

LeadConvertX: A Multimodal Temporal Heterogeneous Graph Transformer for Explainable CRM Lead Conversion Prediction

Aarush Kukade, Advait Deogade, Atharva Mane, Dr. Saurabh Saoji

Department of Information Technology,
Nutan Maharashtra Institute of Engineering and Technology (NMIET), Pune, India

Abstract- In the modern digital banking era, effective marketing lead generation depends on leveraging heterogeneous, multimodal customer data. Traditional predictive models primarily rely on static, flat tabular attributes, overlooking the relational and temporal information inherent in customer transaction histories and support networks. This paper proposes MTHGT, a Multimodal Temporal Heterogeneous Graph Transformer framework that integrates multimodal data—including structured CRM attributes, sequential transactional records, and unstructured call transcripts—to predict customer lead conversion in banking. The proposed system models customers, transactions, locations, and events as nodes in a heterogeneous graph with relationships based on transactional similarity, campaign logs, and temporal history. Using a multimodal embedding strategy, the model learns customer representations via Graph Transformer layers with type-aware, distance, and temporal bias encodings. Empirical results on the Multimodal Banking Dataset (MBD; 85,620 client-month nodes, 2.26% positive rate) demonstrate that graph-based models outperform tabular baselines on ranking (HGT ROC-AUC of 0.7809 ± 0.0092 and MTHGT ROC-AUC of 0.7763 ± 0.0160 vs. Logistic Regression ROC-AUC of 0.7397 ± 0.0002). Furthermore, MTHGT improves F1-score over HGT (0.0778 ± 0.0225 vs. 0.0651 ± 0.0050) and exposes dynamic modality attributions (CRM features: 25%, dialogue text: 36%, temporal transactions: 39%), enabling explainable CRM lead scoring. The paper details system design, dataset structure, implementation, graph construction methodology, performance evaluation, and outlines a roadmap to bridge the tabular baseline gap using Focal Loss and behavioral k-NN edges.

Keywords— Graph Neural Networks, Graph Transformers, Multimodal Learning, Banking CRM, Lead Generation, Heterogeneous Graph, Explainable AI.

I. INTRODUCTION

The increasing digitization of banking operations has generated vast volumes of multimodal customer data. This data includes structured attributes (demographics, account balances, loan statuses), temporal sequences (daily transactions, monthly activity levels, campaign touchpoints), and unstructured interactions (call-center transcripts, support chat logs, email threads). Traditional machine learning models, such as logistic regression, random forests, or gradient boosted trees, perform well on structured tabular datasets but struggle to capture relational dependencies and unstructured insights simultaneously.

Recent advancements in Graph Neural Networks (GNNs) enable representation learning over connected data, modeling how customer relationships and interactions influence lead conversion likelihood. GNNs aggregate features from neighboring nodes, allowing banks to capture patterns such as peer similarity, campaign propagation, and transactional relationships. However, standard GNN architectures remain

limited in their ability to model long-range relationships, typed edge interactions, temporal sequences of events, and business decision explainability.

This paper introduces a multimodal Graph Transformer framework, **MTHGT** (Multimodal Temporal Heterogeneous Graph Transformer), for marketing lead prediction in the banking sector. The approach integrates multiple data modalities and leverages relational graph structures to improve predictive accuracy and interpretability. We compare the proposed model against classical machine learning (Logistic Regression, Random Forest, XGBoost), homogeneous GNNs (GCN, GraphSAGE, GAT), and heterogeneous GNNs (HGT) on the real-world Multimodal Banking Dataset (MBD) across two complementary tracks:

- **Track A — Tabular and Homogeneous GNN Baselines:** A comparison of traditional machine learning and standard GNN models (GCN, GraphSAGE, GAT) evaluating structured features and similarity-based graph connections.
- **Track B — Heterogeneous Graph Transformer Models:** Evaluation of the type-aware Graph Transformer

(HGT) and the proposed MTHGT model, which fuses structured, text, and temporal transaction sequences into a unified dynamic graph structure.

