# Introduction: GNNs in Recommendation Systems

Graph Neural Networks (GNNs) have gained significant traction in recommendation systems due to their ability to model complex relationships in user-item interaction graphs. While GNNs are widely used in various fields such as natural language processing (NLP) and computer vision, their application in recommendation systems comes with unique challenges. In contrast to typical tasks where GNNs are applied (e.g., node classification or graph-level predictions), recommendation systems involve heterogeneous, dynamic, and large-scale graphs. As a result, significant adaptations and structural improvements are required for GNNs to efficiently handle sparse and dynamic interaction data.

Key differences include:

1. **Heterogeneous Graphs**: Unlike homogeneous graphs in social networks or citation graphs, recommendation graphs consist of distinct types of nodes (users and items) with different interaction semantics.

2. **Scalability**: GNNs need to be scalable enough to handle millions of users and items, which is more challenging compared to smaller, dense graphs.

3. **Cold Start Problem**: GNN-based recommendation systems need to be designed to effectively handle new users/items with little interaction data.

To adapt GNNs for recommendation systems, researchers have focused on structural improvements like using bipartite graphs, designing efficient message-passing mechanisms, and integrating side information (e.g., user profiles or item attributes).

## Common Paradigms in GNN-based Recommendation Systems

1. **Bipartite Graphs**: Recommendation systems often use a bipartite graph structure, where one set of nodes represents users, and the other set represents items. The edges between them represent interactions (e.g., clicks, ratings). This setting naturally fits the GNN framework, as message passing can propagate signals between users and items through these interactions.

2. **Session-based Recommendations**: In session-based recommendation tasks, the interactions are sequential, requiring GNNs to consider the temporal order of interactions. Temporal graphs or dynamic GNNs are often employed in this case.

3. **Social-based Recommendations**: Social relationships (e.g., friendships) between users can be modeled as additional edges in the graph, allowing GNNs to aggregate information not only from past user-item interactions but also from the user's social network.

## Current Research Challenges

1. **Scalability**: Handling large-scale datasets with millions of nodes and edges is a primary challenge. Techniques like graph sampling, clustering, and hierarchical GNNs have been introduced to reduce computational overhead.

2. **Dynamic Graphs**: User-item interactions are dynamic in nature, with new interactions happening continuously. Static GNN models struggle with this, leading to the development of dynamic GNN architectures, which update the learned embeddings as new data arrives.

3. **Cold Start Problem**: Handling users or items with few interactions is a well-known challenge. Hybrid models that incorporate side information (e.g., user/item attributes or external metadata) have been explored to alleviate this issue.

4. **Long-tail Recommendation**: The interaction data is usually skewed, where a few popular items receive most of the interactions. GNNs need to be adapted to avoid overfitting to these popular items while still providing relevant recommendations for less popular ones.

## Commonly Used Research Methods & GNN Structures (Formulas)

1. **Message Passing Mechanism**:

   The general message-passing process in GNNs involves two key steps:

   ○ **Aggregation**: The message passed to node $i$ from its neighbors $\mathcal{N}(i)$ is aggregated:

   $$m_i^{(l)} = \text{AGGREGATE}^{(l)} \left( \{h_j^{(l-1)} : j \in \mathcal{N}(i)\} \right)$$

   where $h_j^{(l-1)}$ represents the embedding of node $j$ at layer $l - 1$, and $\mathcal{N}(i)$ denotes the set of neighbors of node $i$.

   ○ **Update**: The node's embedding is updated based on the aggregated message:

   $$h_i^{(l)} = \text{UPDATE}^{(l)} \left( h_i^{(l-1)}, m_i^{(l)} \right)$$

   where $h_i^{(l)}$ is the updated embedding at layer $l$.

2. **Bipartite Graph Neural Network**:

   For recommendation systems, a bipartite graph GNN typically follows:

   $$h_u^{(l)} = \sigma \left( W_u^{(l)} \cdot \text{AGGREGATE}^{(l)} \left( \{h_i^{(l-1)} : i \in \mathcal{N}(u)\} \right) \right)$$

   $$h_i^{(l)} = \sigma \left( W_i^{(l)} \cdot \text{AGGREGATE}^{(l)} \left( \{h_u^{(l-1)} : u \in \mathcal{N}(i)\} \right) \right)$$

   where $h_u^{(l)}$ and $h_i^{(l)}$ are user and item embeddings at layer $l$, respectively, and $W_u^{(l)}$ and $W_i^{(l)}$ are layer-specific weight matrices.

3. **Graph Attention Networks (GAT) in Recommendations**:

   GATs are often used to weigh the importance of different neighbors during aggregation:

   $$e_{ij}^{(l)} = \text{LeakyReLU} \left( a^\top [W^{(l)} h_i^{(l-1)} || W^{(l)} h_j^{(l-1)}] \right)$$

   $$\alpha_{ij}^{(l)} = \frac{\exp(e_{ij}^{(l)})}{\sum_{k \in \mathcal{N}(i)} \exp(e_{ik}^{(l)})}$$

   The final aggregation is computed as:

   $$h_i^{(l)} = \sigma \left( \sum_{j \in \mathcal{N}(i)} \alpha_{ij}^{(l)} W^{(l)} h_j^{(l-1)} \right)$$

   where $\alpha_{ij}^{(l)}$ denotes the attention coefficient.

4. **Hybrid Models**:

   Hybrid models incorporate both collaborative filtering (CF) signals and content-based features:

   $$h_u^{CF} = \text{CF}_u \quad \text{and} \quad h_i^{CF} = \text{CF}_i$$

   The final user-item interaction score is a combination of collaborative and content signals:

   $$\hat{y}_{ui} = f \left( h_u^{CF}, h_i^{CF}, h_u^{content}, h_i^{content} \right)$$

   where $f$ could be a simple dot product or a more complex fusion mechanism.

# Conclusion

The use of GNNs in recommendation systems requires careful adaptation to handle the unique challenges of large, dynamic, and heterogeneous user-item graphs. Structural modifications such as bipartite graph modeling, dynamic graph updates, and attention mechanisms have proven effective in making GNNs more suitable for the recommendation domain. Despite progress, challenges such as scalability, the cold start problem, and long-tail recommendations remain active areas of research.