

GMTRouter: Personalized LLM Router over Multi-turn User Interactions

Anonymous ACL submission

Abstract

Large Language Model (LLM) routing has demonstrated strong capability in balancing response quality with computational cost. As users exhibit diverse preferences, personalization has attracted increasing attention in LLM routing, since even identical queries may require different models to generate responses tailored to individual needs. However, existing approaches are not fully personalized and often fail to faithfully capture the complex interactions between users and LLMs. Moreover, user preference data is typically scarce and inconsistent in format, which limits the effectiveness of methods that directly leverage user-specific data. To address these challenges, we propose *GMTRouter*, which represents multi-turn user-LLM interactions as a heterogeneous graph with five node types: user, LLM, query, response and turn, thereby maximally preserving the rich relational structure of the interaction. Through a lightweight inductive graph learning framework combined with a tailored user-conditioned graph sampling mechanism, *GMTRouter* learns to capture user preferences from few-shot data, enabling effective personalization. Extensive experiments demonstrate that *GMTRouter* outperforms the strongest baselines, achieving up to a 0.105 absolute improvement in accuracy and a 0.12 improvement in AUC. More importantly, we further demonstrate that *GMTRouter* can adapt to new users using only few-shot data, without extensive fine-tuning.

1 Introduction

With the rapid development of the field of Large Language Models (LLMs), an increasing number of models with varying sizes, computational costs, and domain expertise have become available (Singhal et al., 2023; Luo et al., 2022). This makes LLM routing particularly important, as it enables the recommendation of appropriate LLMs for diverse user

queries while balancing response quality with computational cost (Šakota et al., 2024; Stripelis et al., 2024). Such routing techniques are increasingly adopted in modern LLMs, including GPT-5 (OpenAI, 2025). At the same time, as more users engage with LLM routing services, differences in individual preferences become increasingly prominent: even identical queries may require different models to generate responses tailored to each user (Li et al., 2024b; Salehi et al., 2024). Therefore, this paper aims to highlight a pressing research question: *Can we design a personalized routing framework that aligns LLM selection with individual user preferences based on their interaction histories?*

Existing research has proposed various architectures for LLM routing frameworks: FrugalGPT introduces a BERT-based router that determines whether to switch to a larger LLM (Chen et al., 2023), while C2MAB-V constructs a bandit-based router to balance exploration and exploitation when selecting an LLM (Dai et al., 2024). GraphRouter formulates routing as a node classification task over a graph of queries, tasks, and LLMs (Feng et al., 2024). However, existing methods largely overlook the importance of extracting structured preference information from users’ interaction histories: they are not fully personalized and often fail to faithfully model multi-turn conversations between users and LLMs, which represent the most common form of user-LLM interaction in real-world scenarios (Zhang et al., 2025a; Li et al., 2025b). Moreover, in real-world scenarios, the preference data provided by a single user is typically scarce and inconsistent in format (Escamocher et al., 2024; Li et al., 2024a). This makes it challenging for methods that directly leverage user-specific data to learn user profiles (Salemi et al., 2024a; Gao et al., 2024) or use such data as a retrieval source to support routing (Au et al., 2025), thereby limiting their effectiveness.

To address these challenges, we introduce **GMTRouter (Graph Multi-Turn Router)**, a het-

User ID	Query	Selected LLM	Response	Feedback
User 1	[Turn 1]"Please Explain ... ?" [Turn 2]"Can a Process ... ?"	GPT-4-1106-Preview	[Turn 1]"Exothermic and endothermic ..." [Turn 2]"Yes, a process ..."	[Turn 1]"rating: 3.0" [Turn 2]"rating: 4.5"
User 2	[Turn 1]"Compose a blog ..."	Claude-V1	[Turn 1]"Title: Aloha Spirit ..."	[Turn 1]"ranking: Claude-V1 > Koala-13B"
User 2	[Turn 1]"Compose a blog ..."	Koala-13B	[Turn 1]"Aloha, fellow travelers ! ..."	[Turn 1]"ranking: Claude-V1 > Koala-13B"
User 3	[Turn 1]"Compose an email ..." [Turn 2]"Rewrite your ..."	Vicuna-13B	[Turn 1]"Subject: An Exciting ..." [Turn 2]"Subject: A Gental ..."	[Turn 1]"response: Subject: Embrace ..." [Turn 2]"response: Subject: Soaring to ..."

Figure 1: **Multi-turn user-LLM Interaction History Table**. Each row captures a multi-turn interaction with associated user feedback. User feedback can take various forms, including ratings, rankings, ground-truth responses.

erogeneous graph-based LLM router based on multi-turn user interactions for personalized LLM routing. GMTRouter first sensitively identifies key entities within the user-LLM interaction process: users, LLMs, queries, and responses. By modeling these entities as different types of nodes and encoding their textual information into node embeddings, it maximally preserves the semantic information from the original data. To faithfully model the relational structure of multi-turn user-LLM interactions, GMTRouter organizes these diverse node types into a heterogeneous graph that captures complex relational dependencies. Each single-turn interaction is treated as a fundamental unit, and a virtual node, referred to as a *turn node*, is introduced to aggregate local information within each interaction round. We further transform user preference feedback into node features, enabling preference information to propagate across the graph. Moreover, rather than training the model directly on large historical datasets, GMTRouter employs a novel **user-conditioned graph sampling** mechanism in conjunction with an inductive graph learning framework to **enhance the model’s ability to capture user preferences from few-shot data**. This design allows effective test-time personalization even under sparse interaction histories, such as cold-start scenarios involving new users. In summary, our main contributions are as follows:

- To our knowledge, we are among the first to introduce a personalized LLM routing task based on multi-turn user interactions, providing new insights for this rapidly growing research field.
- We propose a novel personalized LLM routing framework, which faithfully models user-LLM interactions as a heterogeneous graph, and learns to capture user preferences from few-shot data within a lightweight inductive framework.
- Through experiments on multiple datasets spanning diverse tasks, GMTRouter outperforms the

strongest baselines, achieving up to a 0.105 absolute improvement in accuracy and a 0.12 improvement in AUC. Moreover, we demonstrate that our method can efficiently adapt to unseen users with only a few interaction examples, without requiring retraining.

2 Preliminaries

2.1 Task Formulation

We introduce the personalized LLM routing task in this section. We focus on the multi-turn interaction scenario between users and LLMs with feedback (Wang et al., 2023b; Shi et al., 2024). Within a dialogue session, a user u repeatedly interacts with an LLM m : in each turn, the user issues a query q , the LLM provides a response r , and the user in turn supplies a piece of feedback f . Such feedback can take multiple forms, including: (1) scalar scores (e.g., numerical ratings), (Wang et al., 2023c, 2024b); (2) preference rankings (e.g., choosing among multiple responses), (Yang et al., 2024; Sun et al., 2025); (3) ground-truth responses (e.g., directly providing the correct answer) (Gao et al., 2024; Salemi et al., 2024a). We structure these interactions into an **Interaction History Table**, illustrated in Figure 1, where each entry records the user ID, the selected LLM, the multi-turn queries and generated responses, and the corresponding user feedback, thereby maximally preserving the rich relational information of the interaction.

Our personalized LLM routing task is then modeled as follows: Given m users $\{u_1, \dots, u_m\}$ and n LLM candidates $\{m_1, \dots, m_n\}$, as well as their historical interaction records:

$$\mathcal{H} = \{(u_i, m_i, \{(q^{(t)}, r^{(t)}, f^{(t)})\}_{t=1}^{T_i})\},$$

where u_i is the user, m_i is the selected LLM, and each record contains a multi-turn sequence of queries $q^{(t)}$, responses $r^{(t)}$, and feedback $f^{(t)}$

Table 1: **The consistency of LLM preferences between users is significantly lower than the consistency within a single user’s preferences.** The self-spearman score is substantially higher than the spearman scores computed across different users.

Metric	Self	Global	Intra-cluster	Inter-cluster
Spearman	0.793	0.524	0.573	0.442
Percent	100%	66.0%	72.3%	55.7%

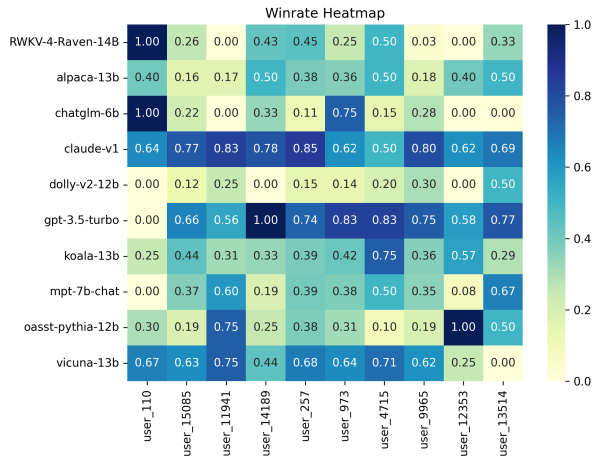


Figure 2: **Significant differences exist in LLM preferences across users.** The figure shows a heatmap of win rates for the 10 most popular LLMs across 10 active users in ChatBot Arena. The uneven color intensity within each row visually highlights the pronounced preference differences between users.

for $t = 1, \dots, T_i$. When a user u raises a new query q , the router is required to select an LLM $m \in \{m_1, \dots, m_n\}$ to generate a response r that best aligns with the user preferences, which is measured through the feedback f provided by the user.

2.2 Motivation

In this section, we highlight the significant differences in LLM preferences across users in the real world (Chevi et al., 2025; Wang et al., 2024a), emphasizing the importance of personalized LLM routing for enhancing user experience. We use the ChatBot Arena dataset (Chiang et al., 2024) to illustrate our findings, which contains extensive anonymized multi-turn conversations from numerous users, with pairwise human preference labels between various LLMs, enabling the study of real-world user–LLM interactions. From this dataset, we select 10 active users, each with at least 50 records, for detailed analysis. For each user, we randomly split their data into two halves and compute the win rates of each LLM within each half. We use Spearman correlation to quantify the consistency of preference rankings over LLMs (De Win-

ter et al., 2016; Hauke and Kossowski, 2011). We then compute the Spearman correlation between the two halves to quantify their self-consistency in preferences over LLMs (Chevi et al., 2025; Jiang et al., 2025), reporting the average as a baseline for comparison with inter-user preference consistency. Next, based on the similarity of queries in each user’s interaction history, we cluster users into three groups (Zeng et al., 2024; Li et al., 2025a), and compute pairwise Spearman correlation scores among users globally, within clusters, and across clusters (Cavallo, 2019; De Winter et al., 2016), reporting the corresponding averages as summarized in Table 1. We observe that global consistency in LLM preferences among users is substantially lower than individual self-consistency, reaching only 65.99% of the latter. Even within the same cluster, the Spearman score is only 72.28% of the self-consistency, highlighting the diversity of user preferences toward LLMs (Sun et al., 2025; Salemi et al., 2024a). To further visualize these differences, we select the 10 most frequently used models across these 10 users and present a win-rate heatmap in Figure 2, offering an intuitive depiction of the variability in user preferences. Therefore, to address the substantial inconsistency of LLM preferences across users, we propose **GMTRouter**, a framework that enables the personalized recommendation of suitable LLMs tailored to each user’s individual preferences.

