

A Node-Aware Dynamic Quantization Approach for Graph Collaborative Filtering

Lin Li

Wuhan University of Technology
Wuhan, China
cathylilin@whut.edu.cn

Chunyang Li

Wuhan University of Technology
Wuhan, China
l-cy@whut.edu.cn

Yu Yin

Wuhan University of Technology
Wuhan, China
Huawei Technologies Co., Ltd
Shanghai, China
mktb@whut.edu.cn

Xiaohui Tao

University of Southern Queensland
Springfield, Australia
Xiaohui.Tao@unisq.edu.au

Jianwei Zhang

Iwate University
Morioka, Japan
zhang@iwate-u.ac.jp

ABSTRACT

In the realm of collaborative filtering recommendation systems, Graph Neural Networks (GNNs) have demonstrated remarkable performance but face significant challenges in deployment on resource-constrained edge devices due to their high embedding parameter requirements and computational costs. Using common quantization method directly on node embeddings may overlooks their graph based structure, causing error accumulation during message passing and degrading the quality of quantized embeddings. To address this, we propose Graph based Node-Aware Dynamic Quantization training for collaborative filtering (GNAQ), a novel quantization approach that leverages graph structural information to enhance the balance between efficiency and accuracy of GNNs for Top-K recommendation. GNAQ introduces a node-aware dynamic quantization strategy that adapts quantization scales to individual node embeddings by incorporating graph interaction relationships. Specifically, it initializes quantization intervals based on node-wise feature distributions and dynamically refines them through message passing in GNN layers. This approach mitigates information loss caused by fixed quantization scales and captures hierarchical semantic features in user-item interaction graphs. Additionally, GNAQ employs graph relation-aware gradient estimation to replace traditional straight-through estimators, ensuring more accurate gradient propagation during training. Extensive experiments on four real-world datasets demonstrate that GNAQ outperforms state-of-the-art quantization methods, including BiGeaR and N2UQ, by achieving average improvement in 27.8% Recall@10 and 17.6% NDCG@10 under 2-bit quantization. In particular, GNAQ is capable of maintaining the performance of full-precision models while reducing their model sizes by 8 to 12 times; in addition, the training time is twice as fast compared to quantization baseline methods.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CIKM '25, November 10–14, 2025, Seoul, Republic of Korea

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-2040-6/2025/11...\$15.00
<https://doi.org/10.1145/3746252.3761244>

CCS CONCEPTS

• Information systems → Recommender systems.

KEYWORDS

Graph Collaborative Filtering; Quantization-aware Training; Dynamic Quantization

ACM Reference Format:

Lin Li, Chunyang Li, Yu Yin, Xiaohui Tao, and Jianwei Zhang. 2025. A Node-Aware Dynamic Quantization Approach for Graph Collaborative Filtering. In *Proceedings of the 34th ACM International Conference on Information and Knowledge Management (CIKM '25)*, November 10–14, 2025, Seoul, Republic of Korea. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3746252.3761244>

1 INTRODUCTION

The rapid proliferation of online services and social platforms has led to an explosion of user-item interaction data, making recommender systems indispensable for personalized information filtering [31]. Collaborative filtering (CF) models, especially those based on graph neural networks (GNNs) [11, 15, 16, 28, 32, 35], have achieved remarkable success in capturing complex user-item relationships and improving the precision of the recommendation. However, the increasing scale of GNN-based recommender systems poses significant challenges in terms of computational cost and memory consumption, especially when deploying models on resource-constrained devices.

Model quantization [9] has emerged as a promising solution to reduce the storage and computation requirements of deep learning models by representing parameters and activations with low-precision values. Although quantization-aware training (QAT) [8] and post-training quantization (PTQ) [7, 19, 26, 34] have been extensively studied in computer vision and natural language processing, their application to GNNs for collaborative filtering remains under-explored [2, 6]. The main object of quantization of graph neural networks in recommendation systems is node encoding. As shown in Figure 1, unlike CNNs, RNNs, or Transformers [27], where the parameters are concentrated mainly in linear, normalization, or attention layers, GNNs for recommendation tasks primarily store the parameters in the embeddings of nodes [11, 28]. During training, node embeddings are optimized via message passing on the

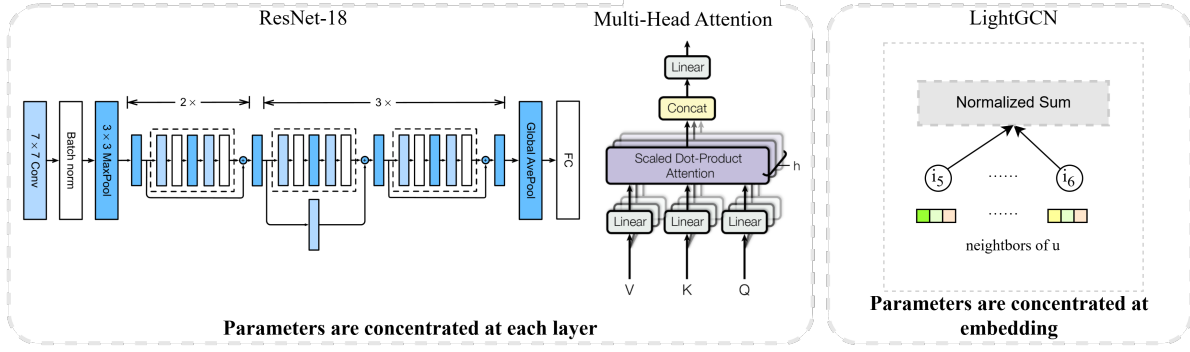


Figure 1: Embeddings as Parameters of LightCNN in RS.

graph [13, 33]. The final node embeddings serve as representations for users and items, and user-item relevance is computed via dot product or fully connected layers. Thus, quantization in this context essentially refers to quantizing node embeddings [2, 6].

The distribution of node embeddings is difficult to capture due to their role in fitting interaction relationships [14, 24]. Directly applying quantization methods from other domains to GNNs in recommendation tasks ignores the unique structure and can lead to error accumulation during message passing, ultimately degrading the quality of quantized embeddings. Moreover, using straight-through estimators [1] for gradient approximation may result in parameter updates that are inconsistent with the quantization objective, slowing convergence. Existing QAT methods for GNNs often treat all nodes and parameters uniformly [20], neglecting the heterogeneity and structural information of graph data. This can lead to suboptimal quantization, as the various roles and connectivity patterns of nodes are not adequately considered [37].

