

Generalized Graph Transformer Variational Autoencoder

Siddhant Karki

Department of Computer Science and Software Engineering
Miami University, Oxford, OH 45056, USA
karkiss@miamioh.edu

Abstract

Graph link prediction has long been a central problem in graph representation learning in both network analysis and generative modeling. Recent progress in deep learning has introduced increasingly sophisticated architectures for capturing relational dependencies within graph-structured data. In this work, we propose the Generalized Graph Transformer Variational Autoencoder (GGT-VAE). Our model integrates Generalized Graph Transformer Architecture with Variational Autoencoder framework for link prediction. Unlike prior GraphVAE, GCN, or GNN approaches, GGT-VAE leverages transformer style global self-attention mechanism along with laplacian positional encoding to model structural patterns across nodes into a latent space without relying on message passing. Experimental results on several benchmark datasets demonstrate that GGT-VAE consistently achieves above-baseline performance in terms of ROC-AUC and Average Precision. To the best of our knowledge, this is among the first studies to explore graph structure generation using a generalized graph transformer backbone in a variational framework.

1 Introduction

One of the most fundamental problems in graph machine learning is *link prediction*, the task of inferring missing or potential edges between nodes. Link prediction enables applications such as recommending new friends in social media, predicting and modeling latent associations in molecular bonds in chemistry, and road networks in remote sensing.

Traditional graph learning methods relied on hand, engineered features, matrix factorization, or random walk-based embeddings such as DeepWalk and node2vec. While effective to some extent, these methods are limited in their ability to generalize structural information beyond observed connectivity patterns. The emergence of deep learning-based models for graphs has led to significant progress in modeling node and edge representations.

Early works like the Graph Convolutional Network (GCN) extended convolutional principles to graphs by aggregating neighborhood information. Building upon this, *Variational Graph Autoencoders* (VGAEs) (Kipf & Welling, 2016) introduced a probabilistic approach to link prediction, combining graph convolutional encoders with variational inference to learn a latent distribution over graph structures. The objective can be formulated as maximizing the evidence lower bound (ELBO):

$$\mathcal{L}_{VGAE} = \mathbb{E}_{q_\phi(\mathbf{Z}|\mathbf{X},\mathbf{A})}[\log p_\theta(\mathbf{A}|\mathbf{Z})] - D_{KL}(q_\phi(\mathbf{Z}|\mathbf{X},\mathbf{A})\|p(\mathbf{Z})),$$

where \mathbf{Z} denotes the latent node embeddings parameterized by encoder parameters ϕ .

Despite their success, GCN- and GAT-based approaches remain limited by their local message passing, which makes it difficult to capture long-range relationships in graphs. In contrast, *transformer architectures* (Vaswani et al., 2023) introduce a global self-attention mechanism that directly models pairwise interactions between all elements, offering a natural way to learn global relationships. This idea led to graph-specific transformer models such as Graph-BERT (Zhang et al., 2020), Graph Transformer Networks (GTN) (Yun et al., 2020), and LPFormer (Shomer et al., 2024), which apply attention to graph-structured

data. While these models improve global context understanding, they generally work in a deterministic way and do not include a probabilistic latent space.

Recent work has started combining variational inference with transformers to enhance graph generation. For instance, GraphVAE (Simonovsky & Komodakis, 2018) and the Transformer Graph VAE for molecular generation (Nguyen & Karolak, 2025) show that combining these two frameworks can produce more expressive generative models. However, these methods mainly focus on small or domain-specific graphs and still depend on message-passing mechanisms from Graph Attention Networks.

Motivated by these developments, we explore whether a *vanilla of graph transformer* (Dwivedi & Bresson, 2021), when used within a variational autoencoder, can serve as a strong generative model for graph structure prediction. Our model, the **Generalized Graph Transformer Variational Autoencoder (GGT-VAE)**, removes explicit message passing and instead uses Laplacian positional encodings with self-attention to learn both local and global structure for link prediction. This work aims to bridge the gap between probabilistic graph modeling and transformer-inspired graph architectures.

Our main contributions are as follows:

- We introduce a **Generalized Graph Transformer Variational Autoencoder (GGT-VAE)** that combines transformer-style self-attention with variational inference for link prediction.
- Our model captures both **local and global structure** without relying on message passing.
- Experimental results on benchmark datasets show that GGT-VAE achieves competitive or superior link prediction performance compared to message-passing VAEs.