3 GMTRouter: Router Over Multi-turn User Interactions

Method Overview As shown in Figure 3, GMTRouter operates in three stages: (a) *Node Embeddings Initialization*: It first identifies the key entities in the Interaction History Table—users, LLMs, queries, responses, and feedback—modeling them as nodes and encoding the textual information into node embeddings to maximally preserve the information of the interaction process. (b) *Heterogeneous Graph Construction*: Based on the relational structure of user–LLM interactions, these nodes are connected to form a heterogeneous graph, which captures rich relational dependencies. (c) *Inductive GNN Training*: Finally, we employ a user-conditioned graph sampling mechanism and an inductive graph training framework to learn user preferences from few-shot data, enabling effective personalization under sparse interaction histories.

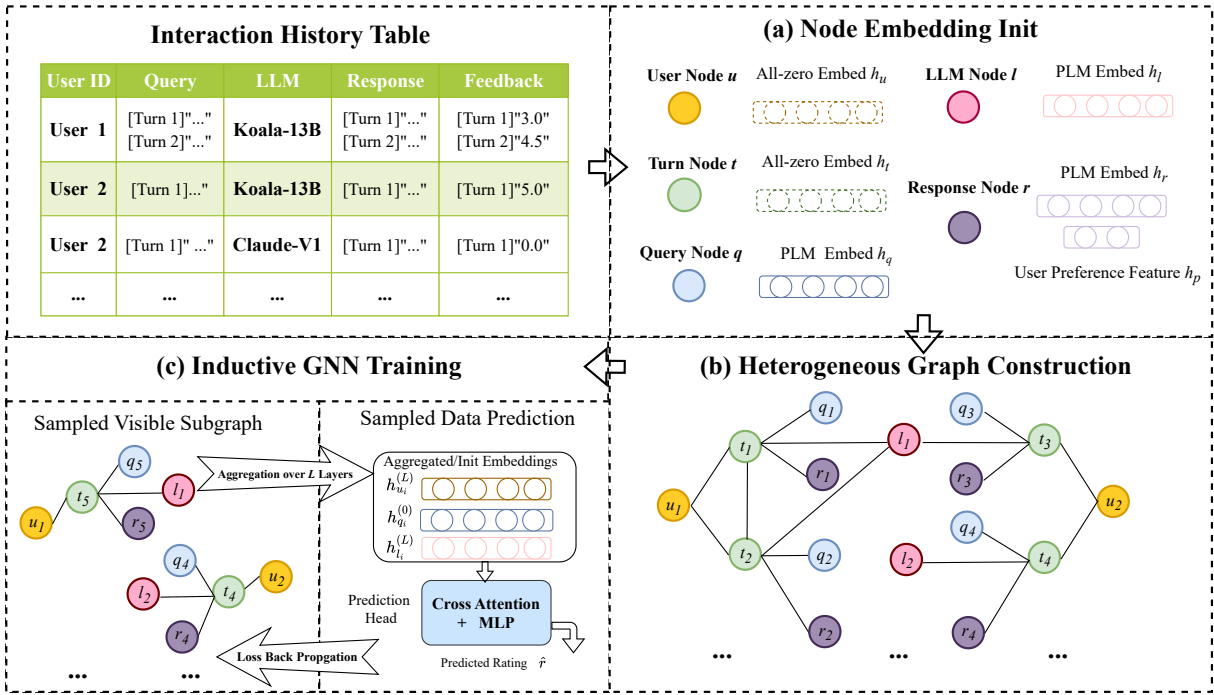


Figure 3: **Overview of GMTRouter.** (a) GMTRouter first extracts key entities: users, LLMs, queries, responses, and feedback, from the Interaction History Table and encodes their textual information using a PLM. (b) It then organizes these entities into a heterogeneous graph to faithfully model the relational structure of user–LLM interactions. (c) Within a lightweight inductive framework, GMTRouter learns to capture user preferences from few-shot data.

3.1 Node Embeddings Initialization.

First, our framework focuses on comprehensively extracting the information of various entities involved in the user–LLM interaction process from the Interaction History Table, along with their relational structures. As illustrated in part (a) of Figure 3, we extract four types of entities: user u , LLM m , query q , and response r , and formalize them as four corresponding node types. Their textual information is encoded using a pretrained language model (PLM) to obtain the initial node embeddings (Wang et al., 2022, 2023a), thereby preserving the semantic information from the original data. Specifically, we encode the query and response texts as their initial embeddings, denoted as h_q and h_r . In addition, we transform various forms of user feedback into numerical ratings and project them into a User Preference Feature h_p , which serves as another attribute on the response nodes. Concretely, ranking feedback is discretized into numerical ratings to ensure that higher-ranked responses receive higher scores (Banditwattana-wong and Masdisornchote, 2025); for ground-truth response feedback, we compute the geometric distance between the embeddings of the ground-truth and the generated response as the rating criterion (Salemi et al., 2024a). For LLM nodes, instead of simply using their names or IDs (Ding et al., 2024; Chen et al., 2023), we encode the model

overviews provided by AI/ML API platforms¹ as their node embeddings h_m , which typically include key information such as model size, usage cost, and domain-specific capabilities, thereby enriching the node embeddings with important background knowledge. Finally, for user nodes, we do not assume the existence of text-based user profiles, as such information is often scarce and noisy in real-world applications (Su et al., 2024; Alzubaidi et al., 2023); therefore, we initialize user embeddings h_u as zero vectors.

3.2 Heterogeneous Graph Construction.

Next, we organize these nodes into a heterogeneous graph to model the relational structure of user–LLM interactions (Zhang et al., 2025b; Schlichtkrull et al., 2017). We consider each single-round user–LLM interaction as a fundamental unit and introduce a kind of virtual node, *the turn node*, to aggregate the information within each interaction round. As illustrated in part (b) of Figure 3, within each interaction round, the associated user, LLM, query, and response nodes are connected to a corresponding turn node that aggregates information from that round. For multi-turn conversations, the turn nodes corresponding to each round are sequentially connected in dialogue order, facilitating information propagation across turns. The turn

¹<https://aimlapi.com/models/>

node embedding h_t is initialized as zero vectors. The resulting heterogeneous graph captures the rich relational dependencies inherent in user–LLM interactions, where turn nodes aggregate local information within each dialogue round and propagate it to user nodes, thereby facilitating the global aggregation of user preference information.

3.3 Inductive GNN Training

After constructing the user–LLM interaction histories into a heterogeneous graph, we train our GNN model on it. Notably, GMTRouter is a **general framework** that can incorporate any heterogeneous GNN as its backbone. We denote the GNN backbone used in our method as *GNN* throughout the rest of the paper. To address scenarios with sparse user history (Su et al., 2024), instead of training the model directly on large amounts of historical data (Lin et al., 2021; Wang et al., 2025), We employ an inductive framework along with **user-conditioned graph sampling** during training, enabling GMTRouter to **capture a user’s preferences from only a few interaction records**.

User-conditioned Graph Sampling As illustrated in the left of (c) in Figure 3, during each training epoch, we sample k interaction histories for each user to construct a visible subgraph \mathcal{G}_{sub} from the heterogeneous graph for message passing, and further sample data outside the visible subgraph as the prediction targets. We then use only these small sampled visible subgraphs and perform message aggregation separately for each user to update the node embeddings. Formally, at layer l , node v ’s embedding $h_v^{(l)}$ is updated by aggregating messages from its neighbors via the backbone GNN’s message-passing mechanism:

$$h_v^{(l)} = \text{Norm} \left(\text{Dropout} \left(\text{GNN}^{(l)}(h_v^{(l-1)}, \mathcal{G}_{\text{sub}}) \right) \right) \quad (1)$$

This restriction on the amount of data involved in message passing encourages the model to learn how to infer user preferences from very limited signals and to generalize efficiently to new users.

LLM Routing with a Prediction Head. After completing L layers of message aggregation, we obtain the updated node representations $h^{(L)}$. We then employ a **Prediction Head** module f_{pred} for preference prediction. As illustrated in the right of (c) in Figure 3, the Prediction Head takes the updated user embedding $h_u^{(L)}$, LLM embedding $h_m^{(L)}$,

and the query embedding $h_q^{(0)}$ from PLM as input. It applies a cross-attention module, where the LLM embedding attends to the fused user–query context to extract relevant preference signals. The module outputs a scalar score $s_{u,q,m}$ for each LLM candidate, representing the likelihood that user u would prefer m to answer query:

$$s_{u,q,m} = f_{\text{pred}}(h_u^{(L)}, h_q^{(0)}, h_m^{(L)}) \quad (2)$$

These scores are then used to rank LLM candidates under the same (u, q) condition. We normalize both the predicted scores and the ground-truth ratings, and apply a criterion function to compute the training loss, which is subsequently used to update the model parameters.

During inference, when a user raises a new query, we first sample k interaction histories of that user to construct the visible subgraph and update the node embeddings. Then, the LLM candidate is selected from the candidate set \mathcal{M} as the one with the highest predicted score:

$$m^* = \arg \max_{m \in \mathcal{M}} f_{\text{pred}}(h_u^{(L)}, h_q^{(0)}, h_m^{(L)}) \quad (3)$$

4 Experiment Setup

4.1 Datasets and data processing

We evaluate our approach on five datasets covering diverse tasks, spanning both real-world and synthetic settings. The two real-world datasets are: (1) **ChatBot Arena** (Chiang et al., 2024), which consists of anonymized multi-turn conversations with human preference labels, and (2) **LaMP** (Salemi et al., 2024b), from which we select the Personalized Scholarly Title Generation task. The remaining three are synthetic datasets: (3) **MT-Bench** (Zheng et al., 2023), assessing multi-turn reasoning and conversational ability; (4) **GSM8K** (Cobbe et al., 2021), focused on grade school math problem solving; and (5) **MMLU** (Hendrycks et al., 2021a,b), measuring general knowledge and multi-domain reasoning. Detailed descriptions of the datasets and preprocessing are provided in Appendix B.1 and B.4.

