

An Unified Search and Recommendation Foundation Model for Cold-Start Scenario

Yuqi Gong
Ant Group
Beijing, China
gongyuqi.gyq@antgroup.com

Xichen Ding
Ant Group
Beijing, China
xichen.dxc@antgroup.com

Yehui Su
Ant Group
Beijing, China
suyehui.syh@antgroup.com

Kaiming Shen
Ant Group
Beijing, China
kaiming.skm@antgroup.com

Zhongyi Liu
Ant Group
Hangzhou, China
zhongyi.lzy@antgroup.com

Guannan Zhang
Ant Group
Hangzhou, China
zgn138592@antgroup.com

ABSTRACT

In modern commercial search engines and recommendation systems, data from multiple domains is available to jointly train the multi-domain model. Traditional methods train multi-domain models in the multi-task setting, with shared parameters to learn the similarity of multiple tasks, and task-specific parameters to learn the divergence of features, labels, and sample distributions of individual tasks. With the development of large language models, LLM can extract global domain-invariant text features that serve both search and recommendation tasks. We propose a novel framework called S&R Multi-Domain Foundation, which uses LLM to extract domain invariant features, and Aspect Gating Fusion to merge the ID feature, domain invariant text features and task-specific heterogeneous sparse features to obtain the representations of query and item. Additionally, samples from multiple search and recommendation scenarios are trained jointly with Domain Adaptive Multi-Task module to obtain the multi-domain foundation model. We apply the S&R Multi-Domain foundation model to cold start scenarios in the pretrain-finetune manner, which achieves better performance than other SOTA transfer learning methods. The S&R Multi-Domain Foundation model has been successfully deployed in Alipay Mobile Application's online services, such as content query recommendation and service card recommendation, etc.

CCS CONCEPTS

- **Computing methodologies** → **Machine learning algorithms;**
- **Information systems** → **Data mining; Retrieval models and ranking.**

KEYWORDS

search and recommendation, LLM, multi-domain recommendation

ACM Reference Format:

Yuqi Gong, Xichen Ding, Yehui Su, Kaiming Shen, Zhongyi Liu, and Guannan Zhang. 2023. An Unified Search and Recommendation Foundation Model for Cold-Start Scenario. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management (CIKM '23)*, October 21–25, 2023, Birmingham, United Kingdom. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3583780.3614657>

1 INTRODUCTION

Modern commercial recommendation systems and search engines are widely used in online service platforms such as YouTube, TikTok, Taobao, Alipay, etc. Search and recommendation can facilitate users' behaviors to browse instant videos, buy products, use services, and make payments using E-wallets. The similarity between Search (S) and Recommendation (R) makes jointly modeling S&R a promising research topic. Some work [17, 18] propose to enhance recommendation by learning from a unified sequence of search and recommendation behaviors. Others [1, 14] propose to improve personalized search by adding users' multi-interests in recommendation. Methods to model search and recommendation tasks jointly are also proposed in [21–24, 26]. JSR framework in [23] simultaneously learns retrieval and recommendation models with shared item set and optimizes a joint loss function. Researchers in [22] applied two-level transformer encoders, text encoder to learn documents and queries, session encoder to model the integrated sequence of search and browsing behaviors. [26] applied GNN to learn node embedding of user&item, and treat search query as a special attribute of edges in the graph. In industrial scenarios, there are several benefits to model S and R jointly. First, there are multiple search and recommendation scenarios in a single mobile application. The training data collected in a single domain can't fully reflect users' complete intents and is sub-optimal compared to modeling them jointly. Secondly, majority of items are shared between search and recommendation. Once users have impressions of any products, videos, and services in a recommendation scenario, they should be able to retrieve the item in Search later for repurchase or reuse purposes. Despite the similarity between Search and Recommendation, there are also difficulties in modeling them jointly, such as the data imbalance issue of multiple domains, the heterogeneous issue of different item sets (videos, products, and services), the negative transferring issue. With the latest developments in large language models (LLM) [2, 4, 5, 20, 27], the pretrain-finetune framework [7, 13, 18, 22, 26] greatly improves the performance of

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 '23, October 21–25, 2023, Birmingham, United Kingdom

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-0124-5/23/10...\$15.00
<https://doi.org/10.1145/3583780.3614657>

downstream tasks. We are inspired by the strong expressive power of natural language features and propose to build the Search and Recommendation Foundation model on top of LLMs, which extract low-level domain-invariant text features of the query (Q) and item (I). The major difference between our S&R foundation model and traditional multi-domain multi-task models is how we use the domain-invariant text features to help constrain the divergence of different tasks, which alleviates data imbalance, negative transferring, and item heterogeneous issues. To summarize, our proposed S&R Foundation model has the following key contributions:

- We apply LLMs in S&R Multi-Domain Foundation model, and extract domain invariant text features to help mitigate the negative transferring and item heterogeneous issues in the multi-domain settings.
- We novelly proposed the Aspect Gating Fusion (Domain-Specific Gating) to fuse the ID feature, text features from LLMs, and sparse features. The Domain Adaptive Multi-Task module is also used to extract the domain-specific query and item towers' representations.
- For the cold start of new scenarios, we have conducted extensive experiments both offline and online, to show the effectiveness of supervised fine-tuning of our S&R Foundation model in downstream tasks, which is now fully deployed online and serving in Alipay's mobile application.

