
| FID | $2.53e\text{-}3 \pm 5.42e\text{-}5$ | $2.61e\text{-}3 \pm 9.19e\text{-}5$ | $\mathbf{0.99e\text{-}3 \pm 1.56e\text{-}5}$ | $1.18e\text{-}3 \pm 1.20e\text{-}5$ | $1.19e\text{-}3 \pm 2.76e\text{-}5$ | $\mathbf{1.07e\text{-}3 \pm 3.63e\text{-}5}$ |

## DBLP

| Metrics | Author and Conference | | | Author, Conference, and Pub. Type | | |
|---|---|---|---|---|---|---|
| | VGAE | VGAE-H | HG2NP | VGAE | VGAE-H | HG2NP |
| Degree Dist. | $0.41 \pm 0.002$ | $0.43 \pm 0.007$ | $\mathbf{0.16 \pm 0.004}$ | $0.36 \pm 0.003$ | $0.35 \pm 0.006$ | $\mathbf{0.01 \pm 0.001}$ |
| Clust. Coeff. | $- \pm -$ | $1.08 \pm 0.007$ | $\mathbf{3.6e\text{-}4 \pm 4.7e\text{-}5}$ | $- \pm -$ | $0.92 \pm 0.018$ | $\mathbf{6.0e\text{-}5 \pm 1.3e\text{-}5}$ |
| Spect. Dens. | $0.48 \pm 0.023$ | $0.52 \pm 0.020$ | $\mathbf{0.07 \pm 0.006}$ | $0.38 \pm 0.015$ | $0.36 \pm 0.006$ | $\mathbf{1.2e\text{-}3 \pm 3.2e\text{-}4}$ |
| Type Degree | $0.30 \pm 0.004$ | $0.29 \pm 0.016$ | $\mathbf{0.04 \pm 0.001}$ | $0.22 \pm 0.002$ | $0.19 \pm 0.003$ | $\mathbf{1.5e\text{-}3 \pm 9.3e\text{-}5}$ |
| W-Dist. | $28.07e\text{-}6 \pm 3.50e\text{-}6$ | $42.38e\text{-}6 \pm 3.80e\text{-}6$ | $\mathbf{0.00 \pm 0.000}$ | $58.91e\text{-}6 \pm 4.65e\text{-}6$ | $63.32e\text{-}6 \pm 6.31e\text{-}6$ | $\mathbf{0.00 \pm 0.000}$ |
| 1-NN Accuracy | $0.92 \pm 0.112$ | $0.98 \pm 0.011$ | $\mathbf{0.61 \pm 0.024}$ | $1.00 \pm 0.002$ | $1.00 \pm 0.000$ | $\mathbf{0.47 \pm 0.053}$ |
| FID | $3.01e\text{-}3 \pm 1.54e\text{-}5$ | $3.23e\text{-}3 \pm 2.44e\text{-}5$ | $\mathbf{0.87e\text{-}3 \pm 6.41e\text{-}5}$ | $3.45e\text{-}3 \pm 3.72e\text{-}5$ | $3.53e\text{-}3 \pm 2.07e\text{-}5$ | $\mathbf{0.14e\text{-}3 \pm 6.11e\text{-}5}$ |

# SINGLE GRAPH AS TRAIN DATA

**IMDB**

| Metrics | Year | | |
|---|---|---|---|
| | **VGAE** | **NetGAN** | **HG2NP** |
| Degree Dist. | $0.37 \pm 0.004$ | $0.37 \pm 0.004$ | $\mathbf{0.23 \pm 0.002}$ |
| Clust. Coeff. | $0.95 \pm 0.016$ | $0.73 \pm 0.008$ | $\mathbf{6.8e\text{-}5 \pm 7.2e\text{-}6}$ |
| Spect. Dens. | $0.61 \pm 0.178$ | $0.64 \pm 0.168$ | $\mathbf{0.22 \pm 0.131}$ |
| Type Degree | $0.30 \pm 0.006$ | $0.24 \pm 4.6e\text{-}4$ | $\mathbf{0.07 \pm 4.9e\text{-}4}$ |
| W-Dist. | $5.55e\text{-}5 \pm 5.01e\text{-}6$ | $4.52e\text{-}5 \pm 5.06e\text{-}6$ | $\mathbf{1.90e\text{-}5 \pm 3.83e\text{-}6}$ |
| 1-NN Accuracy | $0.71 \pm 0.026$ | $0.78 \pm 0.063$ | $\mathbf{0.55 \pm 0.151}$ |
| FID | $2.61e\text{-}3 \pm 7.44e\text{-}5$ | $2.32e\text{-}3 \pm 1.07e\text{-}5$ | $\mathbf{1.21e\text{-}3 \pm 5.85e\text{-}5}$ |

**DBLP**

| Metrics | Author and Conference | | |
|---|---|---|---|
| | **VGAE** | **NetGAN** | **HG2NP** |
| Degree Dist. | $0.41 \pm 0.003$ | $0.73 \pm 0.000$ | $\mathbf{0.17 \pm 0.001}$ |
| Clust. Coeff. | $-\pm-$ | $2.00 \pm 0.000$ | $\mathbf{0.01 \pm 3.3e\text{-}4}$ |
| Spect. Dens. | $0.68 \pm 0.080$ | $0.75 \pm 0.000$ | $\mathbf{0.29 \pm 0.205}$ |
| Type Degree | $0.30 \pm 0.004$ | $0.36 \pm 4.1e\text{-}4$ | $\mathbf{0.040 \pm 3.1e\text{-}4}$ ] |
| W-Dist. | $2.19e\text{-}5 \pm 3.38e\text{-}6$ | $3.05e\text{-}5 \pm 3.61e\text{-}6$ | $\mathbf{1.28e\text{-}5 \pm 3.58e\text{-}6}$ |
| 1-NN Accuracy | $0.98 \pm 0.011$ | $0.99 \pm 0.005$ | $\mathbf{0.63 \pm 0.033}$ |
| FID | $3.04e\text{-}3 \pm 2.44e\text{-}5$ | $3.11e\text{-}3 \pm 1.20e\text{-}5$ | $\mathbf{1.76e\text{-}3 \pm 15.95e\text{-}5}$ |

INTRODUCTION
OO

MOTIVATION
O

DATASET
OOO

METHODOLOGY
OOO

RESULTS
OO●

## CONCLUDING REMARKS

- ◄ We see that for majority of the cases, our hierarchical model outperforms the baselines.
- ◄ The idea of node pooling is used to sample node feature vectors and assign them to the nodes, thereby reconstructing the heterogeneous graphs.
- ◄ An important limitation, and potential future work, is to design our approach to explicitly work with edge types.
- ◄ Another potential improvement is to have iterative refinement of feature vector assignments. Vectors once assigned are not revisited in HG2NP. Assignment is one shot in the sense that all of them are done in parallel in one go.
- ◄ Overall, our work HG2NPprovides a hierarchical approach that successfully generates heterogeneous graphs with promising results for the domains that we experimented with.

# Thank you