Objectives

The objectives of this research are:

- **Multimodal Integration:** To design a Graph Transformer-based system architecture that fuses structured CRM features, precomputed dialogue embeddings, and temporal transaction sequences into a unified node representation.
- **Heterogeneous Graph Construction:** To build a dynamic banking graph schema capturing connections among client-month records, transactional events, and geographical contexts.
- **Comprehensive Benchmarking:** To evaluate graph transformer performance against traditional tabular models and standard GNN baselines using ROC-AUC, PR-AUC, F1-score, and Balanced Accuracy under identical splits and seeds.
- **Actionable Explainability:** To extract modality importance weights and attention scores, providing explainable lead predictions for bank marketing agents.
- **Robust Verification:** To validate the system under identical data splits and multiple seeds to guarantee statistical reproducibility.

II, LITERATURE REVIEW

The field of banking Customer Relationship Management (CRM) and marketing intelligence has evolved from simple heuristic-based segmentation to complex deep learning models on graphs. We review the literature across seven main areas to establish the research gap.

Traditional Machine Learning in CRM Lead Scoring

Earlier research on bank marketing relied heavily on tabular machine learning models. Moro et al. (2014) [1] applied logistic regression and decision trees to the UCI Bank Marketing dataset, demonstrating that tabular classifiers can achieve high overall accuracy but struggle to adapt to dynamic customer relationships. Chen and Guestrin (2016) [23] introduced XGBoost, which remains a gold standard for tabular data, utilizing gradient-boosted decision trees. While tree-based models excel at handling non-linear boundary splits of demographic and financial ratios, they suffer from a fundamental limitation: they treat each customer as an independent row, completely ignoring the social, geographical, and transactional linkages between customers.

Homogeneous Graph Neural Networks

To exploit relational structures, GNNs were introduced to capture neighborhood context. Kipf and Welling (2017) [5] established Graph Convolutional Networks (GCN) using localized first-order approximations of spectral graph convolutions. Hamilton et al. (2017) [6] proposed GraphSAGE, which replaced full-graph training with neighborhood sampling, making GNNs inductive and scalable to millions of nodes. Veličković et al. (2018) [4] introduced Graph Attention Networks (GAT), enabling nodes to dynamically weight their neighbors during aggregation. While Homogeneous GNNs model customer similarity graphs successfully, they treat all relationships as a single type, failing to distinguish between diverse connection pathways (e.g., sharing a physical location vs. sharing a transactional product).

Heterogeneous Graph Learning

Real-world CRM systems are heterogeneous, containing multiple node types (customers, campaigns, accounts) and edge types (has_account, responded_to). Wang et al. (2019) [20] introduced the Heterogeneous Graph Attention Network (HAN), using meta-paths to capture semantic relationships. Hu et al. (2020) [7] designed the Heterogeneous Graph Transformer (HGT), modeling node- and edge-type specific attention. HGT projects different types of nodes and edges into type-dependent spaces, dynamically weighting the message passing without requiring manual meta-path design. Despite these advances, standard HGT does not natively fuse multiple modalities (such as sequential transactions and chat texts) into the graph representation.

Graph Transformers

Graph Transformers have emerged to solve GNN limitations like over-smoothing and over-squashing. Dwivedi and Bresson (2020) [6] proposed a generalization of Transformer networks to graphs using Laplacian eigenvectors as positional encodings. Ying et al. (2021) [5] proposed Graphormer, showing that spatial, edge, and centrality encodings allow pure transformers to outperform deep GNNs. Rampásek et al. (2022) [7] formulated GPS (General Powerful Scalable Graph Transformer), which modularizes local message-passing and global attention. Kim et al. (2022) [8] proposed TokenGT, tokenizing nodes and edges together. While highly powerful, these architectures are primarily benchmarked on molecular graphs (e.g., ZINC) and academic citation graphs (e.g., OGB), lacking direct optimization for multimodal, dynamic CRM networks.