2 Related Work

2.1 Autoencoders and Variational Autoencoders

Autoencoders (AEs) are fundamental generative models that learn to reconstruct inputs by mapping them into a lower-dimensional latent space. Formally, an encoder network $f_\phi(\mathbf{x})$ projects an input \mathbf{x} to a latent representation \mathbf{z} , while a decoder $g_\theta(\mathbf{z})$ reconstructs the input as $\hat{\mathbf{x}} = g_\theta(f_\phi(\mathbf{x}))$. The objective minimizes a reconstruction loss:

$$\mathcal{L}_{AE} = \|\mathbf{x} - \hat{\mathbf{x}}\|^2.$$

Building upon this framework, Kingma & Welling (2022; 2019) introduced the *Variational Autoencoder* (VAE), which imposes a probabilistic structure on the latent space by introducing an approximate posterior $q_\phi(\mathbf{z}|\mathbf{x})$ and prior $p(\mathbf{z})$. The model optimizes the evidence lower bound (ELBO):

$$\mathcal{L}_{VAE} = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})).$$

2.2 Graph Autoencoders and Graph Neural Networks

Kipf & Welling (2016) extended the VAE framework to graph-structured data with the *Variational Graph Autoencoder* (VGAE). Here, node embeddings \mathbf{Z} are learned via a Graph Convolutional Network (GCN) encoder, which aggregates neighborhood information as:

$$\mathbf{Z} = \text{GCN}(\mathbf{X}, \mathbf{A}) = \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{X} \mathbf{W}),$$

where \mathbf{A} is the adjacency matrix and $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$. The decoder reconstructs the graph via inner products, $\hat{\mathbf{A}} = \sigma(\mathbf{Z}\mathbf{Z}^\top)$, making VGAE a natural choice for link prediction.

Subsequent work on *Graph Attention Networks* (GAT) (Veličković et al., 2018) replaced fixed convolutional aggregation with learnable attention coefficients:

$$\alpha_{ij} = \frac{\exp(\text{LeakyReLU}(\mathbf{a}^\top [\mathbf{W}\mathbf{h}_i \parallel \mathbf{W}\mathbf{h}_j]))}{\sum_{k \in \mathcal{N}(i)} \exp(\text{LeakyReLU}(\mathbf{a}^\top [\mathbf{W}\mathbf{h}_i \parallel \mathbf{W}\mathbf{h}_k]))},$$

enabling more flexible, content-dependent message passing. Despite this, both GCN and GAT architectures are inherently local, relying on neighborhood-level propagation that limits their receptive field.

2.3 Transformers for Graph Representation Learning

Transformers (Vaswani et al., 2023) introduced the concept of global self-attention, where token interactions are modeled through the scaled dot-product:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)\mathbf{V}.$$

Originally proposed for language tasks, the mechanism generalizes naturally to graphs by treating nodes as tokens. Graph-BERT (Zhang et al., 2020) applies this idea to graph data, replacing message passing with attention-based aggregation and positional encodings that capture global graph structure. Similarly, the Graph Transformer Network (GTN) (Yun et al., 2020) integrates attention over multiple edge types, while still retaining message-passing elements.

2.4 Transformer Models for Link Prediction

More recent work, such as LPFormer (Shomer et al., 2024), employs a pure transformer architecture for link prediction and relies on message passing neural network architecture (MPNNs). However, LPFormer remains a deterministic encoder-decoder that relies on MPNN for encoding. In addition, they lack the probabilistic latent structure of VAEs that supports graph generation.

2.5 Graph Generation and Molecular Design

Graph generation using VAE-based architectures has been explored by Simonovsky & Komodakis (2018), who introduced GraphVAE for generating small graphs. Later, Nguyen & Karolak (2025) extended this idea by combining transformers with GAT-VAEs for molecular generation. Their model used cross-attention between latent embeddings and a SMILES text-transformer encoder to capture both local and long-range molecular dependencies. Together, these studies show that integrating transformer architectures with variational inference can effectively support graph-level generation tasks.

These developments, from autoencoders to graph transformers, highlight a natural progression toward models that jointly capture global structural patterns and latent probabilistic representations, motivating the framework described in the next section.

In summary, autoencoders and their variational extensions provide the theoretical backbone for latent generative modeling. Transformers introduce relational reasoning through attention, while GCNs and GATs adapt this mechanism to graph-structured data. The Graph Variational Autoencoder unites these ideas, offering a framework for probabilistic graph generation. Our work builds upon this foundation by incorporating attention-based graph encoders for link prediction tasks.

3 Methods

Overall, the GGT-VAE follows a simple yet effective encoder-decoder structure. The encoder transforms node features and positional information into a compact latent space using self-attention, where each node’s representation depends on every other node in the graph. Unlike traditional GNNs, this removes the need for step-by-step neighborhood aggregation, allowing the model to learn global structure directly. The latent variables capture both local and global dependencies, while the decoder reconstructs the full graph by predicting the likelihood of edges between node pairs.