2 PROPOSED MODEL

2.1 Problem Formulation

Given a set of K search and recommendation tasks $\{D_k\}_{k=1}^K$, D_k denotes the dataset for the k -th task. We let $\mathcal{U} = \{u_1, u_2, \dots, u_N\}$ denote the user set, $\mathcal{I} = \{i_1, i_2, \dots, i_M\}$ denote the item set and $\mathcal{Q} = \{q_1, q_2, \dots, q_T\}$ denote the search query set. In real-world scenarios, items in search and recommendation usually come from different domains and are heterogeneous. Some items are shared across multiple domains and some items belong to each specific domain. And we let $\mathcal{I} = \mathcal{I}_1 \cup \mathcal{I}_2 \cup \dots \cup \mathcal{I}_K$ denote the union of all items in K domains, which contains M items in total. We aim to jointly train a search and recommendation (S&R) foundation model $M_{Foundation}^{S\&R}$ in the multi-task setting and predict the probability of user u_l click the item i_l given input query q_l as $p(y_l^{ctr} = 1 | u_l, q_l, i_l)$. And for search scenarios, additional query-item relevance score is also predicted as $p(y_l^{sim} = 1 | q_l, i_l)$. For cold start of a new search or recommendation scenario D^* , we restore parameters of embedding tables and partial network structures from the pretrained S&R foundation model $M_{Foundation}^{S\&R}$, and then apply supervised fine-tuning on the downstream tasks, such as click through rate (CTR) prediction, query-item relevance prediction, etc. For the search task D_k^S , we let $D_k^S = \{x_l = (u_l, q_l, i_l), y_l\}_l$, which denotes the search ranking task given the triple input of (user, query, item) as (u_l, q_l, i_l) . For the recommendation task D_k^R , we set search query set \mathcal{Q} as emptyset \emptyset in $D_k^R = \{x_l = (u_l, q_l = \emptyset, i_l), y_l\}_l$.

2.2 S&R Multi-Domain Foundation Model

As illustrated in Figure 1, the S&R Multi-Domain Foundation model has three main components: the User-Query-Item encoding module, the Aspect Gating Fusion module, and the Domain-Adaptive

Multi-Task module. Firstly, raw features of user, query and item pass through the embedding layers, and we extract the ID embedding, token-level text embedding and sparse features' embedding. We apply LLM to extract domain-invariant text features of query and item towers, which minimize the divergence of features' distribution cross multiple domains. Secondly, the Aspect Gating Fusion module is designed to merge different groups of ID, text, sparse features' embedding. The fusion network is to balance the relative importance of ID, text, and sparse features. Very few training samples contain ID features of cold start items and can't represent them well, and generic text features play more important role. Finally, we feed the concatenated embedding of user, query and item towers to the Domain Adaptive MTL module. The module has two outputs representing the click through rate (CTR) prediction task and the query-item relevance prediction task. The final loss function is the sum of CTR prediction loss \mathcal{L}^{ctr} , relevance prediction loss \mathcal{L}^{sim} and domain adaptive regularization \mathcal{L}^{reg} .

2.2.1 User Query and Item Encoding. We extract three towers for user, query and item respectively. For the user tower, $e_u^{ID} \in \mathbb{R}^D$ denotes user id embedding. $e_u^{NH} = [x_1, \dots, x_s, \dots, x_{NH}]$ denotes the unified sequence of both search and recommendation clicks in chronological order. Each behavior x_s is encoded as multiple layers of MLPs with inputs of ID feature, sparse feature of behavior type S or R, and other sparse features of attributes, $x_s = FC(e_s^{ID} \oplus e_s^{type} \oplus e_s^{attr})$. For the query (Q) and item (I) features, we extract both domain-invariant text features, such as tokens in search query and items' title, and the domain-specific sparse features. The tokens of Q and I go through the same tokenizer and we get the tokenized id sequences as integer tensors e_q^{Token} and e_i^{Token} . $e_q^{Token} = [e_q^1, e_q^2, \dots, e_q^{L_q}] \in \mathbb{R}^{L_q \times D}$ denotes the query's tokenized id tensor of length L_q , and $e_i^{Token} = [e_i^1, e_i^2, \dots, e_i^{L_i}] \in \mathbb{R}^{L_i \times D}$ denotes the item's tokenized id tensor of length L_i . For ID feature, we also embed the search query as ID feature $e_q^{ID} \in \mathbb{R}^D$, and item ID as $e_i^{ID} \in \mathbb{R}^D$. For the sparse features, we embed sparse features of Q as e_q^S and sparse features of I as e_i^S . Finally, we get the feature groups of query tower as $e_q = [e_q^{ID}, e_q^{Token}, e_q^S]$ and the feature groups of item tower as $e_i = [e_i^{ID}, e_i^{Token}, e_i^S]$.