To address these limitations, we propose a novel Graph Node-Aware Quantization with dynamic step size (GNAQ) approach for GNN-based collaborative filtering. Our approach introduces three key innovations: (1) node-aware parameter sharing quantization [18], which adapts quantization parameters to the local structure of each node; (2) dynamic step size quantization [5, 26], which learns quantization intervals during training to better fit the distribution of node embeddings; and (3) graph interaction-aware quantization updates, which leverages neighborhood information to refine quantized representations. By integrating these components, GNAQ enables precise and adaptive quantization, leading to improved recommendation performance and efficiency.

Extensive experiments are conducted on four widely used public datasets, demonstrating that our GNAQ consistently outperforms state-of-the-art quantization methods in terms of Recall and NDCG metrics, while maintaining low computational overhead. Our contributions can be summarized as follows. 1) We identify the limitations of existing GNN quantization methods in collaborative filtering and highlight the importance of structure-aware quantization. 2) We propose GNAQ, a novel quantization framework that incorporates node-aware parameter sharing, dynamic step size learning, and

graph interaction-aware updates. 3) The comprehensive experimental results show the effectiveness and efficiency of GNAQ in multiple benchmark datasets.

2 RELATED WORK

To establish comprehensive technical context for our work, we first review fundamental components spanning three critical dimensions: graph-based collaborative filtering architectures that form our methodological foundation, general model quantization techniques that provide theoretical underpinnings for parameter compression, and specialized quantization approaches for graph neural networks that reveal existing limitations in structural awareness. This tripartite analysis will systematically position our work within the current research landscape.

2.1 Graph Collaborative Filtering

Recent advances in graph neural networks have established Graph Collaborative Filtering (GCF) as a dominant paradigm for modeling user-item interactions. Early works like NGCF [26] pioneered the construction of bipartite interaction graphs to capture high-order collaborative signals through multi-layer message passing. Subsequent innovations such as LightGCN [11] demonstrated that simplified architectures with feature transformation removal could achieve superior performance, highlighting the critical role of graph structural learning in recommendation tasks. The emergence of graph contrastive learning techniques further enhanced GCF’s capability to handle data sparsity, with methods like SGL [30] and HCCF [22] leveraging self-supervised learning to refine node representations through augmented graph views.

Despite these advancements, conventional GCF architectures face inherent scalability limitations due to their reliance on full-precision node embeddings. The storage complexity grows linearly with the number of entities, requiring gigabyte-level memory for million-scale user/item scenarios. This creates significant deployment barriers for edge devices while neglecting the latent structural patterns that could guide efficient parameter allocation. Existing solutions predominantly focus on architectural modifications rather than fundamentally addressing the precision-redundancy in node representations, leaving the potential of structure-aware quantization largely unexplored.

Table 1: Symbol Summary.

Notation	Description	Notation	Description
U	Set of users	q_i	Low-precision embedding representation of entity i
I	Set of items	a_i	Scaling factor for the quantization interval
N	Total number of entities	e'_i	Extended quantized embedding representation
X	User-item interaction matrix	E	Full-precision embedding table
P	Full-precision GNN model with floating-point parameters.	E_{quant}	Low-precision embedding table after quantization
Q	Quantized GNN model with low-precision parameters.	E'	Extended quantization embedding table
e_i	Full-precision embedding representation of entity i	E'_{quant}	Embedding table in the forward propagation process

2.2 Model Quantization

Model quantization aims to reduce the memory footprint and computational cost of deep learning models by representing weights and activations with low-precision values. Quantization-aware training (QAT) and post-training quantization (PTQ) are two mainstream approaches. QAT simulates quantization effects during training, allowing the model to adapt to low-precision constraints, while PTQ applies quantization after model training, often with minimal or no retraining. Although quantization has been extensively studied in computer vision [4, 12] and NLP [36], its application to GNNs, especially for recommender systems, remains limited.

2.3 Quantization for GNNs

Recent work has explored quantization techniques for GNNs to enable efficient inference on large-scale graphs. N2UQ introduces non-uniform-to-uniform quantization for GNNs, while BiGeaR proposes a binarized graph representation learning method with multi-faceted quantization reinforcement. However, most existing methods treat all nodes and parameters uniformly, ignoring the heterogeneity of graph structures. Degree-Quant and SGQuant [6] attempt to address this by considering node degrees or specialized quantization, but still lack fine-grained adaptation to node-specific characteristics.

Despite recent progress, there remains a gap in structure-aware quantization methods that can dynamically adapt quantization parameters to the local graph structure with respect to node embeddings. Our approach addresses this gap by proposing a node-aware, dynamically adaptive quantization framework for GNN-based collaborative filtering.

3 PROBLEM DESCRIPTION

The problem description in the paper is as follows. We define $U = \{u_1, u_2, \dots, u_n\}$ and $I = \{i_1, i_2, \dots, i_m\}$ as the sets of users and items, respectively, where n and m denote the number of users and items. The total number of entities is $N = n + m$. Based on the historical interaction data between users and items, we define the interaction matrix $X = \{x_{ui} | u \in U, i \in I\}$, where $x_{ui} = 1$ indicates that user u has interacted with item i , and $x_{ui} = 0$ otherwise. P and Q represent the full-precision graph neural network model and the quantized graph neural network model, respectively.

Given the above definitions, the quantization task of graph neural network models in the context of collaborative filtering recommendation can be formally described as follows:

Input: A graph neural network model P with full-precision floating-point parameters; User set U ; Item set I ; User-item interaction matrix X .

Output: A quantized graph neural network model Q with compressed parameters.

To ensure consistency in notation, Table 1 lists the definitions of high-frequency symbols used in this paper.

4 GNAQ METHODOLOGY

Considering the structural characteristics of GNNs and the specifics of recommendation scenarios helps to improve the quantization effect of GNN-based recommendation models. Our approach Graph Node-Aware Quantization-aware training (GNAQ) with dynamic step size and asymmetric quantization for collaborative graph filtering recommendation. As shown in Figure 2, this approach is based on quantization-aware training and parameter sharing, transforming the search for quantization parameters into learning of the parameters of the shared model. The optimization of quantization step size is incorporated into the objective of quantized model parameter tuning, enabling effective dynamic adjustment of the step size. Furthermore, a backpropagation method for quantization functions based on graph interaction aggregation is designed, fully considering the mechanism of message passing of GNNs. Finally, to address Top N rank recommendation scenarios, a ranking loss is added to the quantization training loss function to enhance the ranking ability of the quantized model. In summary, the main features of our proposed GNAQ are as follows.

Parameter sharing for quantization-aware training: GNAQ leverages the rich node association information in GNNs and implements quantization-aware training based on parameter sharing.

Dynamic step size update via message passing: The quantization step size is concatenated with the embeddings of the nodes, and dynamic updates are achieved through the GNN message passing mechanism. The quantized encoding is updated via graph interaction aggregation, effectively updating quantized values.