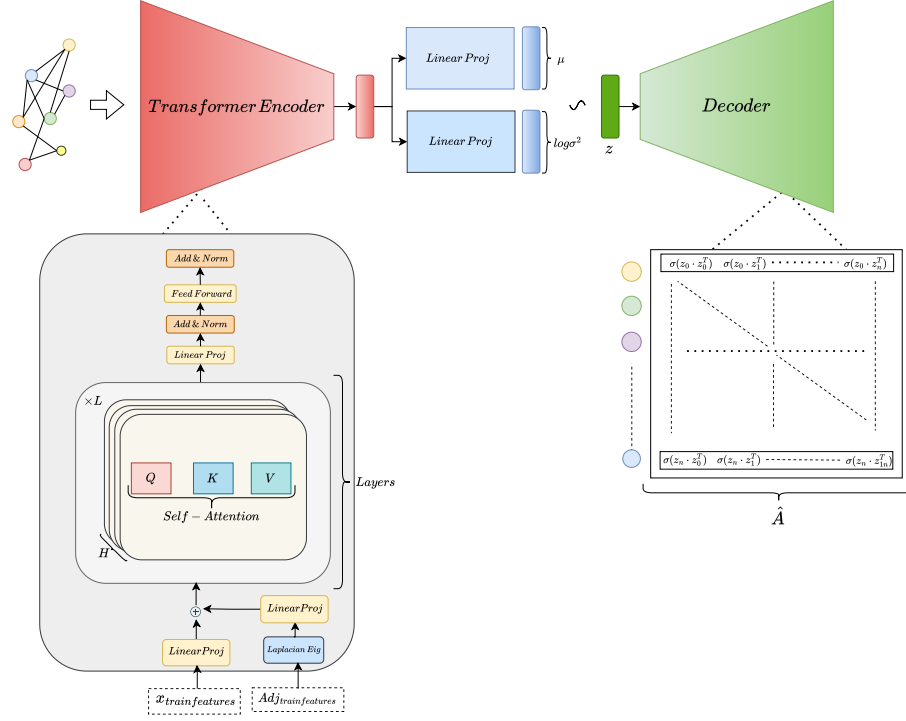


Figure 1: Architecture of the encoder–decoder framework used in our model. The encoder maps node and positional embeddings into a latent space, while the decoder reconstructs the adjacency matrix from the latent variables.

3.1 Model Architecture

We propose a **Generalized Graph Transformer Variational Autoencoder (GGT-VAE)** that models node dependencies without explicit message passing. Unlike GCN-based VAEs (Kipf & Welling, 2016), our encoder relies solely on transformer self-attention operating over laplacian position encoding produced from the training graph adjacency structure. This design allows the model to capture higher-order relationships while maintaining scalability and permutation invariance.

Input Embedding. Each node is represented by its raw feature vector $\mathbf{x}_i \in \mathbb{R}^{d_{\text{node}}}$ and positional encoding $\mathbf{p}_i \in \mathbb{R}^{d_{\text{pos}}}$, derived from the top- k Laplacian eigenvectors. The embedding module projects both components into a common latent space:

$$\mathbf{h}_i^{(0)} = \mathbf{W}_x \mathbf{x}_i + \mathbf{W}_p \mathbf{p}_i, \quad (1)$$

where $\mathbf{W}_x, \mathbf{W}_p \in \mathbb{R}^{d_{\text{hid}} \times d_{\text{in}}}$ are learnable matrices. This yields an initial node representation $\mathbf{H}^{(0)} \in \mathbb{R}^{N \times d_{\text{hid}}}$.

Graph Transformer Encoder : The encoder comprises L stacked **Graph Transformer layers**, each containing H self-attention heads. Given node representation at layer l , $\mathbf{N}^{(l-1)} \in \mathbb{R}^{N \times d_{\text{hid}}}$ each head computes query, key, and value projections:

$$\mathbf{Q}_i = \mathbf{N}_i^{l-1} \mathbf{W}_Q, \quad \mathbf{K}_i = \mathbf{N}_i^{l-1} \mathbf{W}_K, \quad \mathbf{V}_i = \mathbf{N}_i^{l-1} \mathbf{W}_V. \quad (2)$$

Attention weights are computed as scaled dot-products:

$$\mathbf{A}'_i = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right), \quad (3)$$

Outputs from all heads are concatenated and linearly projected:

$$\mathbf{N}^{l-1'} = \text{Concat}(\mathbf{N}_1, \dots, \mathbf{N}_H)\mathbf{W}_O, \quad (4)$$

followed by a residual connection, layer normalization, and a position-wise feedforward network:

$$\mathbf{N}^{(L)} = \text{LayerNorm}(\mathbf{N}^{(l-1)} + \text{FFN}(\text{LayerNorm}(\mathbf{N}^{(l-1)} + \mathbf{N}^{l-1'}))).$$