LLM as Domain-Invariant Feature Extractor

We apply the pretrained Large Language Model, such as BERT [6], GPT [2], ChatGLM [8, 25], to extract domain-invariant text features on both query tower and item tower, represented as $\phi_{lm}(Q) = \phi_{lm}(e_q^{Token}) \in \mathbb{R}^{L_q \times D}$ and $\phi_{lm}(I) = \phi_{lm}(e_i^{Token}) \in \mathbb{R}^{L_i \times D}$. After mean pooling of the encoding layer, followed by shared linear projection, we get the domain-invariant text representation of query and item as $E_{lm}(Q) = W_{lm} \times \text{MEAN}(\phi_{lm}(e_q^{Token})) \in \mathbb{R}^H$, $E_{lm}(I) = W_{lm} \times \text{MEAN}(\phi_{lm}(e_i^{Token})) \in \mathbb{R}^H$. $W_{lm} \in \mathbb{R}^{H \times D}$ denotes the linear projection layer shared between query tower and item tower, H denotes the hidden size of the learned representations. The language models' representation is useful for cold start scenarios, especially when we have few training samples to update the ID feature of new item i^* and new search query q^* . We also apply linear projections $W_{ID}^q, W_{ID}^i \in \mathbb{R}^{H \times D}$ to ID feature of query and item tower, and get the ID representation of query and item as $E_{ID}(Q), E_{ID}(I) \in \mathbb{R}^H$. For the sparse features we have separate network (usually multiple layers of MLPs) to encode

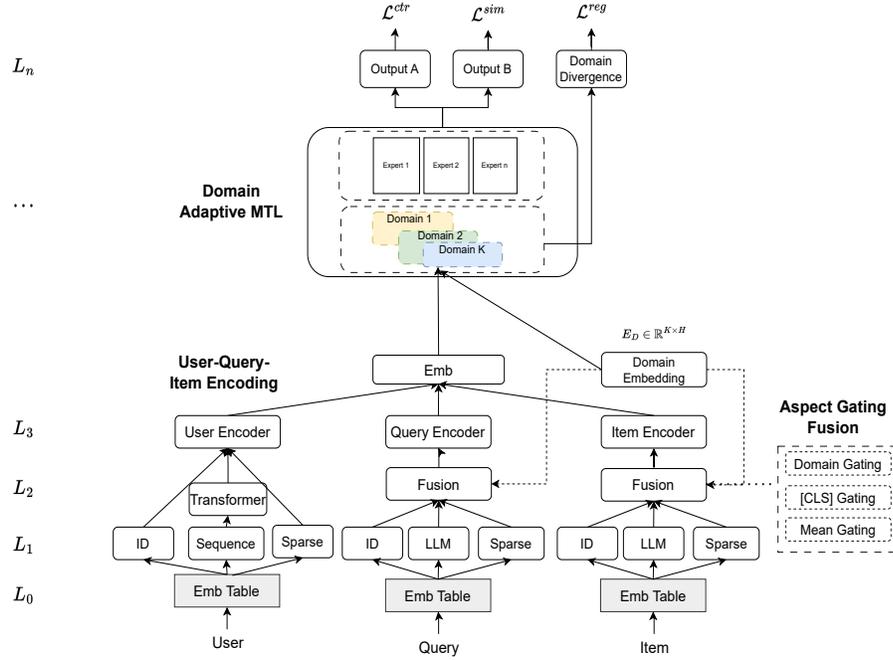


Figure 1: SR Multi-Domain Foundation Model Architecture

query and item as $E_S(Q), E_S(I) \in \mathbb{R}^H$. Finally, we get the feature groups of query tower as $[E_{ID}(Q), E_{Im}(Q), E_S(Q)]$ and item tower as $[E_{ID}(I), E_{Im}(I), E_S(I)]$.