Ranking loss in quantization training: The loss function includes a ranking loss, which effectively improves the ranking ability of the quantized model.

4.1 Quantization Based on Parameter Sharing

Parameter sharing is a method to reduce the total number of parameters by reusing the same parameters in different parts of a neural network model. This method defines an allocation matrix to

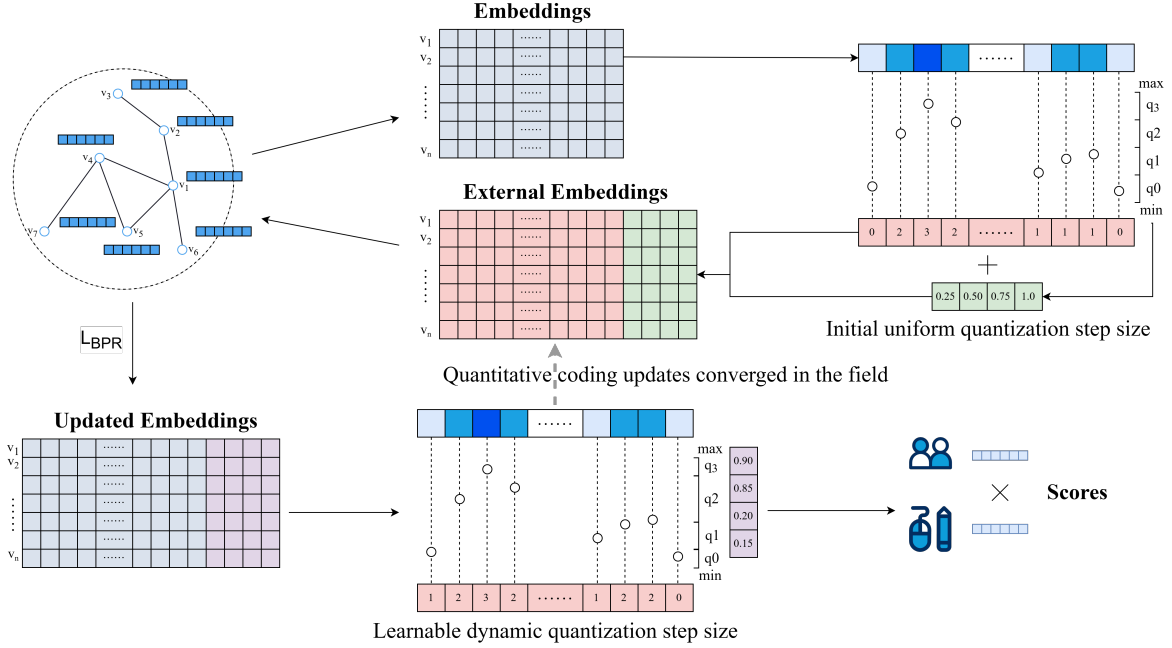


Figure 2: Overall Framework of Our GNAQ.

associate the encoding of each node with the encoding in the meta-embedding. During training, the encoding representation of each node is composed by combining the meta-embedding according to the allocation matrix. Through the backpropagation mechanism of the graph neural network, the update of the meta-embedding parameters is realized. Similarly, if the allocation matrix is used as the quantization encoding representation of the node and the meta-embedding is used as the quantization step size of each node, the quantization step size of the node can be determined using the backpropagation of the graph neural network through the idea of parameter sharing.

Our approach realizes graph node-aware dynamic quantization step size asymmetric n -bit quantization by assigning $2n$ quantization step sizes to each node. Using the idea of parameter sharing, the quantization step size of each node is concatenated into the quantization encoding of the node, and through the information transfer of the graph neural network, the dynamic update of the node quantization encoding and the quantization step size is achieved. Unlike traditional quantization methods that set fixed quantization functions to quantize parameters within a fixed range to specified quantization values, GNAQ deals with parameter quantization by using an allocation matrix similar to that in parameter sharing methods. During the learning process of the model for the allocation matrix and the quantization step size, the quantization function is learned dynamically, and the quantization of the graph neural network is carried out efficiently.

4.2 Initialization of Quantization Step Size

Our approach uses the node encoding table $E_{full} \in \mathbb{R}^{N \times d}$ of the full-precision model as the initial values of the quantized model.

Then, the maximum and minimum values in each node encoding are obtained. According to the maximum and minimum values of each node encoding, the initial quantization step size gap_i of each node encoding is calculated as:

$$gap_i = \frac{\max(e_i) - \min(e_i)}{2^b}, \quad (1)$$

where e_i is the full-precision encoding of node i ; N represents the total number of nodes; n represents the quantization precision, that is, the bit width used for quantization; max represents the maximum value by row; min is the minimum value by row.

The initial quantization step size is uniform. Based on the maximum and minimum values of each node encoding and the quantization bit width b , the value range of each node encoding is evenly divided into 2^n intervals. The initial quantization function based on these intervals is:

$$Q_i(a) = \begin{cases} 0, & \min_i \leq a < \min_i + gap_i \\ 1, & \min_i + gap_i \leq a < \min_i + 2 * gap_i \\ \vdots & \\ 2^n - 1, & \min_i + (2^n - 1) * gap_i \leq a \leq \max_i \end{cases} \quad (2)$$

The above formula represents the quantization function for quantizing the embedding representation $e_i \in \mathbb{R}^d$ of the i -th node. \min_i is the minimum value in e_i , \max_i is the maximum value in e_i , and gap_i is the quantization interval obtained by calculating according to Equation 1. The quantization function corresponding to each node encoding is executed to obtain the quantization encoding q_i and E_{quant} of each node:

$$q_i = [Q_i(e_i[0]), Q_i(e_i[1]), \dots, Q_i(e_i[d-1])] \quad (3)$$

$$E_{\text{quant}} = [q_0, q_1, \dots, q_{N-1}]^T \quad (4)$$

4.3 Update of Quantization Interval Scaling

The midpoint value of each quantization interval is used as the scaling factor of the corresponding node encoding. The scaling factor of the quantization interval for a node is computed as:

$$s_i = \left[\min_i + \frac{\text{gap}_i}{2}, \min_i + \frac{3 * \text{gap}_i}{2}, \dots, \min_i + \frac{(2^{n+1} - 1) * \text{gap}_i}{2} \right] \quad (5)$$

$$S = [s_0, s_1, \dots, s_{N-1}]^T \quad (6)$$

where $a_i \in \mathbb{R}^{2^n}$. The scaling factor of the quantization interval is represented by a low-precision floating-point number (FP8 [23]), which ensures the efficiency of subsequent calculations. Then, the scaling factor of the quantization interval, defined in Equations 7 to 8, is concatenated into the node encoding to form the extended quantization encoding table E.

$$e'_i = [q_i, s_i] \quad (7)$$

$$E = [e'_0, e'_1, \dots, e'_{n-1}], \quad (8)$$

where $e'_i \in \mathbb{R}^{d+2^n}$, $E \in \mathbb{R}^N \times (d+2^n)$. Using the extended quantization encoding, the quantization interval scaling factor can be updated with the help of the message-passing mechanism of the graph neural network, thus achieving dynamic quantization step sizes. In this way, different quantization functions can be constructed for each node encoding, improving the quantization's perception of node information and thus enhancing the quantization effect.