Data Processing For ChatBot Arena, we discretize the pairwise preferences to serve as the ratings for responses. For the other datasets, we adopt the data collected in Ong et al. 2024, which generated responses to all questions using "GPT-4-1106-preview" (Achiam et al., 2023) and "Mixtral-8x7B-Instruct-v0.1" (Jiang et al., 2024) Based on this, we

convert these datasets into multi-user personalized datasets. Specifically, for each response, we consider the following four dimensions: (a) Quality: For open-ended questions, we use the GPT-4 scores provided by Ong et al., 2024; for objective questions, we directly evaluate the correctness. (b) Cost: We calculate the cost of generating each response based on the API pricing provided by AI/ML API platform. (c) Response Length: We compute the token length of each response using the Contriever tokenizer (Izacard et al., 2021). (d) Rare Words: We count the number of rare words in each response using the *wordfreq* package (Speer, 2022).

We obtain the final rating of a response by computing a weighted sum of these four metrics. Different users are assigned different weightings to reflect their individual preferences over these dimensions (Feng et al., 2024, 2025). The specific weights used are provided in Appendix B.3.

Data Splitting For all datasets, we partition the data into training, validation, and test sets with a 7:1:2 ratio. For the GMTRouter, we further adopt an additional splitting strategy: we sample 30% of the users and restrict their data to the test sets only, in order to evaluate the generalization ability of our method to new users unseen during training.

4.2 Baselines

We compare our GMTRouter against the following baselines:

Prompt-based: (1) Vanilla LLM. We incorporate the query and the descriptions of candidate LLMs into the prompt, and feed it into LLaMA-3.1-70B (Grattafiori et al., 2024) to select the LLM. **(2) Personalized LLM.** Building on the vanilla LLM, we retrieve the ten most relevant interaction histories from the training set for a given user query and incorporate them into the prompt. Leveraging in-context learning (Dong et al., 2022), the LLM is then guided to perform personalized routing.

Representative Router: (3) GraphRouter. (Feng et al., 2024) A graph-based model that formulates routing as a node classification task over a graph of queries, tasks, and LLMs with learned edge interactions. **(4) FrugalGPT (Chen et al., 2023)** utilizes a PLM to predict the score of the generation result of all LLMs given a query, and then selects the LLM with the highest score within a given cost. **(5) RouteLLM (Ong et al., 2024).** Learns to route queries among a weak-strong pair of LLMs. Following the official setup, we design

the weak model as the one with the lower average win rate in the dataset, and the strong model as the one with the higher win rate.

Memory-based: (6) MA-GNN (Chen Ma, 2020). A memory-augmented GNN that models both *short-term* and *long-term* user interests through item-level message passing and a dedicated memory module. **(7) TIGER (Shashank Rajput, 2023).** A generative retrieval-based sequential recommender that models item sequences via semantic discrete codes.

4.3 Metrics

We evaluate the performance of all methods using two metrics: **(1) Accuracy** measures how often the model correctly identifies the most preferred LLM to answer a given query from a specific user. **(2) AUC-ROC** (Area Under the Receiver Operating Characteristic Curve) measures the model’s ability to correctly rank candidate LLMs according to user preferences.

4.4 Implementation Details

We employ Contriever (Izacard et al., 2021) as the PLM to obtain the initial node embeddings. We adopt the Heterogeneous Graph Transformer (HGT) (Ziniu Hu, 2020) as our model backbone due to its strong capability to maintain dedicated representations for different types of nodes. Additionally, we experiment with various other GNNs as the backbone to investigate their impact on GMTRouter’s performance. Experimental details are provided in Appendix E.4. We set the visible data size per user to $k = 10$ during both training and inference and adopt Entropy Loss as our loss function. In Appendix C, we will experimentally analyze the impact of different values of k on our method, and hyperparameter details are provided in Appendix A.1.

5 Experiment Results

5.1 Comparison with Baselines

We compare GMTRouter with multiple baselines across five datasets in Table 2. We observe that GMTRouter outperforms the strongest baseline on the vast majority of evaluation metrics, achieving up to 0.105 absolute improvement in accuracy and 0.12 in AUC, which demonstrates the superiority of our framework. For prompt-based baselines, although incorporating user interaction histories into prompts improves performance over the Vanilla

Table 2: **GMTRouter consistently outperforms baselines across all datasets.** Bold and underline denote the best and second-best results. The results are averaged over multiple runs. Since RouteLLM and FrugalGPT are inherently binary routers, we evaluated them only in the binary setting from our datasets.

Method	Chatbot-Arena		MT-Bench		GSM8K		MMLU		LaMP	
	ACC	AUC	ACC	AUC	ACC	AUC	ACC	AUC	ACC	AUC
Vanilla LLM	0.525	0.741	0.481	0.457	0.546	0.533	0.473	0.475	0.334	0.584
Personalized LLM	0.646	0.780	0.437	0.491	0.553	0.536	0.675	0.678	0.312	0.605
GraphRouter	0.771	<u>0.869</u>	0.568	0.550	0.717	0.792	0.699	0.746	0.345	0.652
FrugalGPT	0.562	0.622	0.551	0.552	0.504	0.515	0.545	0.575	-	-
RouteLLM	0.492	0.485	0.480	0.475	0.499	0.498	0.532	0.513	-	-
MA-GNN	0.673	0.775	0.679	0.739	0.636	0.648	0.702	0.758	0.347	0.661
TIGER	0.739	0.735	0.656	0.691	0.639	0.683	0.710	0.764	0.339	0.698
Ours	<u>0.774</u>	0.875	0.784	0.859	0.773	0.859	0.771	0.870	<u>0.349</u>	0.662
Ours (new user)	0.780	0.858	<u>0.759</u>	0.824	<u>0.756</u>	<u>0.833</u>	<u>0.751</u>	<u>0.831</u>	0.400	0.673

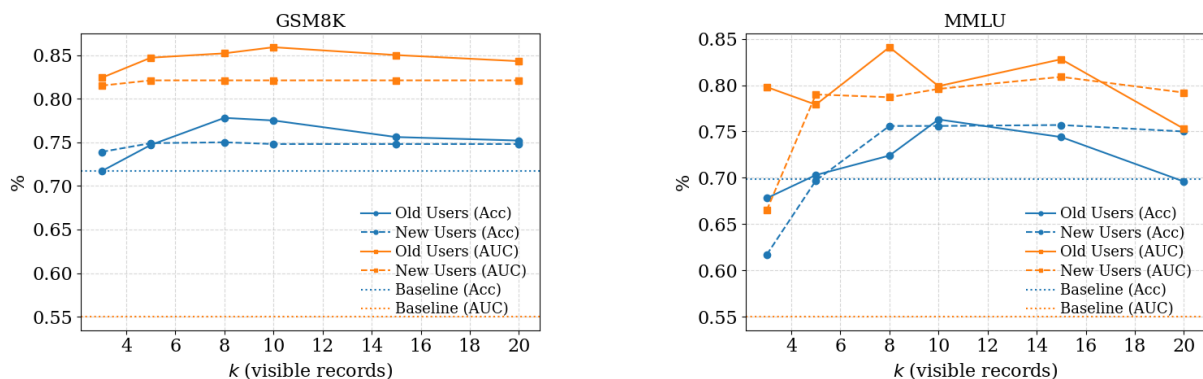


Figure 4: **Result comparison between old-user and new-user settings for GSM8K (left) and MMLU (right).** The dashed line represents the GraphRouter baseline. The personalized performance under the new-user setting is comparable to that under the old-user setting, highlighting the strong generalization capability of our method.

LLM, these methods still fall significantly behind GMTRouter. This result highlights the limited ability of LLMs to directly extract user preference patterns from raw interaction data. Moreover, our approach consistently outperforms router-based baselines, including GraphRouter, which has shown strong performance in non-personalized LLM routing tasks. These results validate the importance of leveraging structured information from user-LLM interaction data, together with explicit user preference signals, to better align model selection with diverse user needs. Compared to memory-based methods, GMTRouter benefits from learning over a fine-grained and highly structured heterogeneous graph, leading to superior performance. Furthermore, even when **30% of users are absent from the training set**, our method achieves performance comparable to the standard setting, underscoring its strong generalization capability to unseen users.

Our Framework is Lightweight With only 27.4M trainable parameters and a 109.6MB model size, our framework remains compact compared to existing routing models. During training, only

4.3GB of GPU memory is needed, making it feasible to train on a single modern GPU without specialized hardware.

5.2 Generalization to New Users

We further investigate the personalized capability of our method in few-shot scenarios with new users. Specifically, we evaluate on the GSM8K and MMLU by sampling 30% users from each dataset and varying the number of visible data $k \in \{3, 5, 8, 10, 15, 20\}$. Figure 4 presents averaged results of the sampled users under two settings: (i) the old user setting, where their records are included in the training set, and (ii) the new user setting, where they appear only in the validation and test sets. We observe that new users achieve results comparable to old users, and their performance curves consistently peak far above the GraphRouter baseline. These findings demonstrate that **our approach effectively learns to capture user preferences from few-shot data and can adapt to new users without requiring extensive fine-tuning.**

Table 3: **Ablation of design components.** We compare the full model with three variants: (1) removing the user preference features, (2) replacing the prediction head with a dot-product, (3) not using user embeddings during prediction. The best and second-best results are highlighted in **bold** and underline, respectively.

Method	Chatbot-Arena		MT-Bench		GSM8K		MMLU	
	Accuracy	AUC	Accuracy	AUC	Accuracy	AUC	Accuracy	AUC
w/o h_p	0.768	0.872	0.569	0.507	0.715	0.784	0.494	0.613
Dot-product	0.777	0.868	<u>0.730</u>	<u>0.795</u>	0.629	0.724	0.681	0.746
w/o h_u	0.771	<u>0.873</u>	0.569	0.631	<u>0.725</u>	<u>0.814</u>	<u>0.701</u>	<u>0.771</u>
GMTRouter	<u>0.774</u>	0.875	0.784	0.859	0.773	0.859	0.771	0.870

5.3 Ablation Studies

To evaluate the effectiveness of each design component of the GMTRouter, we conduct ablation studies along the following aspects.