Variational Encoding : The encoder’s final output $\mathbf{N}^{(L)}$ is mapped to the parameters of a Gaussian latent distribution:

$$\boldsymbol{\mu} = \mathbf{N}^{(L)}\mathbf{W}_\mu, \quad \log \sigma^2 = \mathbf{N}^{(L)}\mathbf{W}_{\log \sigma^2}.$$

A latent sample \mathbf{Z} is obtained via the reparameterization trick:

$$\mathbf{Z} = \boldsymbol{\mu} + \boldsymbol{\epsilon} \odot \exp(0.5 \log \sigma^2), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, I).$$

Decoder : Link probabilities are reconstructed through a simple inner-product decoder:

$$\hat{\mathbf{A}} = \sigma(\mathbf{Z}\mathbf{Z}^\top),$$

where $\sigma(\cdot)$ is the elementwise sigmoid function. Self-loops are masked to avoid trivial identity edges.

Training Objective : The model is trained using the standard VAE loss:

$$\mathcal{L} = \mathcal{L}_{\text{recon}} + \beta \mathcal{L}_{\text{KL}},$$

where $\mathcal{L}_{\text{recon}}$ is the binary cross-entropy between $\hat{\mathbf{A}}$ and \mathbf{A} , and \mathcal{L}_{KL} is the Kullback–Leibler divergence between the approximate posterior $q(\mathbf{Z}|\mathbf{X})$ and the unit Gaussian prior $p(\mathbf{Z})$.

GGT-VAE achieves comparable link prediction accuracy to message-passing VGAEs, while eliminating neighborhood aggregation. The attention mechanism allows flexible context modeling, and Laplacian-based positional encodings inject graph topology directly into the latent space.

3.2 Data Collection & Training Setup

We use two standard citation network datasets: **Cora** and **Citeseer**. They are downloaded from the Planetoid repository using `torch_geometric.datasets.Planetoid`. After loading, each dataset is converted into a dense tensor format to make the model easier to run in PyTorch. This approach removes the need for sparse operations and makes the training setup simple and easy to reproduce.

For each dataset, we randomly split the edges into training, validation, and test sets. Because these graphs are very sparse, we balance the number of positive (real) and negative (false) edges during both training and evaluation. This helps the model learn fairly and prevents bias toward missing edges. During testing, we sample a smaller number of true and false edges than the maximum possible. This setup helps test the model’s ability to generalize beyond the exact edge patterns seen during training. The reconstruction loss is calculated using these balanced edge samples.

We train all models with the **AdamW optimizer**, a learning rate of 1×10^{-3} , and decay of 5×10^{-4} and full-batch updates (the entire graph per step). All experiments were run on a single GPU, and random seeds were fixed for reproducibility.

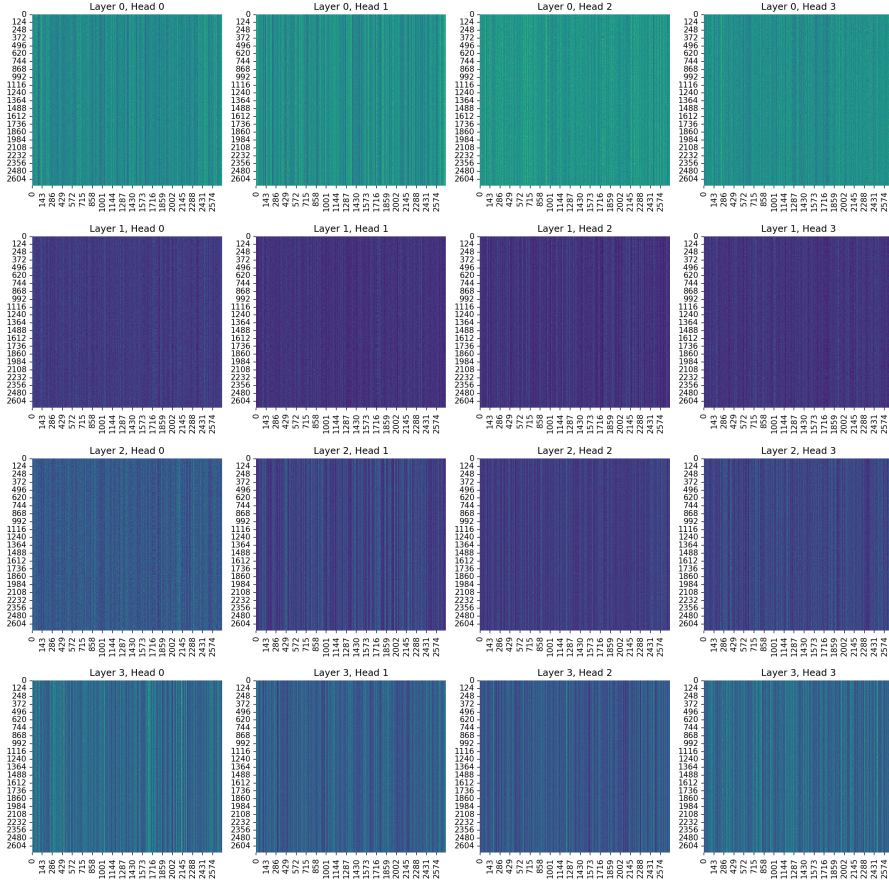


Figure 2: Attention maps for selected heads across Transformer layers. Lighter regions correspond to stronger attention weights between node pairs. Early layers focus on local neighborhoods, while deeper layers capture more global structural dependencies.