After completing the update of the scaling factor of the quantization interval, it is necessary to dynamically adjust the quantization step size. Therefore, the dynamic update mechanism of the quantization step size is described in detail below.

4.4 Dynamic Quantization Step Size

Forward Propagation: During the forward propagation process of the graph neural network, our extended quantization encoding is used as the input, and message passing is performed on the graph. Before inputting the extended quantization embedding into the graph neural network, it is necessary to transform the quantization embedding E_{quant} , mapping the quantization encoding to the quantization interval scaling factor:

$$q'_i = [s_i[q_i[0]], s_i[q_i[1]], \dots, s_i[q_i[d-1]]] \quad (9)$$

$$E'_{\text{quant}} = [q'_0, q'_1, \dots, q'_{N-1}] \quad (10)$$

In this way, the effect of non-uniform quantization is achieved. The transformed quantization encoding is concatenated with the

quantization interval scaling factor to obtain the input encoding table $H^{(0)} \in \mathbb{R}^{N \times (d+2^n)}$ of the graph neural network in Equation 11.

$$H^{(0)} = [E'_{\text{quant}}, S] \quad (11)$$

Then, the input encoding table is used as the initial encoding of the nodes and input into the graph neural network. Through the information transfer of the graph neural network, message passing is performed on the normalized Laplacian matrix to optimize the quantization encoding representation of the nodes and the quantization interval scaling factor. Then, the dynamic quantization step size is achieved through the relationship between the quantization interval scaling factor and the quantization step size.

It can be seen from Equation 5 and 7 that the initial quantization clipping interval of GNAQ is uniform. During the subsequent model training process, the quantization clipping interval will also change by learning the quantization interval scaling factor, thus realizing a learnable quantization clipping interval. Forward propagation across GNN layers adheres to:

$$H^{(l+1)} = \left(D^{-\frac{1}{2}} X D^{-\frac{1}{2}} \right) H^{(l)}, \quad (12)$$

$$H = \frac{1}{l+1} \sum_{l=0}^L H^{(l)}, \quad (13)$$

where X is the user-item interaction data, represented as an adjacency matrix; L is the number of layers of the graph neural network; D is the degree matrix of X , and $D^{-\frac{1}{2}} X D^{-\frac{1}{2}}$ forms a symmetric degree-normalized adjacency matrix; the final node encoding representation H is the average value of each layer's representation.

Update of Quantization Step Size: Since the initial quantization interval scaling factor is added to the node encoding of the graph neural network and calculated together with the encoding during the forward propagation of the graph neural network, the quantization interval scaling factor will be updated during the backpropagation process of the model. Suppose the updated quantization interval scaling factor is S' . The updated quantization step size can be calculated through S' and the current value range of the node encoding. The update method is as follows: First, sort the updated quantization interval scaling factors according to their numerical sizes. Then, take the midpoint positions between adjacent quantization interval scaling factors and the boundaries of the value distribution range as the boundaries of the quantization intervals. There are a total of $2^n - 1$ midpoint values, which can divide the value range of the node encoding into 2^n intervals, and the length of each interval is the updated quantization step size.

As shown in Figure 3, taking 2-bit quantization as an example, each node encoding has four quantization interval scaling factors s_1, s_2, s_3 , and s_4 . Initially, the intervals of the quantization interval scaling factors are uniform, and the quantization step sizes are also equal (i.e., $\text{step}_1 = \text{step}_2 = \text{step}_3 = \text{step}_4$). After model training, the quantization interval scaling factors are updated. First, rearrange the updated quantization interval scaling factors according to their numerical sizes, such as s'_1, s'_2, s'_3 , and s'_4 in Figure 3. Then, the midpoint of two scaling factors is taken as the new boundary of the quantization interval and four new quantization intervals

are obtained. The length of each interval is the updated quantization step size that are no longer equal. In this way, the value of the quantization step size is learned through model training, and the parameter optimization ability of the neural network model is utilized to efficiently optimize the quantization step size.

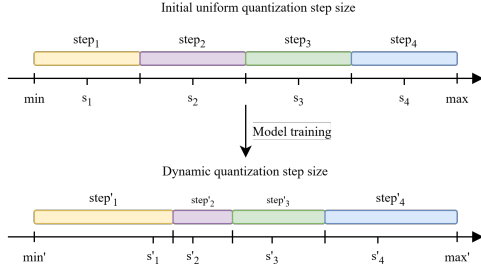


Figure 3: Dynamic Update of Quantization Step Size.

Update of Quantization Encoding Aggregated by Graph Interaction Relationships: As shown in Equation 2, the quantization function of node encoding is a discrete piecewise function. Therefore, at the edges of the piecewise intervals, the quantization function is non-differentiable; within the piecewise intervals, the derivative is constantly 0. This makes it impossible to update the model parameters through backpropagation during the model training process. Traditional quantization methods approximate the gradient of the quantization function through gradient estimation methods such as the straight through estimator to achieve parameter updates. Since our approach performs model quantization based on the idea of parameter sharing, the node quantization encoding is not directly updated during the backpropagation process of the model. Therefore, this paper combines the ideas of parameter sharing and low precision training (LTP [17]), regards the update of the node quantization encoding as the update of the allocation matrix. After the backpropagation of the model, the information of adjacent nodes is aggregated to achieve the update of the node quantization encoding, thus avoiding the estimation of the gradient of the quantization function.

Specifically, after the backpropagation of the model, the encoding and quantization interval scaling factor of each node will be updated. As described in Section 4.4, the quantization step size can be calculated based on the updated quantization interval scaling factor, and then the updated quantization function can be obtained. In GNAQ, the first-order aggregated information of nodes is used as the basis for updating the node quantization encoding, and the calculation formula for the first-order aggregated information of nodes is as follows:

$$\hat{e}_i = \frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} e_j, \quad (14)$$

where \mathcal{N}_i represents the set of neighbor nodes of node i , and $|\mathcal{N}_i|$ represents the number of neighbor nodes of node i . Then, the first-order aggregated information is input into the updated quantization function to obtain the updated node quantization encoding. Using

the first-order aggregated information of nodes as the basis for updating the node quantization encoding can avoid the "Z-shaped" oscillation phenomenon [21] during the parameter update process. Using the node representation aggregated by graph interaction relationships can limit the model to consider the interaction information in the graph during the update of the node quantization encoding, playing a part of the regularization effect.