- **w/o User Preference Feature** To verify the effectiveness of the user preference feature in propagating preference signals during GNN aggregation, we remove this feature in this variant. As a result, node embeddings are updated without incorporating preference ratings.
- **Dot-product Prediction Head** To evaluate whether the cross-attention prediction head captures non-linear interactions more effectively than standard similarity scoring when predicting the optimal model, we replace it in this variant with a simple dot product between the (user + query) and LLM embeddings.
- **w/o User Embedding** To evaluate the effectiveness of user embeddings aggregated from the sampled visible graph for personalized prediction, we replace the user embeddings fed into the prediction head with zero vectors, thereby ablating their influence on the predictions.

The results of our ablation studies are presented in Table 3. As shown, our GMTRouter achieves the best performance on most metrics across all four datasets compared to the other variants, confirming the effectiveness of our design choices.

5.4 Further Experiments

We present additional experimental results in Appendix E to further examine the robustness and generality of GMTRouter. Specifically, we evaluate its generalization performance on a larger set of users from ChatBot Arena (Appendix E.1), assess its robustness under varying levels of noise in the interaction data (Appendix E.2), and compare GMTRouter with personalized generation methods (Appendix E.3).

6 Additional Related Works

LLM Routing. LLM routing focuses on enhancing inference efficiency and response quality by

assigning queries to the most appropriate model (Yue et al., 2025; Zhang et al., 2025c). Recent work frames routing as learning with cost-quality tradeoffs (Kadavath et al., 2022; Dekoninck et al., 2024): RouteLLM learns from preference data (Ong et al., 2024), and RouterBench offers standardized routing benchmarks (Hu et al., 2024). BEST-Route jointly selects LLM and generation count at test-time via a bandit controller (Ding et al., 2025). However, existing approaches are not fully personalized and fail to exploit user information from interaction histories as well as the structure of multi-turn dialogues.

Heterogeneous Graph Learning. HetGNNs are designed to model heterogeneous graphs by capturing complex multi-type interactions among various nodes and edges (Chien et al., 2021; Feng et al., 2019). HAN uses hierarchical attention over metapaths (Wang et al., 2019), while MAGNN and HeCo improve metapath aggregation and cross-view contrast (Fu et al., 2020; Wang et al., 2021). Transformers such as HGT provide inductive, relation-aware message passing with temporal encoding (Ziniu Hu, 2020). This enables rich relational structures in user-LLM interactions while leveraging inductive training to enhance generalization on sparse data from new users.

7 Conclusion

In this work, we introduced GMTRouter, a heterogeneous graph-based framework for personalized LLM routing. By modeling multi-turn user-LLM interactions as a heterogeneous graph and propagating preference signals across node types, our method effectively captures user-specific patterns even from few-shot data. Experiments across five benchmarks confirm that GMTRouter consistently surpasses strong baselines, while adapting efficiently to new users without retraining. These results highlight the value of structured interaction modeling for advancing preference-aware LLM routing and point to promising future directions in scalable, user-aligned LLM deployment.

8 Limitations

While GMTRouter demonstrates significant improvements in personalized routing, we identify a few limitations that offer avenues for future research. First, the initial embeddings for LLM nodes are currently derived from static descriptions provided by API platforms. While effective, these representations may require manual updates if a model undergoes significant versioning or performance shifts. Second, although we evaluated our framework across five diverse benchmarks spanning reasoning, math, and general knowledge, its performance in highly specialized niche domains, such as advanced medical or legal sub-specialties, remains to be further explored. Finally, the current implementation focuses on individual user personalization; extending the framework to model collective or group-level preferences in collaborative environments represents a promising direction for future work.

9 Ethical Considerations

Our work does not pose significant ethical risks or involve sensitive data collection. The experiments were conducted using anonymized, publicly available datasets such as ChatBot Arena ensuring that no personally identifiable information (PII) was processed or compromised. Furthermore, GMTRouter serves as an algorithmic optimization layer for model selection rather than a content generation engine; as such, it does not generate new text or introduce additional biases beyond those already inherent in the underlying candidate LLMs. Our use of an inductive learning framework further supports privacy-preserving principles, as it enables effective personalization from minimal, local interaction histories rather than requiring the construction of extensive, long-term global user profiles.

References

2024. Qwen2 technical report.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

AI@Meta. 2024. *Llama 3 model card*.

Laith Alzubaidi, Jinshuai Bai, Aiman Al-Sabaawi, José I. Santamaría, A. Albahri, B. S. Al-dabbagh,

M. Fadhel, M. Manoufali, Jinglan Zhang, Ali H. Altimemy, Ye Duan, Amjed Abdullah, Laith Farhan, Yi Lu, Ashish Gupta, Felix Albu, Amin Abbosh, and Yuantong Gu. 2023. A survey on deep learning tools dealing with data scarcity: definitions, challenges, solutions, tips, and applications. *Journal of Big Data*, 10:1–82.

Steven Au, Cameron J Dimacali, Ojasmitha Pedirappagari, Namyong Park, Franck Dernoncourt, Yu Wang, Nikos Kanakaris, Hanieh Deilamsalehy, Ryan A Rossi, and Nesreen K Ahmed. 2025. Personalized graph-based retrieval for large language models. *arXiv preprint arXiv:2501.02157*.

T. Banditwattanawong and Masawee Masdisornchote. 2025. Unbiased machine learning-assisted approach for conditional discretization of human performances. *PeerJ Comput. Sci.*, 11:e2804.

B. Cavallo. 2019. Functional relations and spearman correlation between consistency indices. *Journal of the Operational Research Society*, 71:301 – 311.

Lingjiao Chen, Matei Zaharia, and James Zou. 2023. Frugalgpt: How to use large language models while reducing cost and improving performance. *arXiv preprint arXiv:2305.05176*.

Yingxue Zhang Jianing Sun Xue Liu Mark Coates Chen Ma, Liheng Ma. 2020. Memory augmented graph neural networks for sequential recommendation. *AAAI 2020*.

Rendi Chevi, Kentaro Inui, T. Solorio, and Alham Fikri Aji. 2025. How individual traits and language styles shape preferences in open-ended user-llm interaction: A preliminary study. *ArXiv*, abs/2504.17083.

Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph Gonzalez, and Ion Stoica. 2024. Chatbot arena: An open platform for evaluating llms by human preference. *ArXiv*, abs/2403.04132.

Eli Chien, Chao Pan, Jianhao Peng, and Olgica Milenkovic. 2021. You are allset: A multiset function framework for hypergraph neural networks. *arXiv preprint arXiv:2106.13264*. ICLR 2022.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

Xiangxiang Dai, Jin Li, Xutong Liu, Anqi Yu, and John Lui. 2024. Cost-effective online multi-llm selection with versatile reward models. *arXiv preprint arXiv:2405.16587*.

Joost CF De Winter, Samuel D Gosling, and Jeff Potter. 2016. Comparing the pearson and spearman correlation coefficients across distributions and sample