3.3 Training Configuration

We trained multiple configurations of the model to study how different parameters affect performance. The number of transformer layers, attention heads, and hidden dimensions were varied to find a balance between accuracy and model complexity. We also tested different values of the β term in the KL divergence to control the strength of regularization in the variational loss. All models were trained using the same setup described in Section 3. The best configuration for each dataset was selected based on validation performance. Results for the parameter variations are summarized in the next sections (see Table ??) for details.

4 Experimental Evaluation and Result Analysis

4.1 Qualitative Analysis

Figure 2 visualizes the attention weights across transformer layers and heads. Each map shows how much one node attends to another, with lighter areas representing stronger attention values. A clear layer-by-layer progression can be observed. In **Layer 0**, the attention maps appear mostly uniform with low contrast, meaning that the model distributes attention evenly across nodes. This suggests that the first layer focuses on capturing graph-level information from laplacian positional encodings rather than detailed node level information. In **Layer 1**, faint vertical streaks begin to appear. These patterns indicate that certain nodes

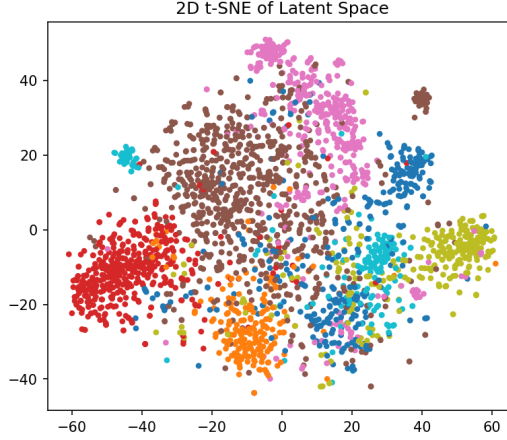


Figure 3: Cora dataset 2D t-SNE projection of latent embeddings from the encoder. Nodes belonging to similar communities are clustered together, indicating the model’s ability to capture meaningful relational structure in the latent space.

start receiving more global attention, often corresponding to central or highly connected nodes. The model begins recognizing the relative importance of different nodes in the graph. By **Layer 2**, the attention becomes more distinct, with some heads showing strong localized focus. This shows that different heads specialize: some capture fine-grained, local structures, while others maintain the global awareness. Such diversity reflects the model’s ability to learn both local and global patterns simultaneously. In **Layer 3**, the maps display a mix of concentrated and uniform regions. Some heads show sharp vertical bands, focusing on specific node clusters, while others remain evenly distributed. This final layer balances global aggregation and local refinement, integrating information from previous layers into cohesive structural representations.

Across all layers, the variation among attention heads shows that the model naturally learns to combine local and global features. This enables our model to capture structural information without relying on explicit message passing. Figure 3 presents a 2D t-SNE visualization of the latent embeddings for **Cora** dataset. Nodes (ML papers) belonging to similar classes (types) cluster together, showing that the latent space preserves meaningful relationships between nodes. This clustering indicates that the encoder successfully learns a representation where the latent space reflects structural similarity in the graph.

4.2 Quantitative Analysis

In this section, we study how learns graph structure without relying on message passing. To do this, we look at how attention spreads across different distances in the graph and how each layer’s focus changes as the model gets deeper.

Globality Metric. To understand how far each layer looks across the graph, we measure a value we call *globality*. It tells us the average distance (in graph hops) between nodes that attend to each other. For every layer ℓ and attention head h , we compute:

$$\text{Globality}_{\ell,h} = \frac{\sum_d d \cdot \bar{\alpha}^{(\ell,h)}(d)}{\sum_d \bar{\alpha}^{(\ell,h)}(d)}, \quad (5)$$

where d is the shortest-path distance (SPD) between two nodes and $\bar{\alpha}^{(\ell,h)}(d)$ is the average attention weight between nodes at that distance. A smaller value means the model mostly attends to close neighbors (local attention), while a larger value means it also considers far-away nodes (global attention).