4.5 Loss Function

To enhance the quantized model's ability to predict and rank interaction relationships, GNAQ incorporates the Bayesian Personalized Ranking Loss (BPR [25]) and LambdaLoss [29] to optimize the model parameters. The BPR loss is formulated as:

$$\mathcal{L}_{\text{BPR}} = \sum_{(u, i^+, i^-) \in \mathcal{B}} -\ln \sigma(\hat{y}_{ui^+} - \hat{y}_{ui^-}) + \lambda \|\Theta\|^2 \quad (15)$$

$$\mathcal{L}_{\text{Lambda}} = -\sum_{i=1}^K \sum_{j=1}^K \eta_{ij} \cdot \log(1 + e^{-(s_i - s_j)}) \quad (16)$$

$$\mathcal{L} = \mathcal{L}_{\text{BPR}} + \mathcal{L}_{\text{Lambda}}, \quad (17)$$

where \mathcal{B} represents the set of data in each batch; (u, i^+, i^-) is the triple representation of training data, with u denoting the user, i^+ representing the item that has an explicit interaction with the user u and i^- representing the item that has no explicit interaction with the user. $\|\Theta\|^2$ is the regularization term, and λ is the hyperparameter. η_{ij} represents the ranking scores of the i -th and j -th items in the prediction sequence. If the i -th item should be ranked before the j -th item, then $\eta_{ij}=1$; conversely, if the i -th item should be ranked after the j -th item, then $\eta_{ij}=0$. s_i and s_j represent the predicted scores of the i -th and j -th items in the prediction sequence, respectively. The BPR loss can optimize the relative preference relationship and effectively handle implicit feedback data. LambdaLoss can directly optimize the ranking metric (NDCG) and re-rank the initially predicted sequence of the model. By combining the BPR loss and LambdaLoss, the model's ability to capture ranking information can be strengthened, effectively improving the model's recommendation performance.

5 EXPERIMENTS

Our approach is evaluated on the Top-K recommendation task with the aim of answering the following research questions:

- **RQ1.** How does GNAQ perform compared to state-of-the-art full-precision and quantization-based models?
- **RQ2.** How is the practical resource consumption of GNAQ?
- **RQ3.** How do the proposed key components affect GNAQ performance?

5.1 Experimental Setup

Datasets. Our evaluation is conducted on the following four datasets: MovieLens [10], Gowalla [3], Yelp2020 and Amazon-Book, among which Yelp2020 and Amazon-Book can be found in the official code repositories of LightGCN [11] and LEGCF [18].

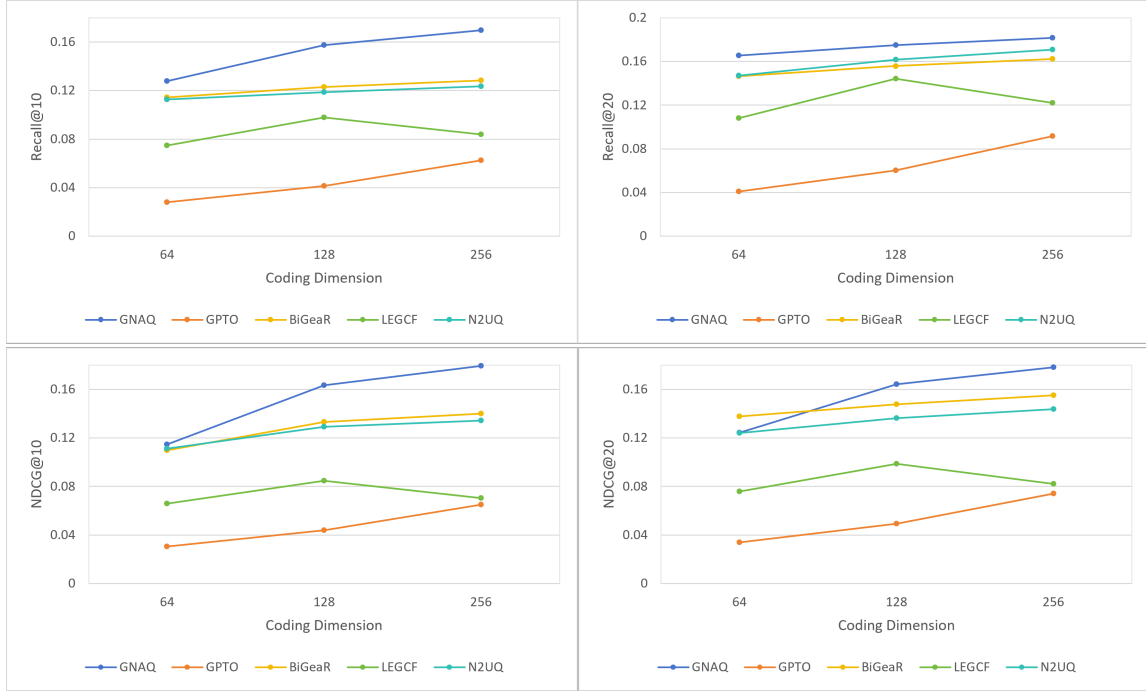


Figure 4: The Effect of Coding Dimension on Recommendation Performance.

Table 2: Parameter settings.

Parameter	MovieLens	Gowalla	Yelp2020	Amazon-Book
batch size	10240	5096	5096	5096
learning rate	1×10^{-3}	1×10^{-3}	5×10^{-4}	5×10^{-4}
embedding dim	[64, 128, 256]	[64, 128, 256]	[64, 128, 256]	[64, 128, 256]
optimizer	Adam	Adam	Adam	Adam
GNN layers	3	3	3	3
λ	5×10^{-4}	5×10^{-4}	5×10^{-4}	5×10^{-4}

Evaluation Metrics. Recall@K and NDCG@K are adopted as evaluation metrics, which are standard in the top-K recommendation tasks. Recall@K measures the proportion of relevant items retrieved in the top-K recommendations, while NDCG@K evaluates the ranking quality by considering the positions of relevant items.

Baselines. We compare our GNAQ against following baselines: **BiGeaR** [2]: presented by Chen et al. in 2022, it is a multi-faceted quantization-enhanced binarized graph representation learning method. **N2UQ** [20]: a quantization method proposed by Liu et al. in 2022, which transforms non-uniform quantization into uniform quantization and uses Generalized Straight. **GPTQ** [7]: proposed by Frantar et al. in 2023, it is a post-training quantization method that uses layer-wise quantization and OBQ (Optimal Brain Quantization) technique to minimize quantization errors. **LEGCF** [18]: put forward by Liang et al. in 2024, it is a lightweight GNN node embedding learning method based on parameter sharing.