713	sizes: A tutorial using simulations and empirical data.	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	765
714	<i>Psychological methods</i> , 21(3):273.	Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle,	766
715	DeepSeek-AI. 2024. Deepseek-v3 technical report .	Aiesha Letman, Akhil Mathur, Alan Schelten,	767
716	<i>Preprint</i> , arXiv:2412.19437.	Alex Vaughan, and 1 others. 2024. The llama 3 herd	768
717	Jasper Dekoninck, Maximilian Baader, and Martin	of models. <i>arXiv preprint arXiv:2407.21783</i> .	769
718	Vechev. 2024. A unified approach to routing and cascading for llms .	William L. Hamilton, Z. Ying, and J. Leskovec. 2017.	770
719	<i>arXiv preprint arXiv:2410.10347</i> .	Inductive representation learning on large graphs.	771
720	Dujian Ding, Ankur Mallick, Chi Wang, Robert Sim,	<i>ArXiv</i> , abs/1706.02216.	772
721	Subhabrata Mukherjee, Victor Ruhle, Laks VS Laksh-	Jan Hauke and Tomasz Kossowski. 2011. Comparison	773
722	manan, and Ahmed Hassan Awadallah. 2024. Hybrid	of values of pearson’s and spearman’s correlation	774
723	llm: Cost-efficient and quality-aware query routing.	coefficients on the same sets of data. <i>Quaestiones</i>	775
724	<i>arXiv preprint arXiv:2404.14618</i> .	<i>geographicae</i> , 30(2):87–93.	776
725	Dujian Ding, Ankur Mallick, Shaokun Zhang, Chi	Dan Hendrycks, Collin Burns, Steven Basart, Andrew	777
726	Wang, Daniel Madrigal, Mirian Del Carmen Hipolito	Critch, Jerry Li, Dawn Song, and Jacob Steinhardt.	778
727	Garcia, Menglin Xia, Laks V. S. Lakshmanan,	2021a. Aligning ai with shared human values. <i>Pro-</i>	779
728	Qingyun Wu, and Victor Rühle. 2025. Best-route: Adaptive llm routing with test-time optimal compute .	<i>ceedings of the International Conference on Learning</i>	780
729	In <i>Proceedings of the 42nd International Conference</i>	<i>Representations (ICLR)</i> .	781
730	<i>on Machine Learning (ICML)</i> . Also available as	Dan Hendrycks, Collin Burns, Steven Basart, Andy	782
731	arXiv:2506.22716.	Zou, Mantas Mazeika, Dawn Song, and Jacob Stein-	783
732	Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan	hardt. 2021b. Measuring massive multitask language	784
733	Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu,	understanding. <i>Proceedings of the International Con-</i>	785
734	Tianyu Liu, and 1 others. 2022. A survey on in-	<i>ference on Learning Representations (ICLR)</i> .	786
735	context learning. <i>arXiv preprint arXiv:2301.00234</i> .	Qijun Hu, Rui Zhang, Wenxuan Ren, Haoran Zhang,	787
736	Guillaume Escamocher, Samira Pourkhajouei, Federico	Minjia Zhang, Xinyu Zhou, Tong Liu, Pengfei Liu,	788
737	Toffano, Paolo Viappiani, and Nic Wilson. 2024. Interactive preference elicitation under noisy preference models: An efficient non-bayesian approach .	Tong Zhang, and Mu Li. 2024. Routerbench: A benchmark for multi-llm routing system .	789
738	<i>Int. J. Approx. Reason.</i> , 178:109333.	<i>arXiv</i>	790
739	Tao Feng, Yanzhen Shen, and Jiaxuan You. 2024.	<i>preprint arXiv:2403.12031</i> .	791
740	Graphrouter: A graph-based router for llm selections.	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Se-	792
741	<i>arXiv preprint arXiv:2410.03834</i> .	bastian Riedel, Piotr Bojanowski, Armand Joulin,	793
742	Tao Feng, Haozhen Zhang, Zijie Lei, Pengrui Han,	and Edouard Grave. 2021. Unsupervised dense infor-	794
743	Mostofa Patwary, Mohammad Shoeybi, Bryan Catan-	mation retrieval with contrastive learning .	795
744	zaro, and Jiaxuan You. 2025. Fusing llm capabilities	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	796
745	with routing data. <i>arXiv preprint arXiv:2507.10540</i> .	sch, Chris Bamford, Devendra Singh Chaplot, Diego	797
746	Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	798
747	Ji, and Yue Gao. 2019. Hypergraph neural networks .	laume Lample, Lucile Saulnier, L�elio Renard Lavaud,	799
748	In <i>Proceedings of the AAAI Conference on Artificial</i>	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	800
749	<i>Intelligence (AAAI)</i> .	Thibaut Lavril, Thomas Wang, Timoth�ee Lacroix,	801
750	Matthias Fey and Jan E. Lenssen. 2019. Fast graph	and William El Sayed. 2023. Mistral 7b . <i>Preprint</i> ,	802
751	representation learning with PyTorch Geometric.	arXiv:2310.06825.	803
752	In <i>ICLR Workshop on Representation Learning on</i>	Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux,	804
753	<i>Graphs and Manifolds</i> .	A. Mensch, Blanche Savary, Chris Bamford, Deven-	805
754	Xinyu Fu, Jiani Zhang, Ziqiao Meng, and Irwin King.	dra Singh Chaplot, Diego de Las Casas, Emma Bou	806
755	2020. Magnn: Metapath aggregated graph neural network for heterogeneous graph embedding .	Hanna, Florian Bressand, Gianna Lengyel, Guil-	807
756	In <i>Proceedings of The Web Conference (WWW)</i> .	laume Bour, Guillaume Lample, L’elio Renard	808
757	Ge Gao, Alexey Taymanov, Eduardo Salinas, Paul	Lavaud, Lucile Saulnier, M. Lachaux, Pierre Stock,	809
758	Mineiro, and Dipendra Misra. 2024. Aligning llm agents by learning latent preference from user edits .	Sandeep Subramanian, Sophia Yang, and 7 others.	810
759	<i>ArXiv</i> , abs/2404.15269.	2024. Mixtral of experts . <i>ArXiv</i> , abs/2401.04088.	811
760	Xinyu Fu, Jiani Zhang, Ziqiao Meng, and Irwin King.	Bowen Jiang, Zhuoqun Hao, Young-Min Cho, Bryan Li,	812
761	2020. Magnn: Metapath aggregated graph neural network for heterogeneous graph embedding .	Yuan Yuan, Sihao Chen, Lyle Ungar, C. J. Taylor, and	813
762	In <i>Proceedings of The Web Conference (WWW)</i> .	Dan Roth. 2025. Know me, respond to me: Benchmarking llms for dynamic user profiling and personalized responses at scale .	814
763	Ge Gao, Alexey Taymanov, Eduardo Salinas, Paul	<i>ArXiv</i> , abs/2504.14225.	815
764	Mineiro, and Dipendra Misra. 2024. Aligning llm agents by learning latent preference from user edits .	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom	816
765	<i>ArXiv</i> , abs/2404.15269.	Henighan, Dawn Drain, Ethan Perez, Nicholas	817
766	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli	818
767	Abhinav Pandey, Abhishek Kadian, Ahmad Al-		819
768	Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,		
769	Alex Vaughan, and 1 others. 2024. The llama 3 herd		
770	of models. <i>arXiv preprint arXiv:2407.21783</i> .		
771	William L. Hamilton, Z. Ying, and J. Leskovec. 2017.		
772	Inductive representation learning on large graphs.		
773	<i>ArXiv</i> , abs/1706.02216.		
774	Jan Hauke and Tomasz Kossowski. 2011. Comparison		
775	of values of pearson’s and spearman’s correlation		
776	coefficients on the same sets of data. <i>Quaestiones</i>		
777	<i>geographicae</i> , 30(2):87–93.		
778	Dan Hendrycks, Collin Burns, Steven Basart, Andrew		
779	Critch, Jerry Li, Dawn Song, and Jacob Steinhardt.		
780	2021a. Aligning ai with shared human values. <i>Pro-</i>		
781	<i>ceedings of the International Conference on Learning</i>		
782	<i>Representations (ICLR)</i> .		
783	Dan Hendrycks, Collin Burns, Steven Basart, Andy		
784	Zou, Mantas Mazeika, Dawn Song, and Jacob Stein-		
785	hardt. 2021b. Measuring massive multitask language		
786	understanding. <i>Proceedings of the International Con-</i>		
787	<i>ference on Learning Representations (ICLR)</i> .		
788	Qijun Hu, Rui Zhang, Wenxuan Ren, Haoran Zhang,		
789	Minjia Zhang, Xinyu Zhou, Tong Liu, Pengfei Liu,		
790	Tong Zhang, and Mu Li. 2024. Routerbench: A benchmark for multi-llm routing system .		
791	<i>arXiv</i>		
792	<i>preprint arXiv:2403.12031</i> .		
793	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Se-		
794	bastian Riedel, Piotr Bojanowski, Armand Joulin,		
795	and Edouard Grave. 2021. Unsupervised dense infor-		
796	mation retrieval with contrastive learning .		
797	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-		
798	sch, Chris Bamford, Devendra Singh Chaplot, Diego		
799	de las Casas, Florian Bressand, Gianna Lengyel, Guil-		
800	laume Lample, Lucile Saulnier, L�elio Renard Lavaud,		
801	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,		
802	Thibaut Lavril, Thomas Wang, Timoth�ee Lacroix,		
803	and William El Sayed. 2023. Mistral 7b . <i>Preprint</i> ,		
804	arXiv:2310.06825.		
805	Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux,		
806	A. Mensch, Blanche Savary, Chris Bamford, Deven-		
807	dra Singh Chaplot, Diego de Las Casas, Emma Bou		
808	Hanna, Florian Bressand, Gianna Lengyel, Guil-		
809	laume Bour, Guillaume Lample, L’elio Renard		
810	Lavaud, Lucile Saulnier, M. Lachaux, Pierre Stock,		
811	Sandeep Subramanian, Sophia Yang, and 7 others.		
812	2024. Mixtral of experts . <i>ArXiv</i> , abs/2401.04088.		
813	Bowen Jiang, Zhuoqun Hao, Young-Min Cho, Bryan Li,		
814	Yuan Yuan, Sihao Chen, Lyle Ungar, C. J. Taylor, and		
815	Dan Roth. 2025. Know me, respond to me: Benchmarking llms for dynamic user profiling and personalized responses at scale .		
816	<i>ArXiv</i> , abs/2504.14225.		
817	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom		
818	Henighan, Dawn Drain, Ethan Perez, Nicholas		
819	Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli		

820	Tran-Johnson, and 1 others. 2022. Language models (mostly) know what they know . <i>arXiv preprint arXiv:2207.05221</i> .	Alireza Salemi, Surya Kallumadi, and Hamed Zamani. 2024a. Optimization methods for personalizing large language models through retrieval augmentation. In <i>Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 752–762.	873
821			874
822			875
823	N. Keskar, Dheevatsa Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang. 2016. On large-batch training for deep learning: Generalization gap and sharp minima. <i>ArXiv</i> , abs/1609.04836.	Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. 2024b. Lamp: When large language models meet personalization. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7370–7392.	876
824			877
825			878
826			
827	Haoxuan Li, Chunyuan Zheng, Wenjie Wang, Hao Wang, Fuli Feng, and Xiao-Hua Zhou. 2024a. De-biased recommendation with noisy feedback . <i>Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> .		879
828			880
829			881
830			882
831			883
832	Xinyu Li, Z. Lipton, and Liu Leqi. 2024b. Personalized language modeling from personalized human feedback . <i>ArXiv</i> , abs/2402.05133.	M. Schlichtkrull, Thomas Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and M. Welling. 2017. Modeling relational data with graph convolutional networks . pages 593–607.	884
833			885
834			886
835	Ying Li, Ye Zhong, Lijuan Yang, Yanbo Wang, and Penghua Zhu. 2025a. Llm-guided crowdsourced test report clustering . <i>IEEE Access</i> , 13:24894–24904.	Anima Singh Raghunandan Hulikal Keshavan Trung Vu Lukasz Heldt Lichan Hong Yi Tay Vinh Tran Jonah Samost Maciej Kula Ed Chi Maheswaran Sathiamoorthy Shashank Rajput, Nikhil Mehta. 2023. Recommender systems with generative retrieval . <i>NeurIPS 2023</i> .	887
836			888
837			889
838	Yubo Li, Xiaobin Shen, Xinyu Yao, Xueying Ding, Yidi Miao, R. Krishnan, and R. Padman. 2025b. Beyond single-turn: A survey on multi-turn interactions with large language models . <i>ArXiv</i> , abs/2504.04717.		890
839			891
840			892
841			893
842	Weiwei Lin, Hao Xu, Jianzhuo Li, Ziming Wu, Zhengyang Hu, Victor I. Chang, and J. Wang. 2021. Deep-profiling: a deep neural network model for scholarly web user profiling . <i>Cluster Computing</i> , 26:1753 – 1766.	Taiwei Shi, Zhuoer Wang, Longqi Yang, Ying-Chun Lin, Zexue He, Mengting Wan, Pei Zhou, S. Jauhar, Xiaofeng Xu, Xia Song, and Jennifer Neville. 2024. Wildfeedback: Aligning llms with in-situ user interactions and feedback . <i>ArXiv</i> , abs/2408.15549.	894
843			895
844			896
845			897
846			898
847	Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. 2022. Biogpt: generative pre-trained transformer for biomedical text generation and mining . <i>Briefings in bioinformatics</i> , 23(6):bbac409.	Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, and 1 others. 2023. Large language models encode clinical knowledge. <i>Nature</i> , 620(7972):172–180.	899
848			900
849			901
850			902
851			903
852	Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E. Gonzalez, Mohammed Kadous, and Ion Stoica. 2024. Routellm: Learning to route llms with preference data . <i>arXiv preprint arXiv:2406.18665</i> .	Robyn Speer. 2022. rspeer/wordfreq: v3.0 .	904
853			905
854			906
855			907
856			908
857	OpenAI. 2025. Gpt-5 system card . Technical report, OpenAI.	Dimitris Stripelis, Zijian Hu, Jipeng Zhang, Zhaozhuo Xu, Alay Shah, Han Jin, Yuhang Yao, Salman Avestimehr, and Chaoyang He. 2024. Polyrouter: A multi-llm querying system . <i>arXiv e-prints</i> , pages arXiv–2408.	909
858			910
859	O. Oyedotun, Konstantinos Papadopoulos, and D. Aouada. 2022. A new perspective for understanding generalization gap of deep neural networks trained with large batch sizes . <i>Applied Intelligence</i> , 53:15621–15637.	Hongzu Su, Jingjing Li, Zhekai Du, Lei Zhu, Ke Lu, and H. Shen. 2024. Cross-domain recommendation via dual adversarial adaptation . <i>ACM Transactions on Information Systems</i> , 42:1 – 26.	911
860			912
861			913
862			914
863			
864	Marija Šakota, Maxime Peyrard, and Robert West. 2024. Fly-swat or cannon? cost-effective language model choice via meta-modeling . In <i>Proceedings of the 17th ACM International Conference on Web Search and Data Mining</i> , pages 606–615.	Yihang Sun, Tao Feng, Ge Liu, and Jiaxuan You. 2025. Premium: Llm personalization with individual-level preference feedback . <i>ArXiv</i> .	915
865			916
866			917
867			
868			918
869	Sogand Salehi, Mahdi Shafiei, Teresa Yeo, Roman Bachmann, and Amir Zamir. 2024. Viper: Visual personalization of generative models via individual preference learning . pages 391–406.	Gemma Team. 2024. Gemma .	919
870			920
871			921
872			922
			923
			924
			925
			926