Temporal Graph Learning

CRM interactions are inherently temporal. Rossi et al. (2020) [10] proposed Temporal Graph Networks (TGN), maintaining a memory module for nodes to capture historical events. Yu et al. (2023) [11] designed DyGFormer, utilizing a transformer block over historical interaction sequences. Huang et al. (2023) [12] introduced the Temporal Graph Benchmark (TGB) to standardize dynamic evaluations. Temporal GNNs model continuous-time changes but rarely support multi-type heterogeneous relations and multimodal inputs in a single model.

Multimodal Learning and Financial Datasets

Mollaev et al. (2024) [14] introduced the Multimodal Banking Dataset (MBD), proving that joint modeling of structured features, temporal event sequences, and dialogue text yields superior predictive power. Financial CRM applications benefit significantly from text analytics, such as BERT-based call transcript classification (Devlin et al., 2018 [22]), and transaction time-series modeling. Integrating these diverse modalities into a unified graph architecture remains an open challenge.

Explainability and Class Imbalance

CRM systems are highly imbalanced, with conversion rates often below 3%. Chawla et al. (2002) [24] proposed SMOTE to over-sample minority classes, but this can distort graph topology. Lin et al. (2017) [25] introduced Focal Loss to down-weight easy negatives. For model transparency, Ying et al. (2019) [26] developed GNNExplainer, identifying compact subgraphs explaining predictions. Lundberg and Lee (2017) [27] developed SHAP for feature attribution, and Sundararajan et al. (2017) [28] designed Integrated Gradients. MTHGT integrates dynamic modality attributions to provide real-time, inherent explanations directly during the forward pass.

Dataset Description

The experiments are conducted on the **Multimodal Banking Dataset (MBD)** [14]. MBD contains real-world anonymized event sequences and monthly product subscription labels. The working graph contains **85,620 client-month nodes** corresponding to **7,135 unique clients**, with a positive label rate of **2.26%** (extreme class imbalance) and **4,751,538 graph edges**.

MBD Dataset Composition

Data Type	Count / Dimension	Key Attributes	Purpose
CRM Attributes	5 attributes	Transaction counts,	Base tabular features

Data Type	Count / Dimension	Key Attributes	Purpose
		balances, month indices, balance flags	
Call Text	32 or 100-dim	Precomputed support-dialogue embeddings	Unstructured conversation features
Transactions	Sequential	Transaction event frequency and trends	Temporal customer context
Graph Edges	4.75 Million	Temporal links, geohash overlaps, transaction event links	Structural connectivity

Graph Node and Relation Statistics

Node/Edge Type	Count	Description
client_month nodes	85,620	Temporal sequence of client records
event_type nodes	104	Links clients to transactional event types
geo_hash nodes	722	Links clients to physical location hashes

III. METHODOLOGY

The proposed MTHGT methodology follows a structured, mathematically rigorous design optimized for modeling multimodal, temporal, and heterogeneous data in banking CRM environments. The workflow is divided into three core stages: feature preparation, graph schema construction, and joint model optimization.

Preprocessing and Feature Alignment

To ensure different raw modalities can be processed by the unified graph transformer, we apply distinct preprocessing pathways for each source:

- **Structured CRM Attributes:** Demographics and account statistics are represented as a vector $\mathbf{x}_{cm} \in \mathbb{R}^5$, consisting of 'trans_count' (total monthly transactions), 'diff_trans_date' (mean days between events), 'month' (calendar month), 'month_index' (relative timeline step), and 'is_balanced' (indicator flag for account balance).

status). These features are scaled using a standard Z-score normalizer:

$$\mathbf{x}_{\text{crm, scaled}} = \frac{\mathbf{x}_{\text{crm}} - \boldsymbol{\mu}_{\text{train}}}{\boldsymbol{\sigma}_{\text{train}}}$$

where $\boldsymbol{\mu}_{\text{train}}$ and $\boldsymbol{\sigma}_{\text{train}}$ are calculated exclusively on the training split folds (0–2) to prevent data leakage.