Parameter Settings. The experimental results of all baseline models were reproduced in this paper, following the optimal settings provided in their publications. The learning rate is from the collection $1e-5, 5e-5, 1e-4, 5e-4, 1e-3$, the optimizer selects Adam. In

our model, the quantization bitwidth is set to 2 bits. To explore the impact of different node embedding dimensions on the quantization performance of graph neural network models, experiments were conducted under three embedding dimensions (64-dimensional, 128-dimensional, and 256-dimensional). Parameters such as batch size and training epochs were adjusted according to the node count and interaction data of different datasets, as detailed in Table 2.

We employ LightGCN [11], the current state-of-the-art graph neural network model for collaborative filtering tasks, as the baseline model. During model quantization, the initial parameters loaded are those of the pre-trained LightGCN model. Our source codes are available at: <https://github.com/WUT-IDEA>.

5.2 Overall Performance Analysis (RQ1)

Extensive experiments were conducted on four datasets with different embedding dimensions to compare our approach against the baseline models. The results are presented in Tables 3 to 5.

Superior Performance of GNAQ: Across almost all evaluation metrics on four datasets, our approach outperforms other baseline models. For example, on the Gowalla dataset (Figure 4), GNAQ achieves the highest Recall@10, Recall@20, and NDCG@10 across all embedding dimensions, except for a slight NDCG@20 deficit compared to BiGeaR at 64 dimensions. It may be because of insufficient information from neighbors, the Role of Relational Aggregation-Based Update (RAU) is limited in its effect during the quantization coding update. The average improvements in Recall@10, Recall@20, NDCG@10, and NDCG@20 are 27.8%, 10.2%, 17.6%, and 8.51%, respectively, demonstrating its effectiveness in preserving recommendation accuracy.

Table 3: Results of Recall@ and NDCG@ for GNAQ vs. Baseline Models with 64-dimensional Embeddings.

Model	Metric	MovieLens		Gowalla		Yelp2020		Amazon-Book	
		Recall	NDCG	Recall	NDCG	Recall	NDCG	Recall	NDCG
LightGCN [11]	@10	0.1551	0.2742	0.1088	0.0884	0.0409	0.0274	0.0204	0.0165
	@20	0.2418	0.2801	0.1557	0.1041	0.0687	0.0366	0.0353	0.0222
BiGeaR [2]	@10	<u>0.1479</u>	<u>0.3723</u>	<u>0.1144</u>	0.1098	<u>0.0452</u>	<u>0.0393</u>	0.0187	<u>0.0161</u>
	@20	<u>0.1796</u>	<u>0.3554</u>	0.1466	0.1378	0.0712	0.0498	0.0266	0.0206
N2UQ [20]	@10	0.1093	0.2494	0.1127	<u>0.1113</u>	0.0324	0.0371	0.0157	0.0168
	@20	0.1692	0.2419	<u>0.1471</u>	<u>0.1204</u>	0.0520	0.0418	<u>0.0282</u>	0.0218
GPTQ [7]	@10	0.0653	0.0897	0.0280	0.0304	0.0106	0.0091	0.0045	0.0052
	@20	0.1066	0.1027	0.0411	0.0340	0.0180	0.0117	0.0079	0.0064
LEGCF [18]	@10	0.0705	0.1435	0.0748	0.0660	0.0223	0.0160	0.0115	0.0100
	@20	0.1160	0.1427	0.1082	0.0759	0.0400	0.0218	0.0204	0.0133
GNAQ (Ours)	@10	0.2048	0.4934	0.1277	0.1145	0.0540	0.0428	<u>0.0186</u>	0.0157
	@20	0.2055	0.3902	0.1605	<u>0.1243</u>	<u>0.0685</u>	<u>0.0470</u>	0.0319	<u>0.0206</u>
%Improv.	@10	38.4%	12.1%	11.6%	2.88%	19.4%	8.91%	-0.53%	-6.55%
	@20	14.4%	9.79%	12.5%	-9.80%	-3.79%	-5.62%	13.1%	-5.50%

Table 4: Results of Recall@ and NDCG@ for GNAQ vs. Baseline Models with 128-dimensional Embeddings.

Model	Metric	MovieLens		Gowalla		Yelp2020		Amazon-Book	
		Recall	NDCG	Recall	NDCG	Recall	NDCG	Recall	NDCG
LightGCN [11]	@10	0.1658	0.2855	0.1101	0.0901	0.0426	0.0284	0.0215	0.0172
	@20	0.2585	0.2935	0.1576	0.1059	0.0721	0.0382	0.0367	0.0230
BiGeaR [2]	@10	<u>0.1551</u>	<u>0.3894</u>	<u>0.1229</u>	<u>0.1331</u>	<u>0.0503</u>	<u>0.0433</u>	<u>0.0219</u>	0.0177
	@20	0.1788	<u>0.3740</u>	0.1558	<u>0.1478</u>	<u>0.0762</u>	<u>0.0548</u>	<u>0.0379</u>	0.0228
N2UQ [20]	@10	0.1272	0.2889	0.1188	0.1291	0.0350	0.0399	0.0188	<u>0.0198</u>
	@20	<u>0.1861</u>	0.2569	<u>0.1618</u>	0.1363	0.0564	0.0459	0.0334	<u>0.0257</u>
GPTQ [7]	@10	0.0867	0.0949	0.0414	0.0440	0.0131	0.0111	0.0067	0.0076
	@20	0.1398	0.1153	0.0603	0.0494	0.0230	0.0146	0.0119	0.0095
LEGCF [18]	@10	0.0737	0.1523	0.0979	0.0846	0.0310	0.0214	0.0156	0.0134
	@20	0.1241	0.1525	0.1444	0.0988	0.0548	0.0291	0.0259	0.0172
GNAQ (Ours)	@10	0.2287	0.5261	0.1576	0.1635	0.0686	0.0624	0.0242	0.0204
	@20	0.2296	0.4194	0.1750	0.1643	0.0779	0.0643	0.0407	0.0266
%Improv.	@10	47.4%	13.7%	28.2%	22.8%	36.4%	44.1%	10.5%	3.03%
	@20	23.4%	12.1%	8.16%	11.2%	2.23%	17.3%	7.39%	3.50%

Impact of Embedding Dimensions: As the embedding dimension increases, our approach demonstrates stronger performance relative to baseline models. Higher-dimensional embeddings contain more information, and GNAQ’s dynamic quantization step strategy better captures inter-dimensional correlations in node embeddings. For instance, on the Yelp2020 dataset, GNAQ’s NDCG@10 at 256 dimensions is 48% higher than GPTQ, highlighting its advantage in high-dimensional scenarios.