927	Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2022. Text embeddings by weakly-supervised contrastive pre-training. <i>arXiv preprint arXiv:2212.03533</i> .	982	Yihan Zhang, Kai Wang, Zexuan Li, Wenqi Xu, Hao-ran Zhu, and Wei Chen. 2025c. <i>Mixllm: Dynamic routing in mixed large language models</i> . In <i>Proceedings of the 2025 Conference of the North American Chapter of the ACL (NAACL)</i> .	983
928		984		985
929		986		987
930				988
931				989
932	Xiao Wang, Xiangnan He, Yuesong Cao, Meng Liu, and Tat-Seng Chua. 2021. <i>Self-supervised heterogeneous graph neural network with co-contrastive learning</i> . In <i>Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining (KDD)</i> .			990
933				991
934				992
935				993
936				994
937	Xiao Wang, Houye Ji, Chuan Shi, Bai Wang, Peng Cui, Philip S. Yu, and Yanfang Ye. 2019. <i>Heterogeneous graph attention network</i> . In <i>Proceedings of The Web Conference (WWW)</i> .			995
938				996
939				997
940				998
941	Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2023b. <i>Mint: Evaluating llms in multi-turn interaction with tools and language feedback</i> . <i>ArXiv</i> , abs/2309.10691.			999
942				1000
943				1001
944				1002
945	Zhaoyang Wang, Li Li, Ketai He, and Zhenyang Zhu. 2025. <i>User profile construction based on high-dimensional features extracted by stacking ensemble learning</i> . <i>Applied Sciences</i> .			1003
946				1004
947				1005
948				1006
949	Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy Zhang, Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev. 2024b. <i>Helpsteer2: Open-source dataset for training top-performing reward models</i> . <i>ArXiv</i> , abs/2406.08673.			1007
950				1008
951				1009
952				1010
953				1011
954				1012
955	Zhilin Wang, Yi Dong, Jiaqi Zeng, Virginia Adams, Makesh Narsimhan Sreedhar, Daniel Egert, Olivier Delalleau, Jane Polak Scowcroft, Neel Kant, Aidan Swope, and 1 others. 2023c. <i>Helpsteer: Multi-attribute helpfulness dataset for steerlm</i> . <i>arXiv preprint arXiv:2311.09528</i> .			1013
956				1014
957				1015
958				1016
959				1017
960				1018
961	Hongyu Yang, Liyang He, Min Hou, Shuanghong Shen, Rui Li, Jiahui Hou, Jianhui Ma, and Junda Zhao. 2024. <i>Aligning llms through multi-perspective user preference ranking-based feedback for programming question answering</i> . <i>ArXiv</i> , abs/2406.00037.			1019
962				1020
963				1021
964				1022
965				1023
966	Yanwei Yue, Guibin Zhang, Boyang Liu, and 1 others. 2025. <i>Masrouter: Learning to route llms for multi-agent systems</i> . In <i>Proceedings of the 63rd Annual Meeting of the ACL</i> .			1024
967				1025
968				1026
969				1027
970	Hang Zeng, Chaoyue Niu, Fan Wu, Chengfei Lv, and Guihai Chen. 2024. <i>Personalized llm for generating customized responses to the same query from different users</i> . <i>ArXiv</i> , abs/2412.11736.			1028
971				1029
972				1030
973				1031
974	Chen Zhang, Xinyi Dai, Yaxiong Wu, Qu Yang, Yasheng Wang, Ruiming Tang, and Yong Liu. 2025a. <i>A survey on multi-turn interaction capabilities of large language models</i> . <i>ArXiv</i> , abs/2501.09959.			1032
975				1033
976				1034
977				1035
978	Chi Zhang, Junho Jeong, and Jin-Woo Jung. 2025b. <i>Anomaly detection over multi-relational graphs using graph structure learning and multi-scale meta-path graph aggregation</i> . <i>IEEE Access</i> , 13:60303–60316.			1036
979				1037
980				1038
981				1039
				1040
				1041
				1042
				1043
				1044
				1045
				1046
				1047
				1048
				1049
				1050
				1051
				1052
				1053
				1054
				1055
				1056
				1057
				1058
				1059
				1060
				1061
				1062
				1063
				1064
				1065
				1066
				1067
				1068
				1069
				1070
				1071
				1072
				1073
				1074
				1075
				1076
				1077
				1078
				1079
				1080
				1081
				1082
				1083
				1084
				1085
				1086
				1087
				1088
				1089
				1090
				1091
				1092
				1093
				1094
				1095
				1096
				1097
				1098
				1099
				1100
				1101
				1102
				1103
				1104
				1105
				1106
				1107
				1108
				1109
				1110
				1111
				1112
				1113
				1114
				1115
				1116
				1117
				1118
				1119
				1120
				1121
				1122
				1123
				1124
				1125
				1126
				1127
				1128
				1129
				1130
				1131
				1132
				1133
				1134
				1135
				1136
				1137
				1138
				1139
				1140
				1141
				1142
				1143
				1144
				1145
				1146
				1147
				1148
				1149
				1150
				1151
				1152
				1153
				1154
				1155
				1156
				1157
				1158
				1159
				1160
				1161
				1162
				1163
				1164
				1165
				1166
				1167
				1168
				1169
				1170
				1171
				1172
				1173
				1174
				1175
				1176
				1177
				1178
				1179
				1180
				1181
				1182
				1183
				1184
				1185
				1186
				1187
				1188
				1189
				1190
				1191
				1192
				1193
				1194
				1195
				1196
				1197
				1198
				1199
				1200

B Dataset Preparation

B.1 Details for datasets

We evaluate our approach on both real-world and synthetic datasets spanning five distinct tasks to provide a comprehensive assessment.

- **Chatbot Arena (Chiang et al., 2024):** ChatBot Arena is a real-world dataset of anonymized multi-turn conversations with pairwise human preference labels across LLMs. For our experiments, we select the 11 users and 16 LLMs with the largest numbers of interactions. Detailed statistics are provided in Appendix B.2.
- **MT-Bench (Zheng et al., 2023):** MT-Bench is a benchmark for evaluating the reasoning and multi-turn conversational capabilities of LLMs, containing 80 multi-turn questions.
- **GSM8K (Cobbe et al., 2021):** GSM8K is a dataset of grade school-level math word problems, designed to assess LLMs’ mathematical reasoning and problem-solving skills.
- **MMLU (Hendrycks et al., 2021a,b):** MMLU is a comprehensive benchmark covering 57 subjects from professional domains, used to measure general knowledge and multi-domain reasoning abilities of LLMs. We sample 10 questions from each subject for our experiments.
- **LaMP (Salemi et al., 2024b):** LaMP is designed to evaluate language models across multiple dimensions of personalization. We select the "Personalized Scholarly Title Generation" task, which provides pairs of paper titles and abstracts for multiple users and requires predicting the title a user would prefer given an abstract. We convert this task into a personalized routing dataset, with processing details provided in Appendix B.4.

B.2 Dataset Statistics

We preprocess each dataset by extracting user–query–LLM–response tuples and partition them into train, validation, and test sets. To ensure fair evaluation and meaningful personalization, we stratify the splits to maintain balanced user–model preference distributions and avoid degenerate cases (e.g., users consistently preferring a single LLM or lacking query diversity). This setup promotes generalization under cold-start conditions and supports robust evaluation of routing behavior.

For ChatBot Arena, we selected the following users and LLMs:

Users: arena_user_9965, arena_user_15085, arena_user_257, arena_user_13046,

Table 4: Dataset statistics, including the number of entries, users, and LLMs in each split.