- **Sequential Transaction Data:** Historical monthly transaction volumes and trends are represented as sequential time-series matrices. For each client-month node i representing a client at month m , we extract a sequential matrix $\mathbf{X}_{\text{temp}} \in \mathbb{R}^{T \times F}$, where $T = 4$ is the historical lookback window size, and F is the feature dimension of transactions. If a customer has less than 4 months of history, the sequence is padded with zeros and a binary mask vector $\mathbf{m} \in \{0,1\}^T$ is generated to signal invalid steps to downstream temporal layers.
- **Unstructured Call Text:** Client chat transcripts and call-center records are represented via pre-computed embeddings. These vectors $\mathbf{x}_{\text{text}} \in \mathbb{R}^{D_{\text{text}}}$ ($D_{\text{text}} \in \{32,100\}$) represent sentence-transformer projections of support logs corresponding to the active month. If no contact occurred in a given month, a zero vector is assigned.

Heterogeneous Graph Assembly

We model the banking system as a heterogeneous graph $G = (V, E)$ containing multiple node types T_v and edge types T_e .

Node Types (T_v):

- *client_month*: Represents a unique customer record at a specific calendar month. Storing records at the client-month level rather than client-level preserves temporal variations in customer behavior.
- *event_type*: Represents transactional and campaign interaction events (e.g., withdraw, credit deposit, direct deposit, marketing SMS).
- *geo_hash*: Represents standardized spatial divisions using 5-character Geohashes (resolving to $\approx 4.9 \text{ km} \times 4.9 \text{ km}$ areas) where transactions occurred.

Edge Types (T_e):

- (*client_month*,*next_month*,*client_month*): Connects a customer's record at month $t - 1$ to their record at month t , creating a chronological timeline.
- (*client_month*,*prev_month*,*client_month*): The reverse chronological link.
- (*client_month*,*to_event_type*,*event_type*): Directed edge indicating that the client performed a specific transaction event type during the month.
- (*event_type*,*rev_to_event_type*,*client_month*): The reverse transaction-type link.

- (*client_month*,*to_geo_hash*,*geo_hash*): Directed edge indicating a transaction or interaction occurred in that spatial geohash during the month.
- (*geo_hash*,*rev_to_geo_hash*,*client_month*): The reverse spatial mapping.
- (*client_month*,*behavioral_knn*,*client_month*): Directed similarity edges connecting client-months with similar transaction profile distributions, capture homophily (peer-like interaction) without target leakage.

Pipeline Execution Flow

The system processes data in four chronological pipeline steps: 1) *Ingestion & Alignment*: Parquet tables containing customer profiles, transaction details, and dialogue embeddings are joined on the client ID and calendar month. 2) *Preprocessing & Scaling*: Numerical normalizations are applied, sequences are padded, and masked, and text vectors are aligned. 3) *Graph Construction*: Nodes and edges are mapped to integer indices, producing the PyTorch Geometric *HeteroData* object. 4) *Neural Modeling*: Modality-specific encoders project all features to a shared space, the dynamic softmax fusion gate blends them, and stacked Type-Aware Graph Transformer layers aggregate structural representations. 5) *Prediction & Optimization*: The output head computes purchase probabilities, compared against targets via Focal Loss or weighted BCE, and backpropagation updates all parameters.

System Design and Architecture

The system architecture is structured in a multi-layered design to handle raw banking data ingestion and convert it into explainable lead scoring outputs. The architecture consists of four distinct processing layers:

image

Layer Definitions

- **Multimodal Data Ingestion Layer:** Responsible for extracting transaction logs, geographical location coordinates, support dialogue transcripts, and core CRM attributes from file databases. It partitions records by calendar month to enable sequence alignment.
- **Relational Preprocessing & Graph Assembly Layer:** Normalizes numerical values, pads sequential transactions, and matches location coordinates to Geohash-5 indexes. It builds adjacency matrices for temporal chronology, spatial overlaps, transaction links, and behavioral k-NN edges, wrapping them into a PyTorch Geometric graph.
- **MTHGT Neural Compute Engine:** The computational core of the model. It contains the modular feature encoders (MLP for CRM, linear projection for text, and

masked pooling for sequences), the dynamic attention-based softmax fusion gate, and stacked heterogeneous transformer layers with structural and temporal biases.