Sparsity-Dependent Performance: Our approach achieves more significant improvements on datasets with lower sparsity (e.g., MovieLens). Dense interaction graphs in such datasets contain richer high-order interaction information, which GNAQ leverages through dynamic quantization steps to adapt to node embedding distributions, thereby enhancing recommendation quality.

5.3 Resource Consumption Analysis (RQ2)

Table 6 compares the parameter counts and training times of different methods. Key observations include:

Efficiency of Direct Optimization Methods: Directly optimizing quantized based models (e.g., GNAQ, N2UQ) exhibit lower training times compared to knowledge distillation-based approaches (e.g., BiGeaR). For example, GNAQ’s training time on the Amazon-Book dataset is 133,482 seconds, significantly less than BiGeaR’s 298,037 seconds, as it avoids re-training full-precision models.

Balance Between Efficiency and Effectiveness: GNAQ’s parameter size is higher than that of lightweight models (e.g. LEGCF), but achieves superior performance in terms of Recall and NDCG. Compared to BiGeaR, GNAQ uses fewer parameters while outperforming it in most scenarios, demonstrating a better trade-off between efficiency and effectiveness.

Table 5: Results of Recall@ and NDCG@ for GNAQ vs. Baseline Models with 256-dimensional Embeddings.

Model	Metric	MovieLens		Gowalla		Yelp2020		Amazon-Book	
		Recall	NDCG	Recall	NDCG	Recall	NDCG	Recall	NDCG
LightGCN [11]	@10	0.1708	0.2904	0.1305	0.1142	0.0548	0.0386	0.0261	0.0217
	@20	0.2641	0.2988	0.1833	0.1301	0.0887	0.0496	0.0446	0.0286
BiGeaR [2]	@10	<u>0.1631</u>	<u>0.4051</u>	<u>0.1283</u>	<u>0.1400</u>	<u>0.0550</u>	<u>0.0475</u>	<u>0.0240</u>	0.0199
	@20	0.1823	<u>0.3879</u>	0.1622	<u>0.1551</u>	<u>0.0811</u>	<u>0.0597</u>	<u>0.0417</u>	<u>0.0258</u>
N2UQ [20]	@10	0.1340	0.2952	0.1235	0.1344	0.0375	0.0427	0.0225	<u>0.0238</u>
	@20	<u>0.2091</u>	0.2928	<u>0.1707</u>	0.1438	0.0606	0.0493	0.0392	0.0304
GPTQ [7]	@10	0.1097	0.0970	0.0624	0.0649	0.0208	0.0175	0.0132	0.0141
	@20	0.1733	0.1184	0.0918	0.0740	0.0350	0.0226	0.0222	0.0175
LEGCF [18]	@10	0.0784	0.1593	0.0839	0.0704	0.0268	0.0191	0.0162	0.0139
	@20	0.1277	0.1586	0.1222	0.0822	0.0467	0.0255	0.0283	0.0183
GNAQ (Ours)	@10	0.2353	0.5370	0.1698	0.1793	0.0780	0.0703	0.0297	0.0243
	@20	0.2362	0.4288	0.1815	0.1781	0.0869	0.0720	0.0496	0.0318
%Improv.	@10	44.3%	32.6%	32.3%	28.1%	41.8%	48.0%	23.8%	2.10%
	@20	12.9%	10.5%	6.33%	14.8%	7.15%	20.6%	18.9%	23.3%

Table 6: Parameter sizes of different quantization methods and training time (in seconds) under the 256-dimensional encoding.

Model	Number of parameters	MovieLens	Gowalla	Yelp2020	Amazon-Book
BiGeaR [2]	$O(N(L + 1)(32 + d))$	14494	33950	61217	298037
N2UQ [20]	$O(2N(32 + d))$	2188	16002	29162	32604
GPTQ [7]	$O(4N(32 + d))$	32	119	332	1069
LEGCF [18]	$O(500d + 2N)$	6780	34823	37904	75916
GNAQ (Ours)	$O(N(d + 128))$	6880	14815	41765	133482

Table 7: Ablation Experiments

Model	Gowalla		Yelp2020	
	Recall@20	NDCG@20	Recall@20	NDCG@20
GNAQ	0.1750	0.1643	0.0779	0.0643
w/o DQS	0.1613	0.1263	0.0708	0.0431
w/o RAU	0.1424	0.0973	0.0518	0.0274
w/o rank_loss	0.1689	0.1245	0.0740	0.0432

5.4 Ablation Experiments (RQ3)

Ablation studies are conducted on two datasets at 128 embedding dimensions to analyze the contributions of dynamic quantization steps (DQS), relational aggregation-based update (RAU), and ranking-enhanced loss (rank-loss), shown in Table 7.

Effectiveness of Dynamic Quantization Steps (DQS): Removing DQS (w/o DQS) leads to a significant performance degradation (e.g., Gowalla’s NDCG @ 20 drops from 0.1643 to 0.1263), confirming that node-aware dynamic quantization steps effectively adapt to embedding distributions and reduce information loss.

Role of Relational Aggregation-Based Update (RAU): Removing RAU (w/o RAU) causes more severe performance declines (e.g., Yelp2020’s Recall@20 drops to 0.0518), indicating that graph interaction-aware embedding updates are critical for quantization accuracy by guiding parameter updates with structural information.

Impact of Ranking-Enhanced Loss (rank-loss): Removing rank-loss (w/o rank-loss) reduces NDCG@20 while preserving Recall@20, demonstrating its specific contribution to improving ranking accuracy without compromising recall, thus improving overall quality of recommendation.

6 CONCLUSION AND FUTURE WORK

This paper proposes a novel graph node-aware quantization with dynamic step size for collaborative filtering. Our GNAQ effectively captures the structural heterogeneity of graph data and achieves superior performance over state-of-the-art quantization methods. Extensive experiments on four public datasets demonstrate the effectiveness and efficiency of our approach. In the future, we apply GNAQ to actual edge deployment scenarios to achieve efficient inference on resource-constrained devices while maintaining high recommendation performance.

ACKNOWLEDGMENTS

This work is partially supported by NSFC, China (No.62276196) and JSPS Invitational Fellowships for Research in Japan (Short-term).

USAGE OF GENERATIVE AI

The general outline of the entire article was independently completed by us and Generative AI was used to modify the grammar of the sentences and polish the expression of the language.