2.2.2 Aspect Gating Fusion. After low level networks L_0 (embedding tables) and L_1 (feature encoding layers) in Figure 1, we fuse different aspects of query and item as in literature [12]. Each aspect E_a represents some fine-grained properties of query and item, such as ID, text and sparse features. \mathcal{A} denotes the set of aspects we extract from query and item. In S&R scenarios, we set $|\mathcal{A}| = 3$ as ID, text and sparse attributes. Final representations are fused as weighted sum of different aspects' representations.

$$E(Q) = \sum_a w_a(Q) E_a(Q), E(I) = \sum_a w_a(I) E_a(I) \quad \forall a \in |\mathcal{A}|$$

The weight vector $w(Q), w(I) \in \mathbb{R}^{|\mathcal{A}|}$ are outputs of a gating network, and we have different strategies to design the network.

- **Mean-Gating Strategy** Simply mean pooling of different aspects of query and item features as $w_a = \frac{1}{|\mathcal{A}|}$.
- **[CLS]-Gating Strategy** We use randomly initialized embedding $E_{CLS}(Q), E_{CLS}(I) \in \mathbb{R}^H$ to represent classification token [CLS] of query and item respectively.

$$w_a = \frac{e^{E_{CLS} E_a}}{\sum_{a \in |\mathcal{A}|} e^{E_{CLS} E_a}} \in \mathbb{R}^{|\mathcal{A}|}$$

- **Domain-Gating Strategy**

We design the domain gating strategy from the intuition that the fusion network has different weights when merging

different aspects of query and item. To model the differences across domains, we randomly initialize the domain embedding $E_D = [E_{D_1}, E_{D_2}, \dots, E_{D_K}] \in \mathbb{R}^{K \times H}$ as the representations of different domains. And the domain-specific gating is calculated as

$$w_a = \frac{e^{E_{D_k} E_a}}{\sum_{a \in |\mathcal{A}|} e^{E_{D_k} E_a}} \in \mathbb{R}^{|\mathcal{A}|}$$

2.2.3 Domain Adaptive Multi-Task Learning. The input to the Domain Adaptive Multi-Task module is the concatenation of representations of user, query and item towers as $\mathbf{x} = E(U) \oplus E(Q) \oplus E(I)$. For multi-domain setting, a series of multi-task and multi-domain models are proposed, such as SharedBottom[3], MMoE[15], PLE[19], STAR[16], SAMD[11], etc. These models use shared structures (Experts or MLP layers) to model the similarity among different tasks or domains, and use individual structures to learn the domain-specific properties. The difficulty of training the multi-domain models is the domain shift phenomena. For the k -th domain D_k , the marginal distribution of input feature $p(\mathbf{x}_k)$ and the conditional distribution of predicting output y_k as $p(y_k | \mathbf{x}_k)$ has divergence from other domains. The well studied MTL models handle the divergence of conditional distribution. We propose to add a Domain Adaptive Layer to the input features \mathbf{x}_i , which maps the inputs from multiple domains to a common vector space. We reuse the randomly initialized domain embedding $E_D = [E_{D_1}, E_{D_2}, \dots, E_{D_K}] \in \mathbb{R}^{K \times H}$ in section 2.2.2 and concatenate the domain embedding E_{D_k} to feature vector \mathbf{x}_i of instances from the k -th domain D_k , followed by domain-specific linear transformation W_k . Suppose \mathbf{x}_i and \mathbf{x}_j denote two instances from different domains in the same training batch, we can get the domain-adaptive representation $\hat{\mathbf{x}}_i, \hat{\mathbf{x}}_j$ as

$$\hat{x}_i = W_i(x_i \oplus E_{D_i}), \hat{x}_j = W_j(x_j \oplus E_{D_j})$$

We apply domain adaptation [9] techniques to constrain the divergence of distributions from domains i and j as $p(\hat{x}_i)$ and $p(\hat{x}_j)$ as $\mathcal{L}^{reg} = \sum_{i,j \in \{1,2,\dots,K\}} d(p(\hat{x}_i) || p(\hat{x}_j))$. In terms of divergence measurement, we compared different metrics such as Jensen-Shannon Divergence (symmetric KL Divergence), Maximum Mean Discrepancy (MMD) [10] in the experiment section. And we find the Jensen-Shannon Divergence achieves the best performance as

$$\mathcal{L}^{reg} = \sum_{i,j \in \{1,2,\dots,K\}} JS(p(\hat{x}_i) || p(\hat{x}_j))$$

Finally, on top of the Domain Adaptive Layer we stack the standard Multi-Task module, such as MMoE to extract outputs and predict two objectives, CTR prediction y^{ctr} and query-item relevance prediction y^{sim} .

