

## Definition: Graph Out-of-Distribution (GraphOOD)

Graph Out-of-Distribution (GraphOOD) refers to the problem of handling graph data where the testing or deployment graph instances deviate from the training graph distribution. In many real-world applications (e.g., molecular chemistry, social networks), graph data is highly structured and dynamic, and the assumption that test data follows the same distribution as training data often does not hold. This problem is important because graph-based models trained on a specific distribution can generalize poorly to unseen or shifted data, leading to significant performance degradation.

## Current Research Challenges

- Distributional Shifts:** GraphOOD faces the challenge of handling different types of shifts, including:
  - Node distribution shifts:** Differences in the distribution of node features between training and test graphs.
  - Edge distribution shifts:** Changes in the relationships or edge connections between nodes.
  - Structural distribution shifts:** Variations in the global topology of the graph.
- Model Robustness:** Traditional graph models like Graph Neural Networks (GNNs) are often brittle to distributional shifts. Designing models that are robust to these shifts while maintaining high performance on in-distribution data is a key challenge.
- Evaluation Metrics:** Evaluating model performance on GraphOOD problems is difficult because existing metrics like accuracy or F1-score may not reflect a model's robustness to out-of-distribution data. New evaluation metrics are required to assess how well a model generalizes across different graph domains.
- Scalability:** Graphs in real-world applications are often large and complex. Ensuring the methods used for GraphOOD are scalable to large graphs is a critical challenge.

---

## Commonly Used Research Methods in GraphOOD

### 1. Domain Adaptation

**Method:** The goal of domain adaptation is to reduce the discrepancy between the training (source) domain and the testing (target) domain by aligning their feature distributions. This can be done using techniques like **Maximum Mean Discrepancy (MMD)** or **adversarial domain adaptation**.

**Formula:**

$$\text{MMD}(P, Q) = \|\mathbb{E}_{x \sim P}[\phi(x)] - \mathbb{E}_{y \sim Q}[\phi(y)]\|_2^2$$

where  $P$  and  $Q$  are the distributions of source and target domains, and  $\phi(x)$  is a feature mapping function.

**Example:** Consider a chemical compound graph (training) and a biological interaction graph (testing). By minimizing the MMD between the distributions of node embeddings, the model adapts to the new domain (biological graph) without retraining on new data.

### 2. Invariant Risk Minimization (IRM)

**Method:** IRM is designed to learn invariant predictors across multiple environments or distributions. The idea is to find a representation such that the optimal classifier for each environment remains unchanged.

**Formula:**

$$\min_{\phi} \sum_{e \in \mathcal{E}} R^e(w \circ \phi) \quad \text{subject to} \quad w \in \arg \min_w R^e(w \circ \phi) \quad \forall e \in \mathcal{E}$$

where  $\mathcal{E}$  is the set of environments,  $R^e$  is the risk in environment  $e$ ,  $\phi$  is the feature representation, and  $w$  is the classifier.

**Example:** In a graph-based recommendation system, the training data consists of user-item interactions from multiple regions (e.g., different countries). By applying IRM, the model learns to generalize across regions, ensuring that the item recommendations remain consistent even if a user comes from an unseen region.

### 3. Graph Augmentation

**Method:** Graph augmentation methods aim to generate diverse variants of the training graphs to mimic the shifts that can occur in out-of-distribution data. This includes adding noise to edges, modifying node features, or changing the graph's structure.

**Formula:**

$$G' = G + \mathcal{N}(G)$$

where  $G'$  is the augmented graph and  $\mathcal{N}(G)$  represents the noise or transformation applied to the original graph  $G$ .

**Example:** In a citation network graph, where nodes represent papers and edges represent citations, the graph can be augmented by randomly removing edges (representing missing citations) or adding new edges (representing unseen citations) to simulate potential distribution shifts.

### 4. Out-of-Distribution Detection

**Method:** Out-of-distribution detection methods are designed to identify whether a given test graph comes from a different distribution than the training graphs. This can be achieved by estimating the likelihood of the test graph under the training distribution or using graph-based anomaly detection methods.

**Formula:**

$$\mathcal{L}(G) = -\log P(G|\theta)$$

where  $\mathcal{L}(G)$  is the log-likelihood of graph  $G$ , and  $\theta$  represents the parameters of the model trained on the training distribution.

**Example:** In a social network graph, where nodes represent users and edges represent friendships, an out-of-distribution detection model can be used to identify users who exhibit unusual interaction patterns (e.g., isolated clusters or rapid growth in the number of connections), indicating a potential deviation from the normal distribution of interactions.

### 5. Graph Contrastive Learning

**Method:** Graph contrastive learning focuses on learning robust graph representations by maximizing the mutual information between different views (e.g., augmentations) of the same graph. It helps improve generalization on out-of-distribution graphs by learning invariant representations.

**Formula:**

$$\mathcal{L}_{\text{contrastive}} = -\mathbb{E} \left[ \log \frac{\exp(\text{sim}(z_i, z_j))}{\sum_k \exp(\text{sim}(z_i, z_k))} \right]$$

where  $z_i$  and  $z_j$  are the embeddings of two augmented views of the same graph, and  $\text{sim}(z_i, z_j)$  is the similarity between them (e.g., cosine similarity).

**Example:** In a protein-protein interaction graph, different augmentations of the graph (e.g., removing edges or perturbing node features) can be used to train a contrastive model that learns robust protein embeddings, ensuring generalization to new unseen proteins in a different distribution (e.g., organisms with different evolutionary histories).

These methods address different aspects of the GraphOOD problem and illustrate a variety of approaches to improving generalization and robustness when dealing with out-of-distribution graphs.