Dataset	Split	Entries	Users	LLMs
ChatBot Arena	Train	2780	11	16
	Valid	386	11	16
	Test	824	11	16
MT Bench	Train	2240	10	2
	Valid	320	10	2
	Test	640	10	2
GSM8K	Train	18460	10	2
	Valid	2620	10	2
	Test	5300	10	2
MMLU	Train	3970	5	2
	Valid	560	5	2
	Test	1150	5	2
LaMP	Train	6895	10	5
	Valid	985	10	5
	Test	1970	10	5

arena_user_11473,
arena_user_9676,
arena_user_6585,
arena_user_1338

arena_user_3820,
arena_user_6467,
arena_user_5203,

1066
1067
1068
1069

LLMs: koala-13b, vicuna-13b, gpt-3.5-turbo, oasst-pythia-12b, gpt-4, claude-v1, RWKV-4-Raven-14B, palm-2, alpaca-13b, mpt-7b-chat, vicuna-7b, claude-instant-v1, chatglm-6b, fastchat-t5-3b, dolly-v2-12b, stablelm-tuned-alpha-7b

1070
1071
1072
1073
1074

We report in Table 5 the size of each dataset along with the time required to process it into the heterogeneous graph used in our experiments.

1075
1076
1077

B.3 Synthetic User Design

To simulate diverse user preferences, we introduce synthetic users whose routing behavior is governed by a weighted linear utility function over multiple metrics: human preference rating, token count, output diversity, and cost. For each dataset, we manually assign different weights $\{w_{\text{rating}}, w_{\text{tokens}}, w_{\text{diff}}, w_{\text{cost}}\}$ per user to reflect individualized trade-offs, such as favoring cost-efficiency or output diversity over raw model quality. These weights are normalized within each dataset to prevent scale bias.

1078

1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089

B.4 Processing of the LaMP Dataset

We select the "Personalized Scholarly Title Generation" task from the LaMP benchmark (Salemi

1090
1091
1092

Table 5: Computational cost of graph construction across datasets.

Dataset	Data Entries	Avg. Tokens	Encoding Time (s)	Graph Construction Time (s)
ChatBot-Arena	3990	184.41	51.73	1.70
MT-Bench	3200	4511.73	55.68	2.40
GSM8K	26380	112.68	142.84	1.49
MMLU	5680	9.35	4.27	1.56
LaMP	9850	66.82	30.98	1.91

Table 6: Synthetic user weights for MT-Bench dataset.

User	w_{rating}	w_{tokens}	w_{diff}	w_{cost}
user_1	1.42	0.0087	-0.174	-45.23
user_2	1.87	0.0012	0.091	-15.55
user_3	0.96	0.0135	0.045	-48.42
user_4	1.15	-0.0008	-0.220	-10.00
user_5	1.69	0.0024	0.175	-38.50
user_6	1.08	-0.0015	-0.030	-25.12
user_7	0.53	0.0162	0.230	-5.75
user_8	1.34	-0.0005	-0.145	-12.40
user_9	1.98	0.0101	0.087	-25.10
user_10	1.57	0.0024	-0.065	-7.79

Table 7: Synthetic user weights for GSM8K dataset.

User	w_{rating}	w_{tokens}	w_{diff}	w_{cost}
user_1	1.0	20.0	100.0	-0.0
user_2	1.5	18.0	50.0	-1.0
user_3	0.8	22.0	80.0	-0.5
user_4	1.2	17.0	120.0	-0.2
user_5	2.0	15.0	70.0	-0.4
user_6	0.4	6.0	-4.0	-1.0
user_7	0.3	7.0	-5.0	-0.9
user_8	0.6	8.0	-7.0	-1.2
user_9	0.2	9.0	-9.0	-0.8
user_10	0.8	10.0	-3.0	-1.1

Table 8: Synthetic user weights for MMLU dataset.

User	w_{rating}	w_{tokens}	w_{diff}	w_{cost}
user_1	1.0	0.00	0.00	0.0
user_2	1.0	0.00	0.00	-600.0
user_3	1.0	0.00	0.00	-1200.0
user_4	1.0	0.00	0.00	-1800.0
user_5	1.0	0.00	0.00	-2400.0

et al., 2024b) as our new dataset. This task provides pairs of paper titles and abstracts for multiple users and requires predicting the title a user would prefer given an abstract.

Data extraction. We identify the 10 users with the largest amount of data and randomly sample 200 (title, abstract) pairs for each user.

LLM response generation. We use five LLMs with diverse architectures and sizes—deepseek-r1 (DeepSeek-AI, 2024), gemma-2-27b-it (Team, 2024), llama-3.1-8b-instruct (AI@Meta, 2024), qwen2-7b-instruct (qwe, 2024), and mistral-7b-instruct-v0.3 (Jiang et al., 2023)—to generate a predicted title for each abstract.

User rating acquisition. For each paper, we encode both the ground-truth title and all LLM-generated titles using a PLM. We compute the cosine similarity between a generated title and the ground-truth title and treat this score as the user rating.

Dataset filtering and splitting. For each abstract, we identify the LLM with the highest user rating and use it as the routing target, discarding samples

where ties occur. This yields a total of 9,850 instances, which we split into training, validation, and test sets using a 7:1:2 ratio.

C Investigating the Impact of Visible Data Size

We investigate the impact of k visible data per user on the quality of the aggregated node embeddings. Figures 5 and 6 present the test results on four datasets, respectively. As k increases, both accuracy and AUC improve, but beyond $k=10$, the performance begins to plateau or slightly decline, indicating diminishing returns from including additional visible data. This may be due to reduced generalization or potential instability caused by excessively large batch sizes during training (Keskar et al., 2016; Oyedotun et al., 2022). Therefore, we

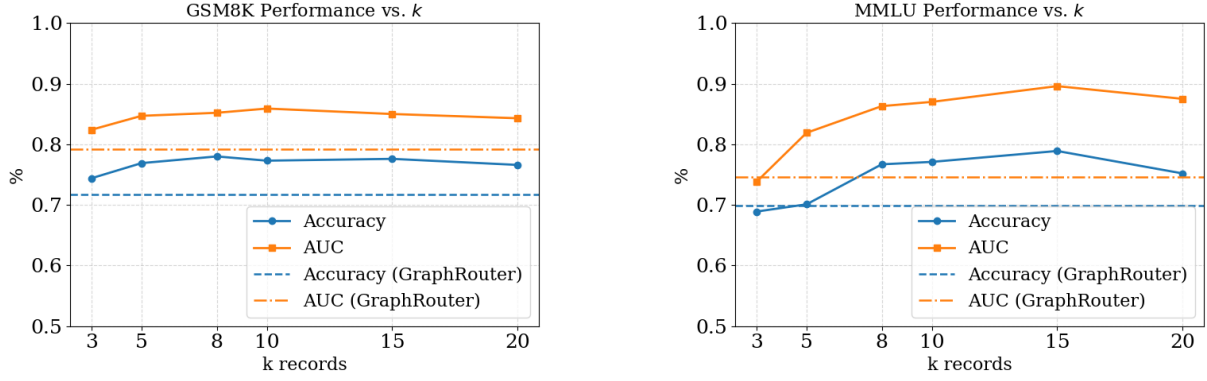


Figure 5: This figure illustrates **the impact of the visible data size k on GMTRouter for GSM8K (left) and MMLU (right)**. The dashed line represents the GraphRouter baseline. As k increases, the performance of our method improves, but it saturates once k reaches 10.

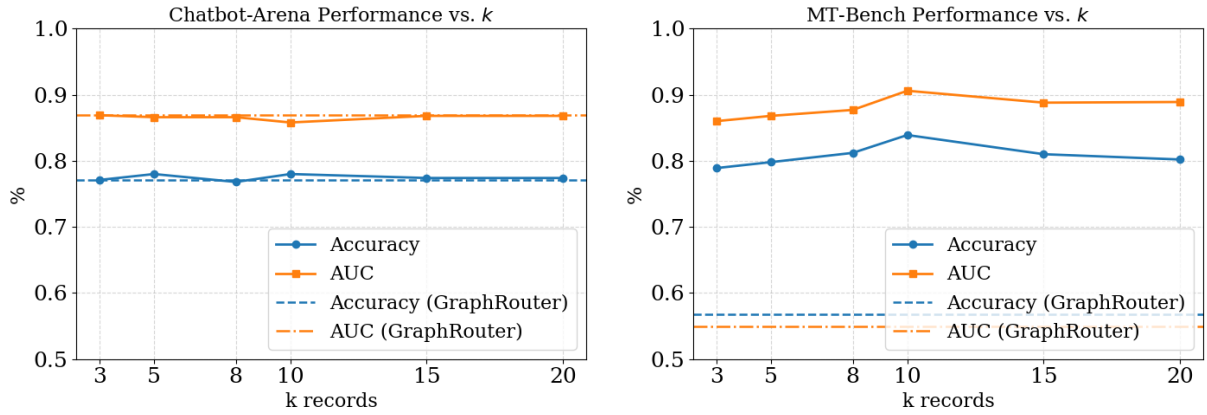


Figure 6: This figure illustrates **the impact of the visible data size k on GMTRouter for ChatBot Arena (left) and MT-Bench (right)**. The dashed line represents the GraphRouter baseline. As k increases, the performance of our method improves, but it saturates once k reaches 10.

1132 choose $k=10$ as a balanced setting for capturing
 1133 user preferences without compromising generaliza-
 1134 tion.

1135 D Additional Results on New User 1136 Experiments

1137 Here, we present the results of the experiment de-
 1138 scribed in Section 5.2 on the ChatBot Arena and
 1139 MT-Bench datasets, as shown in Figure 7.

1140 E Further Experiment

1141 E.1 Extended Experiments on ChatBot Arena 1142 Users

1143 To more thoroughly validate that GMTRouter ex-
 1144 hibits strong generalization to new users, we further
 1145 selected all ChatBot Arena users with more than 15
 1146 historical interactions (161 users in total). We keep
 1147 the original training and test sets used in Section
 1148 5 unchanged and use all newly added data exclu-

Table 9: GMTRouter performance on the ChatBot Arena dataset with newly added users.

New User Num	0	3	161
Accuracy	0.774	0.780	0.790
AUC-ROC	0.875	0.858	0.837

1149 sively as the test set. GMTRouter is then evaluated
 1150 on these newly added users. The experimental re-
 1151 sults are shown in Table 9.

1152 This expanded evaluation enables us to assess
 1153 GMTRouter on a substantially larger user set and
 1154 further demonstrates its effectiveness and general-
 1155 ization capability.