CTR Prediction Click-Through Rate (CTR) Prediction is a common task in both search and recommendation scenarios. We apply a unified scoring function $y_1^{ctr} = f_{\theta}(u_l, q_l, i_l)$ in the S&R foundation framework to predict CTR with the triple inputs of user, item and query as (u, q, i) . For search tasks, users have explicit search query q . And for recommendation tasks users don't have explicit intentions. So we set $q = \emptyset$ as the default embedding in the unified scoring function.

Query-Item Relevance Prediction Query-Item Relevance Prediction is a common task in search scenarios, which predicts the relevance score of query-item pair of (q, i) and train a function $y_1^{sim} = f_{\phi}(q_l, i_l)$ to represent query-item pair's relevance score. The relevance prediction is usually a classification task.

2.2.4 Loss of S&R Foundation model. We train the S&R foundation model in multi-domain multi-task settings, using datasets from K domains. Each domain calculates either or both of two objectives of CTR prediction $y_1^{ctr} = f_{\theta}(u_l, q_l, i_l)$ and relevance prediction $y_1^{sim} = f_{\phi}(q_l, i_l)$, depending on whether the task is search or recommendation. The final objective function consists of three parts, the loss for CTR prediction \mathcal{L}^{ctr} , the loss for relevance prediction \mathcal{L}^{sim} , and the loss for domain adaptive regularizer \mathcal{L}^{reg} .

$$\mathcal{L} = \mathcal{L}^{ctr} + \mathcal{L}^{sim} + \mathcal{L}^{reg}$$

$$\mathcal{L}^{ctr} = \sum_{k \in K} \sum_{l \in N_k^{ctr}} \mathcal{L}_{ce}(f_{\theta}(u_l, q_l, i_l); y_l^{ctr})$$

$$\mathcal{L}^{sim} = \sum_{k \in K} \sum_{l \in N_k^{sim}} \mathcal{L}_{ce}(f_{\phi}(q_l, i_l); y_l^{sim})$$

2.3 Supervised Fine-Tuning Downstream Tasks

The pretrained S&R foundation model can benefit downstream tasks in the pretrain-finetune manner. The downstream model restores parameters from the foundation model, freezes part of the parameters and finetunes the remaining layers. We experiment different ways of freeze-finetune split. Firstly, the freeze-finetune split is between level L_0 and L_1 as in Figure 1. The pretrained embedding in level L_0 is frozen and the remaining layers from L_1 to L_n are finetuned. Secondly, the freeze-finetune split is between level L_1

and L_2 . The embedding in level L_0 as well as the parameters of encoding layers in level L_1 are frozen, and the parameters from level L_2 to L_n are finetuned. Given dataset of new downstream task $D^* = \{(u_l^*, q_l^*, i_l^*), y_l^*\}$, the domain embedding $E_{D^*} \in \mathbb{R}^H$ is randomly initialized and finetuned. In the experiment section, we thoroughly tested the performance of different ways of freeze-finetune split. We also compared the performance of pretrain-finetuning S&R Foundation model $M_{Foundation}^{S\&R}$ with the performance of training single domain model without transfer learning.

3 EXPERIMENT

To test the effectiveness of our proposed S&R Multi-Domain Foundation model, we want to answer the following questions:

- **RQ1:** Whether our joint S&R Multi-Domain Foundation model can achieve SOTA performance compared to other multi-domain and multi-task models?
- **RQ2:** In terms of query and item towers' representations, what's the performance of the domain-invariant text features extracted by LLM and Aspect Gating Fusion network compared to other methods?
- **RQ3:** Whether S&R Multi-Domain Foundation and Supervised Finetuning can help benefit cold start scenarios?

3.1 Experimental Settings

3.1.1 Dataset. We conducted extensive experiments of S&R Foundation model on real-world datasets, including 7 industrial datasets of Alipay Search Ranking and Query Recommendation. The statistics are summarized in table 1. S denotes the search dataset, in which users have explicit search query, such as Query-Item Relevance Prediction, Content Search Ranking, etc. And R denotes the recommendation dataset, in which users don't have explicit intent of search query. There are also some tasks between Search and Recommendation, which we classify as S/R, such as Query Suggest CTR Prediction, in which users have explicit query, and at the same time the task is a CTR prediction task to make recommendation of query suggestions to users.

3.1.2 Comparison Methods. S&R Foundation Model We compared our proposed S&R Multi-Domain Foundation model with SOTA multi-domain and multi-task models, such as Shared Bottom MTL [3], Multi-Gate Mixture of Experts (MMoE) [15], PLE [19], etc. For ablation study, we designed separate experiments to evaluate different modules of the framework, including the User-Query-Item encoding module, Aspect Gating Fusion module and Domain Adaptive Multi Task module. The experiment of S&R Multi-Domain Foundation (MLP) denotes the concatenated user-query-item representations are followed by multiple MLP layers. And the experiment of S&R Multi-Domain Foundation-MMoE-DA-JS denotes the representations are followed by a Domain Adaptive Layer (JS-Divergence) and MMoE multi-task module.