- **Output & Explainability Rig:** Passes final node representations through a classification MLP head to compute lead conversion probabilities. Simultaneously, it extracts the fusion gate weights $(\alpha_i, \beta_i, \gamma_i)$ and edge attention matrices to output explainability reports for banking agents.

Mathematical Model

Multimodal Encoders

For each client-month node i , the encoders project CRM features \mathbf{x}_{crm} , dialogue embeddings \mathbf{x}_{text} , and transaction sequences \mathbf{x}_{temp} into a shared hidden dimension d :

$$\begin{aligned} \mathbf{h}_{\text{crm}} &= \mathbf{W}_2 \cdot \text{ReLU}(\mathbf{W}_1 \mathbf{x}_{\text{crm}} + \mathbf{b}_1) + \mathbf{b}_2 \\ \mathbf{h}_{\text{text}} &= \mathbf{W}_{\text{text}} \mathbf{x}_{\text{text}} + \mathbf{b}_{\text{text}} \\ \mathbf{h}_{\text{temp}} &= \frac{1}{\sum_{t=1}^T (1 - m_t)} \sum_{t=1}^T (1 - m_t) \\ &\quad \cdot \text{ReLU}(\mathbf{W}_{\text{temp}} \mathbf{x}_{\text{temp},t} + \mathbf{b}_{\text{temp}}) \end{aligned}$$

where $m_t \in \{0,1\}$ represents the padding mask for sequence step t .

Dynamic Cross-Modal Fusion

Rather than using static weight assignments, MTHGT dynamically computes fusion weights $\alpha_i, \beta_i, \gamma_i$ for each node i :

$$\begin{aligned} \mathbf{z}_i &= [\mathbf{h}_{\text{crm},i} \parallel \mathbf{h}_{\text{text},i} \parallel \mathbf{h}_{\text{temp},i}] \\ \mathbf{w}_i &= \mathbf{W}_{\text{fusion}} \mathbf{z}_i + \mathbf{b}_{\text{fusion}} \\ &= \text{softmax}(\mathbf{w}_i) \end{aligned}$$

$$\mathbf{h}_{\text{fused},i} = \alpha_i \mathbf{h}_{\text{crm},i} + \beta_i \mathbf{h}_{\text{text},i} + \gamma_i \mathbf{h}_{\text{temp},i}$$

Type-Aware Graph Attention

Let $\mathbf{h}_i^{(l)}$ be the representation of node i at layer l . Under relation type $r = (T(i), R, T(j))$, the Query, Key, and Value projections are computed as:

$$\mathbf{q}_i = \mathbf{W}_{T(i)}^Q \mathbf{h}_i^{(l)}, \quad \mathbf{k}_j = \mathbf{W}_{T(j)}^K \mathbf{h}_j^{(l)}, \quad \mathbf{v}_j = \mathbf{W}_{T(j)}^V \mathbf{h}_j^{(l)}$$

The attention score from node i to neighbor j is defined as:

$$\begin{aligned} &\text{Attn}(i, j) \\ &= \text{softmax} \left(\frac{\mathbf{q}_i \mathbf{W}_R^{\text{Attn}} \mathbf{k}_j^T + \mathbf{b}_{\text{type}}(r) + \mathbf{b}_{\text{dist}}(\text{dist}_{ij}) + \mathbf{b}_{\text{time}}(\tau_{ij})}{\sqrt{d}} \right) \end{aligned}$$

where $\mathbf{W}_R^{\text{Attn}}$ is the relation-specific attention matrix, $\mathbf{b}_{\text{type}}(r)$ is the edge-type bias, \mathbf{b}_{dist} is the structural distance bias, and \mathbf{b}_{time} is the temporal bias based on relative timestamps τ_{ij} .