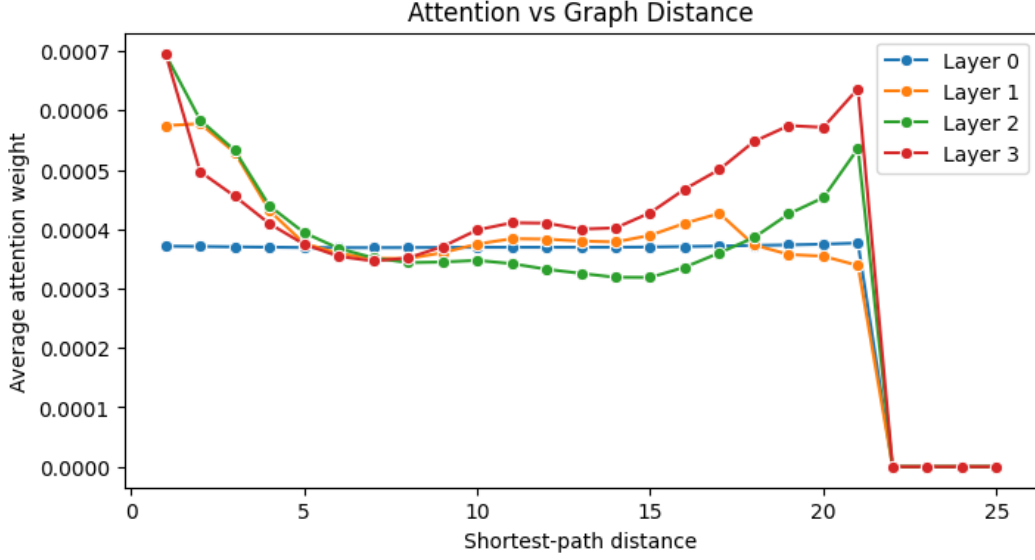


Figure 4: **Attention vs. Graph Distance.** Average attention weight for different shortest-path distances (SPD) on Cora. Unlike message-passing models that only attend to direct neighbors (SPD = 1), the Transformer assigns nonzero weight even to distant nodes (SPD > 10), showing that it captures both local and global structure without relying on adjacency-based message passing.

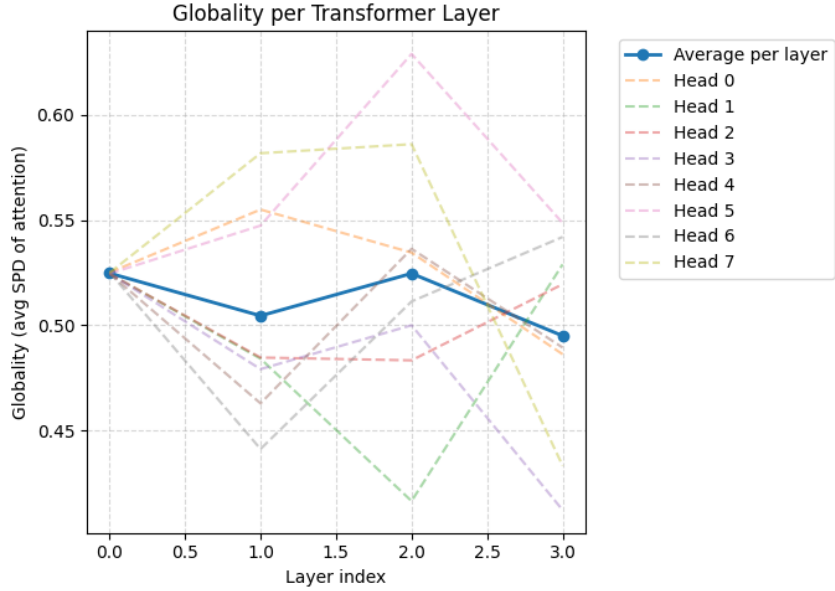


Figure 5: **Normalized globality across layers and heads.** The solid line shows the layer average; dashed lines show individual heads. Globality rises in early layers (broader attention) and decreases in the final layer (local refinement). Diverging head trends show that some heads specialize in global structure while others focus on local details.

We also divide this score by the graph’s diameter to keep the values consistent between 0 and 1:

$$\text{NormalizedGlobality}_{\ell,h} = \frac{\text{Globality}_{\ell,h}}{D_{\text{graph}}}. \quad (6)$$

Table 1: Link prediction on Planetoid datasets. Results are ROC-AUC / AP (%). Higher is better.