Domain-Invariant Text Features and Aspect Gating Fusion To prove the effectiveness of adding domain-invariant text features in S&R Foundation model, we have conducted experiments and ablation studies on different query and item token encoding methods on Alipay Content Query Recommendation dataset of tasks 4 in table 1. In the baseline method, we intentionally leave out the token-embedding of text features and only use ID and sparse features. We also compared randomly initialized token embedding with

Table 1: Statistics of Alipay Search Ranking and Query Recommendation Datasets.

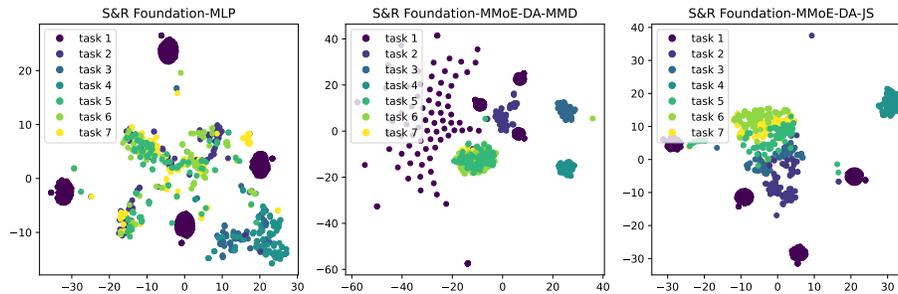
ID	Dataset	S/R	Train	Eval	Test	#Query	#Item
Task 1	Query-Item Relevance Prediction	S	76.2M	12.7M	12.7M	40K	40K
Task 2	Query Suggest CTR Prediction	S/R	145.4M	23.5M	23.5M	0.84M	0.16M
Task 3	Background Word Query Recommendation CTR Prediction	R	146.2M	24.3M	24.3M	-	65K
Task 4	Content Query Recommendation CTR Prediction	R	0.76M	0.09M	0.09M	-	4.6K
Task 5	People Also Ask DeepSuggest	S/R	2.4M	0.38M	0.38M	0.41M	25K
Task 6	Service Card Recommendation	S/R	1.01M	0.17M	0.17M	1.3K	1.6K
Task 7	Content Search Ranking	S	6.13M	1.03M	1.03M	0.27M	0.14M

Table 2: Performance of S&R Multi-Domain Foundation Model.

Method	Task1	Task2	Task3	Task4	Task5	Task6	Task7
S&R Multi-Domain Shared Bottom [3]	0.6483	0.8993	0.7829	0.6575	0.8511	0.8015	0.8561
S&R Multi-Domain MMoE [15]	0.6482	*0.9003	0.7812	0.6650	0.8463	0.7942	0.8599
S&R Multi-Domain PLE [19]	*0.7006	0.8981	0.7815	0.6682	0.8487	0.7978	0.8620
S&R Multi-Domain Foundation (MLP)	0.6827	0.8974	0.7784	0.6683	0.8462	0.7926	0.8629
S&R Multi-Domain Foundation-MMoE-DA-MMD	0.6874	0.8942	*0.7971	0.6793	0.8564	0.8203	0.8569
S&R Multi-Domain Foundation-MMoE-DA-JS	0.6942	0.8973	0.7912	*0.6979	*0.8703	*0.8312	*0.8692
Absolute Improvement	+0.0459	-0.0020	+0.0083	+0.0404	+0.0192	+0.0297	+0.0131

Table 3: Comparison of Query and Item Token Encoding Methods after Fine-tuning Task 4.

ID	Token Embedding	Query/Item Encoder	Finetune	AUC
1	Baseline: Without Token Emb	-	-	0.7524
2	Randomly Initialized	Mean Pooling	-	0.7551
3	Randomly Initialized	Transformer(L=1)	True	0.7544
4	Randomly Initialized	Transformer(L=6)	True	0.7559
5	SR Foundation (LM=Transformer)	Transformer(L=1)	L_0, L_1 :True	0.7562
6	SR Foundation (LM=Transformer)	Transformer(L=1)	L_0 :False, L_1 :True	0.7531
7	SR Foundation (LM=Transformer)	Transformer(L=1)	L_0, L_1 :False	0.7574
8	SR Foundation (LM=BERT)	BERT BASE(L=12)	True	0.7563
9	SR Foundation (LM=BERT)	BERT BASE(L=12)	False	*0.7580
10	SR Foundation (LM=ChatGLM 6B)	ChatGLM 6B Pretrained LLM [8, 25]	False	0.7518
11	SR Foundation (LM=ChatGLM 6B)	ChatGLM 6B Pretrained LLM [8, 25] + prompt	False	0.7503
12	SR Foundation (LM=ChatGLM2 6B)	ChatGLM2 Pretrained LLM [8, 25]	False	0.7502
	Absolute Improvement	-	-	+0.0056