Contrastive Learning (InfoNCE Loss)

To align representation views across different times, we enforce an auxiliary contrastive objective. The positive pair (z_i, z_i^+) represents the same customer at adjacent months, whereas z_k^- represents other negative client samples:

$$\mathcal{L}_{\text{con}} = -\log \frac{\exp(\text{sim}(z_i, z_i^+)/\tau)}{\exp(\text{sim}(z_i, z_i^+)/\tau) + \sum_{k=1}^K \exp(\text{sim}(z_i, z_k^-)/\tau)}$$

Objective Function

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{BCE}} + \lambda \mathcal{L}_{\text{con}}$$

Implementation Details

- **Software Stack:** Python 3.10, PyTorch 2.2.0, PyTorch Geometric, HuggingFace Transformers, Pandas, and Scikit-Learn.
- **Hardware Setup:** Experiments executed on local workstation with NVIDIA RTX 3060 (12GB VRAM) and cloud instances (AWS p3.2xlarge).
- **Training Settings:** Models trained for 40 epochs using Adam optimizer, learning rate 1×10^{-3} , weight decay 1×10^{-4} , and batch size scaling.

System Implementation Effort Summary

Category	Technical Operations	Effort (hrs)
Data Engineering	Parquet ETL, sequence parsing, heterogeneous graph construction	120
Model Development	Multi-encoder modeling, dynamic fusion gates, HGT convolution layers	140
Evaluation & Tests	Cross-validation loops, seed evaluations, baseline implementations	80
Reporting & Paper	Outline preparation, mathematical modeling, result aggregation	40

Performance Evaluation

Track A: Tabular and Standard GNN Baselines

Standard machine learning models (XGBoost, Random Forest, Logistic Regression) perform exceptionally well on threshold-based metrics (PR-AUC, F1-score) because they can adjust decision boundaries directly on structured tabular features. Homogeneous GNNs (GCN, GAT) yield lower performance because they flatten the typed relations, resulting in over-smoothing and loss of edge context.

Track B: Heterogeneous Graph Transformer Models

Heterogeneous models (HGT, MTHGT) achieve superior global ranking performance (ROC-AUC **0.7809** and **0.7763**). In addition, the proposed MTHGT model secures a significantly higher F1-score (**0.0778 ± 0.0225**) than standard HGT (**0.0651 ± 0.0050**).

Consolidated Performance Comparison (MBD Test Set)

Model	ROC-AUC	PR-AUC	F1-Score	Balanced Accuracy	MCC
Logistic Regression	0.7397 ± 0.0002	0.0964 ± 0.0004	0.1136 ± 0.0003	0.6768 ± 0.0002	0.1345 ± 0.0003
Random Forest	0.7615 ± 0.0002	0.1017 ± 0.0001	0.0992 ± 0.0013	0.6874 ± 0.0010	0.1268 ± 0.0008
XGBoost	0.7606 ± 0.0004	0.1038 ± 0.0004	0.0954 ± 0.0007	0.6920 ± 0.0008	0.1255 ± 0.0009
GCN	0.7095 ± 0.0030	0.0484 ± 0.0003	0.0593 ± 0.0028	0.6334 ± 0.0106	0.0769 ± 0.0056
GAT	0.6746 ± 0.0125	0.0451 ± 0.0035	0.0602 ± 0.0119	0.6134 ± 0.0072	0.0698 ± 0.0090
GraphSAGE	0.7713 ± 0.0157	0.0771 ± 0.0113	0.0897 ± 0.0183	0.6964 ± 0.0355	0.1247 ± 0.0258
HGT	0.7809 ± 0.0092	0.0686 ± 0.0066	0.0651 ± 0.0050	0.6744 ± 0.0182	0.1006 ± 0.0091
MTHGT (proposed)	0.7763 ± 0.0160	0.0636 ± 0.0038	0.0778 ± 0.0225	0.6427 ± 0.0902	0.0931 ± 0.0503

Modality Attribution Analysis

A key feature of MTHGT is explainability. The model dynamically learns the importance of each modality for conversion prediction. Averaged across all test predictions, the modality weights are:

- Temporal transaction sequences: 39.3%
- Call transcript dialogue embeddings: 35.9%
- Structured CRM features: 24.8%

IV. CONCLUSION AND FUTURE WORK

This paper presents **MTHGT**, a multimodal heterogeneous graph transformer framework that unifies structured, text, and temporal transaction sequences. The framework outperforms tabular-only models and homogeneous GNNs on global ranking (ROC-AUC) while providing transparent modality attributions.