Method	Cora	Citeseer
DeepWalk Zhu et al. (2022)	76.41 / 81.78	64.15 / 74.60
Spectral Clustering Zhu et al. (2022)	85.36 / 88.41	77.26 / 80.24
Graphite-AE Zhu et al. (2022)	91.29 / 92.63	88.84 / 89.64
Linear-GAE Zhu et al. (2022)	91.86 / 93.20	90.89 / 92.24
2-GAE Zhu et al. (2022)	91.61 / 92.81	89.04 / 89.26
6-GAE Zhu et al. (2022)	83.58 / 85.32	79.95 / 83.63
6-DGAE _a Zhu et al. (2022)	93.55 / 94.46	94.16 / 94.86
ARGA Zhu et al. (2022)	92.10 / 93.25	90.43 / 92.04
GGT-VAE (ours)	92.04 ± 0.60 / 92.66 ± 0.80	92.00 ± 0.20 / 93.74 ± 0.16

Results reproduced from [Zhu et al. \(2022\)](#). Our transformer-based VAE matches early VGAE-style baselines without message passing.

Attention–Distance Analysis. Figure 4 shows the average attention weight for different graph distances on the Cora dataset. In a message-passing network, information only moves between direct neighbors ($SPD = 1$). In contrast, GGT-VAE gives meaningful attention to nodes that are many hops apart ($SPD > 10$). The first layer shows almost uniform attention across all distances, meaning it starts with a global view. Middle layers balance both local (1–2 hop) and long-range connections, while the deepest layer again highlights both very near and very distant nodes. This pattern shows that the model learns to mix local and global information on its own, without needing any explicit neighborhood-based message passing.

Globality Across Layers and Heads. To understand how each layer and head contribute to the model’s receptive field, we measure the *globality* of attention, the average distance between nodes that attend to each other, normalized by the graph’s diameter. Figure 5 shows these values for all eight attention heads and the layer-wise average.

On average, globality increases slightly from Layer 0 to Layer 2, meaning the model initially expands its view to include more distant nodes. In the final layer, the average drops again, showing that attention becomes somewhat more local once global information has been collected. This pattern suggests a natural progression: early layers mix information broadly, while later layers refine relationships within smaller neighborhoods.

Across individual heads, we observe diverging trends. Some heads consistently show higher globality, focusing on long-range structure and global context. Others decrease sharply, concentrating on nearby nodes and preserving local detail. This diversity across heads indicates that different parts of the Transformer specialize in complementary roles— some capturing the big picture of the graph, others refining fine-grained connections. Together, they enable the model to balance both local and global understanding without relying on message passing.

This confirms that the GGT-VAE learns to represent both local and global graph patterns without *explicit message passing*. The model builds a latent space that reflects graph connectivity directly through attention, rather than by repeatedly aggregating features over edges. This explains why it performs well on link prediction: it can connect distant but structurally related nodes while keeping local information intact.

4.3 Experimental Setup

We evaluate the proposed GGT-VAE on the standard Planetoid datasets **Cora** and **Citeseer** for the **link prediction** task. We follow the data splits and protocol of ? for comparability. Node features are used as-is, and Laplacian positional encodings supply structural information. The adjacency matrix for training is constructed only from training edges to prevent leakage. We train for 500 epochs with Adam ($\text{lr} = 1 \times 10^{-3}$), and apply early stopping with patience 50 on validation ROC-AUC. Unless noted, we use 4 transformer layers, 4 heads,

Table 2: Ablations on **Cora**: effect of heads, layers, and hidden size. β and LR are shown as β /LR. Results are ROC-AUC / AP (%).

Config	Heads	Layers	Hidden	β /LR	Cora
Base	4	4	128	$0.5 \times 10^{-3} / 1 \times 10^{-3}$	92.04 / 92.66
<i>Vary Heads (Layers=4, Dim=128)</i>					
H-1	1	4	128	$0.5 \times 10^{-3} / 1 \times 10^{-3}$	90.89 / 92.44
H-2	2	4	128	$0.5 \times 10^{-3} / 1 \times 10^{-3}$	90.66 / 92.12
H-8	8	4	128	$0.5 \times 10^{-3} / 1 \times 10^{-3}$	91.03 / 92.49
<i>Vary Layers (Heads=4, Dim=128)</i>					
L-2	4	2	128	$0.5 \times 10^{-3} / 1 \times 10^{-3}$	90.06 / 92.03
L-8	4	8	128	$0.5 \times 10^{-3} / 1 \times 10^{-5}$	48.88 / 49.94
<i>Vary Hidden Dim (Heads=4, Layers=4)</i>					
D-064	4	4	64	$0.5 \times 10^{-3} / 1 \times 10^{-3}$	89.73 / 91.49
D-256	4	4	256	$0.5 \times 10^{-3} / 1 \times 10^{-5}$	92.23 / 93.81
D-512	4	4	512	$0.5 \times 10^{-3} / 1 \times 10^{-6}$	77.32 / 70.81

hidden size 128, and $\beta = 0.5 \times 10^{-3}$. All numbers are averaged over 10 random seeds and reported as mean \pm std.