**Figure 2: Visualization of SR Foundation Model's Domain-Adaptive Layers****Table 4: Comparison of Aspect Gating Fusion on Task 4.**

Method	AUC	Absolute Gain
Mean-Pooling	0.7385	-
[CLS]-Gating	0.7515	+0.0130
Domain-Gating	*0.7524	+0.0139

Table 5: Comparison of Cold Start Scenarios Task 4 and 6.

	Service Card Rec	Content Query Rec
Single Domain	0.8229	0.7295
SR Fdt->Finetune	0.8446	0.7574
Absolute Improvement	+0.0216	+0.0279

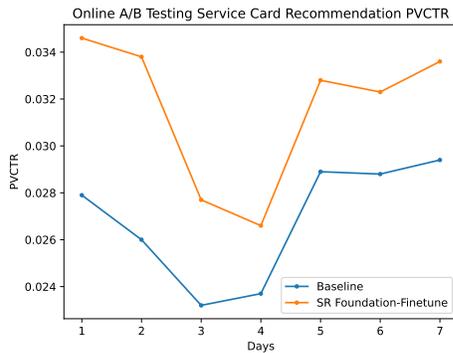


Figure 3: Online AB Testing PVCTR Performance of Service Card Recommendation.

embedding restored from the pretrained S&R foundation model under different configurations. For the encoders, we compared mean pooling, randomly initialized Transformer, BERT, ChatGLM-6B[8, 25], and ChatGLM+prompt, etc. In methods 10-12, we adopt ChatGLM-6B and ChatGLM2-6B to encode the text features of query and items. The implementation details are: we utilized the encoders of ChatGLM-6B and ChatGLM2-6B to convert the input text features and corresponding prompts into 4096-dimensional vectors, which are followed by 2 MLP dense layers and further reduced to 32-dimensional vectors. The prompt function used in our approach is defined as $f_{prompt}(X) = \text{"Extract keywords from sentence [X]"}$. To compare different Aspect Gating Fusion methods, e.g. Mean-Pooling, [CLS]-Gating, Domain-Gating, we conducted ablation studies and the results are listed in table 4.

Supervised Fine-Tuning on Cold Start Scenarios For cold start scenarios, we compared the performance of supervised finetuning (SFT) the foundation model using downstream dataset, with the method of training single domain model on several tasks, including Service Card Recommendation (task 6) and Content Query Recommendation (task 4) as in table 5.

3.2 Experimental Results

3.2.1 S&R Multi-Domain Foundation model. To compare the performance of different multi-domain models, we report AUC performance on 7 search and recommendation datasets in table 2. All the experimented models share same input features, User-Query-Item Encoding module, and Domain-Gating as Aspect Gating Fusion strategy. The baseline for the multi-task learning (MTL) module is the shared bottom model. MMoE-DA-MMD and MMoE-DA-JS represent models that utilize Maximum Mean Discrepancy (MMD) and Jensen-Shannon Divergence (JS-Divergence) to constrain the distributions of domain adaptive layers respectively. The asterisk (*) denotes the best performance achieved in each task, and the absolute improvement represents the absolute improvement of MMoE-DA-JS method compared to baseline. MMoE-DA-JS achieved best performance on 4 tasks: 4, 5, 6, 7 with AUC improvement of +0.0404, +0.0192, +0.0297, +0.0131 respectively. The domain-adaptive layer constrains the embedding representations from different domains in the common vector space. The t-SNE visualization of S&R Foundation model’s domain-adaptive layers is depicted in Figure 2. The

embedding depicted in the first subplot "S&R Foundation MLP" is scattered, and the embedding in the third subplot "S&R Foundation-MMoE-DA-JS" is coherently aligned.

3.2.2 Domain-Invariant Text Features and Aspect Gating Fusion. We report the performance of different methods to encode domain-invariant text features and freeze-finetune split in table 3 on task 4 Content Query Recommendation. Our proposed method of restoring pretrained parameters from BERT BASE (12 layers Transformer) in S&R Foundation, freezing the parameters of the encoder and finetuning the remaining networks achieves the best AUC performance 0.7580, which is 0.0056 absolute gain over baseline model. Comparing different freeze-finetune split (methods 5-7), we can see that freezing pretrained parameters in level L_0 and L_1 (method 7) achieves better performance than other split methods (method 5/6), which is 0.0043 absolute gain in AUC. As for the ablation studies of Aspect Gating Fusion in table 4, the baseline is to simply mean pooling three aspects: ID, text and sparse features. We can see the Domain-Gating achieves best AUC performance 0.7524, which is 0.0139 absolute gain over mean-pooling method.