Future work will target

- **Focal/Cost-Sensitive Loss:** Replacing standard BCE with Focal Loss to bridge the F1/PR-AUC gap against tabular baselines.

- **Behavioral k-NN Graph Construction:** Constructing transaction feature similarity networks to enrich node connectivity.
- **Real-World Metrics:** Transitioning evaluation metrics to Precision@K and Lift to reflect actual banking lead campaign performance.

REFERENCES

1.] S. Moro, P. Cortez, and P. Rita, "A data-driven approach to predict the success of bank telemarketing," *Decision Support Systems*, vol. 62, pp. 22–31, 2014.
2.] O. Shumovskaia et al., "Graph-based Machine Learning for Financial Networks," *International Journal of Data Science and Analytics (IJDSA)*, 2021.
3.] K. Mollaev et al., "Multimodal Banking Dataset for Financial AI," *arXiv preprint arXiv:2410.XXXXX*, 2024.
4.] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph Attention Networks," in *Proc. ICLR*, 2018.
5.] T. N. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks," in *Proc. ICLR*, 2017.
6.] W. L. Hamilton, Z. Ying, and J. Leskovec, "Inductive Representation Learning on Large Graphs," in *Proc. NeurIPS*, 2017.
7.] Z. Hu, Y. Dong, K. Wang, and Y. Sun, "Heterogeneous Graph Transformer," in *Proc. WWW*, 2020.
8.] F. M. Harper and J. A. Konstan, "The MovieLens Datasets: History and Context," *ACM Transactions on Interactive Intelligent Systems*, vol. 5, no. 4, 2015.
9.] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, and J. Leskovec, "Open Graph Benchmark: Datasets for Machine Learning on Graphs," in *Proc. NeurIPS*, 2020.
10.] E. Rossi et al., "Temporal Graph Networks for Deep Learning on Dynamic Graphs," in *Proc. NeurIPS Workshop*, 2020.
11.] L. Yu et al., "Towards Better Dynamic Graph Learning: New Architecture and Unified Library," in *Proc. NeurIPS*, 2023.
12.] S. Huang et al., "Temporal Graph Benchmark for Machine Learning on Temporal Graphs," in *Proc. NeurIPS*, 2023.
13.] L. Müller et al., "Graph Transformers: A Survey," *Computer Science Review*, 2024.
14.] W. Fey and J. E. Lenssen, "Fast Graph Representation Learning with PyTorch Geometric," in *ICLR Workshop*, 2019.

15.] X. Wang et al., “Heterogeneous Graph Attention Network,” in *Proc. WWW*, 2019.
16.] P. Battaglia et al., “Relational Inductive Biases, Deep Learning, and Graph Networks,” *arXiv preprint arXiv:1806.01261*, 2018.
17.] A. Vaswani et al., “Attention Is All You Need,” in *Proc. NeurIPS*, 2017.
18.] J. Devlin et al., “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” in *Proc. NAACL-HLT*, 2019.
19.] N. Chawla et al., “SMOTE: Synthetic Minority Over-sampling Technique,” *Journal of Artificial Intelligence Research*, 2002.
20.] T.-Y. Lin et al., “Focal Loss for Dense Object Detection,” in *Proc. ICCV*, 2017.
21.] R. Ying et al., “GNNExplainer: Generating Explanations for Graph Neural Networks,” in *Proc. NeurIPS*, 2019.
22.] S. Lundberg and S. Lee, “A Unified Approach to Interpreting Model Predictions,” in *Proc. NeurIPS*, 2017.
23.] M. Sundararajan et al., “Axiomatic Attribution for Deep Networks,” in *Proc. ICML*, 2017.
24.] Z. Hou et al., “GraphMAE: Self-Supervised Masked Graph Autoencoders,” in *Proc. KDD*, 2022.
25.] Y. You et al., “GraphCL: Contrastive Self-Supervised Learning of Graph Representations,” in *Proc. NeurIPS*, 2020.