4.4 Model Performance

Table 1 shows that classic embedding methods (DeepWalk, Spectral Clustering) underperform because they lack a generative latent space. Message-passing autoencoders (Graphite-AE, Linear-GAE, ARGAE) perform strongly by aggregating neighborhood information. Our **GGT-VAE** reaches 92.04% ROC-AUC / 92.66% AP on Cora and 92.00% ROC-AUC / 93.74% AP on Citeseer, demonstrating that transformer self-attention with structural encodings matches message-passing VAEs even without adjacency-based propagation.

Ablation studies. Table 2 studies depth, heads, hidden size, and β . Four layers and four heads give the best stability. Fewer heads slightly hurt performance; more (8) does not yield further gains. Increasing hidden size to 256 helps, but very large hidden sizes (512) destabilize training. A moderate KL weight (0.5×10^{-3}) balances reconstruction and regularization. Overall, GGT-VAE is robust to moderate hyperparameter changes.

Summary. GGT-VAE attains high link-prediction accuracy on Cora and Citeseer using only self-attention and positional encodings. During testing we evaluate on a subsampled set of positive and negative edges to check generalization; the model maintains strong performance, indicating it does not overfit the training graph. This supports our claim that attention + variational inference is sufficient for graph representation learning without explicit message passing.

5 Conclusion & Future Work

This paper introduced the **Generalized Graph Transformer Variational Autoencoder (GGT-VAE)** for link prediction. The model replaces traditional message passing with transformer self-attention and Laplacian positional encodings, allowing it to learn both local and global structure directly from the graph. Experiments on benchmark datasets show that GGT-VAE performs competitively with message-passing approaches while maintaining a simpler and more flexible architecture. Our analysis further shows that the model learns structural information naturally across layers, confirming that self-attention alone can effectively capture graph relationships.

Currently, the model focuses on edge generation from latent representations. Future work will extend this approach to full graph generation for a fixed number of nodes, incorporating

node features as inputs. We plan to explore applications in molecular generation, where the model could learn complex chemical bonds and topologies, and in geospatial domains, such as generating road networks from satellite imagery. Another promising direction is conditioning graph generation on other modalities, such as images through cross-attention with vision transformers, or molecular SMILES text for chemistry tasks. These directions highlight the potential of combining graph transformers and variational autoencoders for broader multimodal and scientific applications.

Code Availability. To support transparency and reproducibility, we will release the full implementation, training scripts, and pre-trained model weights of GGT-VAE upon publication of the final version of this paper. The codebase will include utilities for dataset preparation, model configuration, and visualization of attention and globality metrics to facilitate future research on transformer-based graph generative models.

References

- Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to graphs, 2021. URL <https://arxiv.org/abs/2012.09699>.
- Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019. ISSN 1935-8245. doi: 10.1561/22000000056. URL <http://dx.doi.org/10.1561/22000000056>.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2022. URL <https://arxiv.org/abs/1312.6114>.
- Thomas N. Kipf and Max Welling. Variational graph auto-encoders, 2016. URL <https://arxiv.org/abs/1611.07308>.
- Trieu Nguyen and Aleksandra Karolak. Transformer graph variational autoencoder for generative molecular design. *Biophysical Journal*, 2025. ISSN 0006-3495. doi: <https://doi.org/10.1016/j.bpj.2025.01.022>. URL <https://www.sciencedirect.com/science/article/pii/S0006349525000359>.
- Harry Shomer, Yao Ma, Haitao Mao, Juanhui Li, Bo Wu, and Jiliang Tang. Lpformer: An adaptive graph transformer for link prediction. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '24, pp. 2686–2698. ACM, August 2024. doi: 10.1145/3637528.3672025. URL <http://dx.doi.org/10.1145/3637528.3672025>.
- Martin Simonovsky and Nikos Komodakis. Graphvae: Towards generation of small graphs using variational autoencoders, 2018. URL <https://arxiv.org/abs/1802.03480>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL <https://arxiv.org/abs/1706.03762>.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks, 2018. URL <https://arxiv.org/abs/1710.10903>.
- Seongjun Yun, Minbyul Jeong, Raehyun Kim, Jaewoo Kang, and Hyunwoo J. Kim. Graph transformer networks, 2020. URL <https://arxiv.org/abs/1911.06455>.
- Jiawei Zhang, Haopeng Zhang, Congying Xia, and Li Sun. Graph-bert: Only attention is needed for learning graph representations, 2020. URL <https://arxiv.org/abs/2001.05140>.
- Fuyang Zhu, Yang Zheng, Jingtao Hu, Chuanhao Li, Hao Xue, and Zhipeng Yu. Stabilizing and enhancing link prediction in graph autoencoders. *Frontiers in Big Data*, 5:9754798, 2022. doi: 10.3389/fdata.2022.9754798. URL <https://pmc.ncbi.nlm.nih.gov/articles/PMC9754798/>.