REFERENCES

- [1] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. 2013. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432* (2013).
- [2] Yankai Chen, Huifeng Guo, Yingxue Zhang, Chen Ma, Ruiming Tang, Jingjie Li, and Irwin King. 2022. Learning binarized graph representations with multi-faceted quantization reinforcement for top-k recommendation. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 168–178.
- [3] Eunjoon Cho, Seth A Myers, and Jure Leskovec. 2011. Friendship and mobility: user movement in location-based social networks. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. 1082–1090.
- [4] Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. 2016. Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1. *arXiv preprint arXiv:1602.02830* (2016).
- [5] Steven K Esser, Jeffrey L McKinstry, Deepika Bablani, Rathinakumar Appuswamy, and Dharmendra S Modha. [n. d.]. LEARNED STEP SIZE QUANTIZATION. In *International Conference on Learning Representations*.
- [6] Boyuan Feng, Yuke Wang, Xu Li, Shu Yang, Xueqiao Peng, and Yufei Ding. 2020. Sgquant: Squeezing the last bit on graph neural networks with specialized quantization. In *2020 IEEE 32nd international conference on tools with artificial intelligence (ICTAI)*. IEEE, 1044–1052.
- [7] E Frantar, S Ashkboos, T Hoefler, and D Alistarh. [n. d.]. OPTQ: Accurate Quantization for Generative Pre-trained Transformers. 2023. In URL <https://openreview.net/forum>.
- [8] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. 2022. A survey of quantization methods for efficient neural network inference. In *Low-power computer vision*. Chapman and Hall/CRC, 291–326.
- [9] Robert M. Gray and David L. Neuhoff. 1998. Quantization. *IEEE transactions on information theory* 44, 6 (1998), 2325–2383.
- [10] F Maxwell Harper and Joseph A Konstan. 2015. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)* 5, 4 (2015), 1–19.
- [11] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*. 639–648.
- [12] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. 2018. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2704–2713.
- [13] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).
- [14] Yehuda Koren, Robert Bell, and Chris Volinsky. 2009. Matrix factorization techniques for recommender systems. *Computer* 42, 8 (2009), 30–37.
- [15] Ming Li, Lin Li, Xiaohui Tao, and Jimmy Xiangji Huang. 2024. Mealrec+: A meal recommendation dataset with meal-course affiliation for personalization and healthiness. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 564–574.
- [16] Ming Li, Lin Li, Xiaohui Tao, Zhongwei Xie, Qing Xie, and Jingling Yuan. 2024. Boosting healthiness exposure in category-constrained meal recommendation using nutritional standards. *ACM Transactions on Intelligent Systems and Technology* 15, 4 (2024), 1–28.
- [17] Shiwei Li, Huifeng Guo, Xing Tang, Ruiming Tang, Lu Hou, Ruixuan Li, and Rui Zhang. 2024. Embedding compression in recommender systems: A survey. *Comput. Surveys* 56, 5 (2024), 1–21.
- [18] Xurong Liang, Tong Chen, Lizhen Cui, Yang Wang, Meng Wang, and Hongzhi Yin. 2024. Lightweight Embeddings for Graph Collaborative Filtering. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1296–1306.
- [19] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. 2024. Awq: Activation-aware weight quantization for on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems* 6 (2024), 87–100.
- [20] Zechun Liu, Kwang-Ting Cheng, Dong Huang, Eric P Xing, and Zhiqiang Shen. 2022. Nonuniform-to-uniform quantization: Towards accurate quantization via generalized straight-through estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4942–4952.
- [21] Lianbo Ma, Yue Zhou, Jianlun Ma, Guo Yu, and Qing Li. 2024. One-step forward and backtrack: overcoming zig-zagging in loss-aware quantization training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 14246–14254.
- [22] Chaojun Meng, Changfan Pan, Hongji Shu, Qing Wang, Hanghui Guo, and Jia Zhu. 2025. Heterogeneous collaborative filtering contrastive learning for social recommendation. *Applied Soft Computing* 173 (2025), 112934.
- [23] Paulius Micikevicius, Dusan Stolic, Neil Burgess, Marius Cornea, Pradeep Dubey, Richard Grisenthwaite, Sangwon Ha, Alexander Heinecke, Patrick Judd, John Kamalu, et al. 2022. Fp8 formats for deep learning. *arXiv preprint arXiv:2209.05433* (2022).
- [24] Ruilin Pan, Chuanming Ge, Li Zhang, Wei Zhao, and Xun Shao. 2020. A new similarity model based on collaborative filtering for new user cold start recommendation. *IEICE TRANSACTIONS on Information and Systems* 103, 6 (2020), 1388–1394.
- [25] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2012. BPR: Bayesian personalized ranking from implicit feedback. *arXiv preprint arXiv:1205.2618* (2012).
- [26] Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. 2024. OmniQuant: Omnidirectionally Calibrated Quantization for Large Language Models. In *ICLR*.
- [27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [28] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural graph collaborative filtering. In *Proceedings of the 42nd international ACM SIGIR conference on Research and development in Information Retrieval*. 165–174.
- [29] Xuanhui Wang, Cheng Li, Nadav Golbandi, Michael Bendersky, and Marc Najork. 2018. The lambdaloss framework for ranking metric optimization. In *Proceedings of the 27th ACM international conference on information and knowledge management*. 1313–1322.
- [30] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. 2021. Self-supervised graph learning for recommendation. In *Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval*. 726–735.
- [31] Le Wu, Xiangnan He, Xiang Wang, Kun Zhang, and Meng Wang. 2022. A survey on accuracy-oriented neural recommendation: From collaborative filtering to information-rich recommendation. *IEEE Transactions on Knowledge and Data Engineering* 35, 5 (2022), 4425–4445.
- [32] Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. 2022. Graph neural networks in recommender systems: a survey. *Comput. Surveys* 55, 5 (2022), 1–37.
- [33] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S Yu. 2020. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems* 32, 1 (2020), 4–24.
- [34] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. 2023. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*. PMLR, 38087–38099.
- [35] Haoran Yang, Hongxu Chen, Lin Li, Philip S Yu, and Guandong Xu. 2021. Hyper meta-path contrastive learning for multi-behavior recommendation. In *2021 IEEE International Conference on Data Mining (ICDM)*. IEEE, 787–796.
- [36] Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. 2019. Q8bert: Quantized 8bit bert. In *2019 Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing-NeurIPS Edition (EMCC2-NIPS)*. IEEE, 36–39.
- [37] Zeyu Zhu, Fanrong Li, Zitao Mo, Qinghao Hu, Gang Li, Zejian Liu, Xiaoyao Liang, and Jian Cheng. 2023. A²Q: Aggregation-Aware Quantization for Graph Neural Networks. In *ICLR*.