1156 E.2 Experiments with Noisy Data

1157 We conduct experiments to evaluate GMTRouter’s
 1158 personalization performance under varying levels
 1159 of noisy data. Specifically, we use the MT-Bench

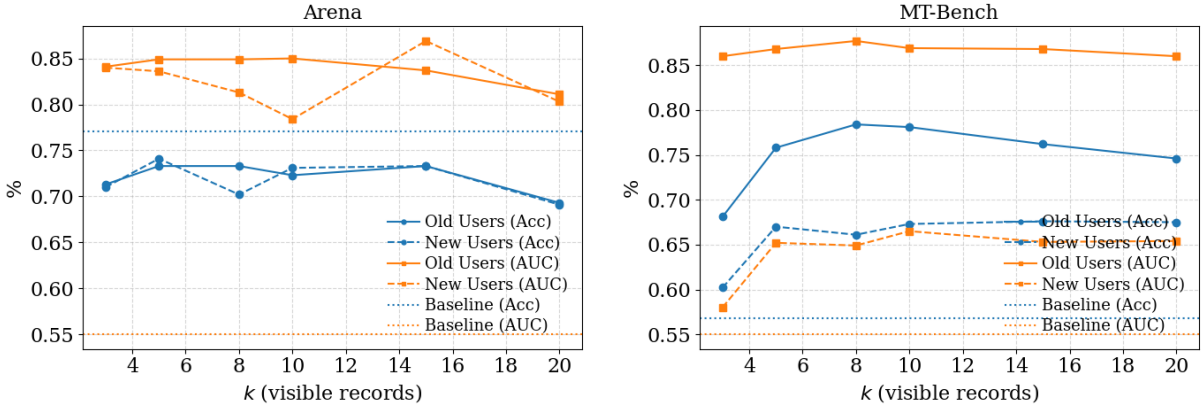


Figure 7: **Result comparison between old-user and new-user settings for ChatBot Arena (left) and MT-Bench (right).** The dashed line represents the GraphRouter baseline. The personalized performance under the new-user setting is comparable to that under the old-user setting, highlighting the strong generalization capability of our method.

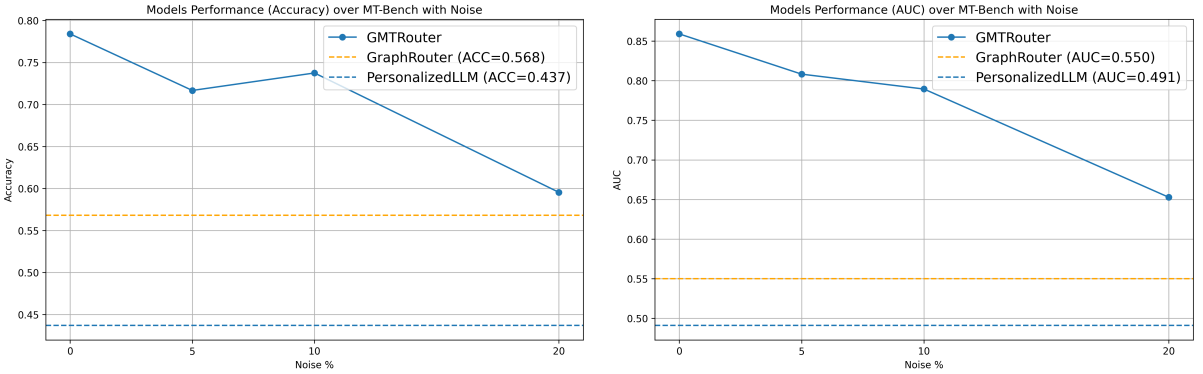


Figure 8: Performance against noisy preference over MT-Bench

dataset and swap the user ratings of $k\%$ of the data ($k = 5, 10, 20$), simulating noise and inconsistencies in preference signals (Li et al., 2024a). The experimental results are shown in Figure 8. Our results indicate that, although performance naturally decreases as noise increases, GMTRouter **degrades gracefully and remains competitive even with 20% noise, outperforming the strongest baseline trained on clean data.** We attribute this robustness to **user-conditioned graph sampling**, which aggregates signals across multiple interactions and mitigates the impact of individual noisy labels.

E.3 Comparison with Personalized Generation Method

To evaluate the personalization capabilities of GMTRouter against personalized generation approaches, we adopt the untuned In-Prompt Augmentation (IPA) based Retrieval Augmented Generation (RAG) from the LaMP benchmark (Salemi

et al., 2024b), employing Contriever (Izacard et al., 2021) as the retrieval backbone.

Experiments were conducted on the LaMP dataset. For the RAG baseline, we utilized the test set’s ‘abstract’ query to retrieve the **top-10** most relevant historical ‘abstracts’ from the specific user’s training data. These retrieved items were formatted as ‘title, abstract’ pairs and integrated into the LLM’s input as few-shot examples. Crucially, while GMTRouter relies on user ratings (calculated by comparing the LLM’s predicted titles against the ground-truth titles) as its supervision signal, we provided the RAG baseline with a **more direct and potent** form of supervision by explicitly incorporating the retrieved ground-truth titles into the few-shot examples.

We employ **average user rating** and **average token usage** as our primary evaluation metrics. The GMTRouter output consists of the raw response generated by the routed LLM, whereas the RAG output is generated by DeepSeek-R1 (DeepSeek-

Table 10: Performance Comparison of GMTRouter and IPA-RAG on the LaMP Benchmark. The **Random Routing** column serves as the **lower bound**, representing the expected user rating achieved by randomly selecting an LLM. Conversely, the **Theoretical Best Routing** column establishes the **upper bound** for the routing task, reflecting the user rating obtained by always selecting the LLM that yields the highest user rating for that specific instance.

Metric	Random	GMTRouter	IPA-RAG	GMTRouter + IPA-RAG	Theoretical Best
AVG User Rating	0.744	0.772	0.775	0.784	0.810
AVG Token Cost	–	293.56	3231.72	2358.89	–

Table 11: **GMTRouter performance using different convolutional layers.** The best and second-best results are highlighted in **bold** and underline, respectively.

Dataset Metric	Chatbot-Arena		MT-Bench		GSM8K		MMLU	
	Accuracy	AUC	Accuracy	AUC	Accuracy	AUC	Accuracy	AUC
GraphSAGE	0.768	<u>0.873</u>	0.569	0.645	0.635	0.648	0.494	0.487
HeteroConv	0.777	0.867	0.569	0.492	0.499	0.603	0.494	0.542
HANConv	0.766	0.776	<u>0.646</u>	<u>0.680</u>	0.774	<u>0.775</u>	<u>0.707</u>	<u>0.746</u>
HGTConv	<u>0.774</u>	0.875	0.784	0.859	<u>0.773</u>	0.859	0.771	0.870

AI, 2024), conditioned on the prompt augmented with the retrieved few-shot examples. Furthermore, we also investigated the personalization capabilities of combining both GMTRouter and RAG. Specifically, GMTRouter is employed to select the optimal LLM, and the RAG is subsequently used to construct the few-shot augmented prompt. The experimental results are presented in Table 10.

Our experiments lead to three key observations:

- **Comparable Personalization:** Both GMTRouter (routing) and the personalized generation approach (IPA-RAG) achieve comparable improvements in personalization capability as measured by the AVG User Rating.
- **Cost Disparity:** The IPA-RAG baseline incurs a significantly higher cost, requiring several times more tokens than GMTRouter, highlighting the efficiency gains offered by the routing mechanism.
- **Synergistic Effect:** Combining the two methods (GMTRouter + IPA-RAG) yields the best empirical performance. This suggests that routing and personalized generation techniques address distinct, complementary facets of personalization.

E.4 Using Different GNNs as the GMTRouter Backbone

We investigate how different GNNs used as the GMTRouter backbone affect its personalization capability. We evaluate four backbone models, including both homogeneous and heterogeneous GNNs: GraphSAGE (Hamilton et al., 2017), HGT (Ziniu Hu, 2020), HAN (Wang et al., 2019), and

HeteroConv, and present the results in Table 11.

We observe that attention-based heterogeneous GNNs (HAN and HGT) consistently outperform the homogeneous GraphSAGE and the simpler aggregation-based HeteroConv backbones. These results further indicate that GMTRouter is not dependent on any particular GNN architecture: multiple attention-based backbones achieve strong performance, suggesting that the gains primarily stem from the graph-based personalized routing framework and data modeling rather than from a specific convolutional operator.

F Baseline Routing Prompts

To benchmark routing strategies, we design two representative prompt templates: one for a vanilla router that selects the best LLM without personalization, and another for a personalized router that incorporates user history and preferences. Both prompts simulate realistic routing scenarios where a system must choose a single LLM for the next turn in a multi-turn dialogue.

G The Use of Large Language Models (LLMs)

During the writing of this paper, we used the GPT-5 Mini model for text polishing and grammatical corrections to enhance the readability of the manuscript.

Table 12: **Prompt Template: Vanilla LLM Routing (No Personalization)**

[Instruction]
 You are an expert routing agent. Your task is to select the most suitable Large Language Model (LLM) to handle the next query in a multi-turn conversation.

[Input Format]
 [Candidate LLM List]
 {{CANDIDATE_LLM_LIST}}
 [Previous Conversation]
 {{PREVIOUS_CONVERSATION}}
 [Current Query]
 {{CURRENT_QUERY}}

[Instructions for Model Selection]

- Consider the query difficulty, the context of the previous conversation, and each LLM’s expertise, cost, and size.
- Choose the single best LLM to respond to the current query.
- Output only the name of the selected LLM in the exact format below.
- Do not provide explanations or commentary.

[Output Format]
 <’{selected_model_name}’>

Table 13: **Prompt Template: Personalized Routing (User History Aware)**

[Instruction]
 You are an expert routing agent. Your task is to select the most suitable Large Language Model (LLM) to handle the next query in a multi-turn conversation, incorporating both model characteristics and personalization signals from the user’s history.

[Input Format]
 [Candidate LLM List]
 {{CANDIDATE_LLM_LIST}}
 [Previous Conversation]
 {{PREVIOUS_CONVERSATION}}
 [Current Query]
 {{CURRENT_QUERY}}
 [User Preference History]
 {{USER_PREFERENCE_HISTORY}}

[Instructions for Model Selection]

- Consider the query difficulty, the context of the ongoing conversation, the LLMs’ specializations, cost, and size.
- Additionally, factor in the user’s historical preferences and ratings to personalize the routing decision.
- Choose the single best LLM to respond to the current query.
- Output only the name of the selected LLM in the exact format below.
- Do not provide explanations or commentary.

[Output Format]
 <’{selected_model_name}’>