3.2.3 Supervised Finetuning in Cold Start Scenarios. To prove the effectiveness of finetuning our pretrained S&R Foundation model, we compared cold start performance of two scenarios, Service Card Recommendation (task 6) and Content Query Recommendation (task 4). They are new scenarios and we only collected a few samples in a short period of time. The samples are splitted as we leave out last one day’s collected data for testing, and use the remaining data for fine-tuning the S&R Foundation. We also train the single domain model as the baseline. From table 5, we can see the fine-tuned S&R Foundation model achieves +0.0216 AUC improvement over single domain model on task 6 and +0.0279 AUC improvement on task 4.

3.2.4 Online AB Testing. To further prove the effectiveness of online performance in cold start scenario, we deployed the fine-tuned S&R Foundation model online in Service Card Recommendation scenario, and compared with baseline, which is the single domain DNN model. The results of the AB Testing from day 1 to day 7 are depicted in Figure 3. The key performance measurement of the cold start scenario is PVCTR (Page View Click Through Rate). And we observed that the fine-tuned S&R Foundation model achieved +17.54% relative gain in PVCTR over baseline. The online AB Testing results showed that our method achieved better performance than baseline consistently in cold start scenario.

4 CONCLUSION

In this paper, we study the problem of training search and recommendation tasks jointly as the S&R Multi-Domain Foundation model, and use domain adaptation techniques to benefit cold start scenario. Our proposed model learns user, query and item representations, applies LLM to encode domain invariant text features and Aspect Gating Fusion to merge ID, text and sparse features. We also conducted extensive experiments on finetuning the foundation models in cold start scenarios, which achieves better performance than the single domain model. The fine-tuned S&R Multi-Domain Foundation model has been successfully deployed online in Alipay’s multiple search and recommendation scenarios.

REFERENCES

- [1] Qingyao Ai, Yongfeng Zhang, Keping Bi, Xu Chen, and W. Bruce Croft. 2017. Learning a Hierarchical Embedding Model for Personalized Product Search. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Shinjuku, Tokyo, Japan) (SIGIR '17). Association for Computing Machinery, New York, NY, USA, 645–654. <https://doi.org/10.1145/3077136.3080813>
- [2] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.).
- [3] Rich Caruana. 1997. Multitask Learning. *Mach. Learn.* 28, 1 (1997), 41–75.
- [4] Zeyu Cui, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022. M6-Rec: Generative Pretrained Language Models are Open-Ended Recommender Systems. *CoRR abs/2205.08084* (2022).
- [5] Zeyu Cui, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022. M6-Rec: Generative Pretrained Language Models are Open-Ended Recommender Systems. *arXiv:2205.08084* [cs.LG]
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2–7, 2019, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, 4171–4186.
- [7] Xichen Ding, Jie Tang, Tracy Liu, Cheng Xu, Yaping Zhang, Feng Shi, Qixia Jiang, and Dan Shen. 2019. Infer Implicit Contexts in Real-Time Online-to-Offline Recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery Data Mining* (Anchorage, AK, USA) (KDD '19). Association for Computing Machinery, New York, NY, USA, 2336–2346. <https://doi.org/10.1145/3292500.3330716>
- [8] Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. GLM: General Language Model Pretraining with Autoregressive Blank Infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 320–335.
- [9] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. 2016. Domain-Adversarial Training of Neural Networks. *J. Mach. Learn. Res.* 17, 1 (jan 2016), 2096–2030.
- [10] Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola. 2012. A Kernel Two-Sample Test. *J. Mach. Learn. Res.* 13, null (mar 2012), 723–773.
- [11] Zhaoxin Huan, Ang Li, Xiaolu Zhang, Xu Min, Jieyu Yang, Yong He, and Jun Zhou. 2023. SAMD: An Industrial Framework for Heterogeneous Multi-Scenario Recommendation. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (Long Beach, CA, USA) (KDD '23). Association for Computing Machinery, New York, NY, USA, 4175–4184. <https://doi.org/10.1145/3580305.3599955>
- [12] Weize Kong, Swaraj Khadanga, Cheng Li, Shaleen Kumar Gupta, Mingyang Zhang, Wensong Xu, and Michael Bendersky. 2022. Multi-Aspect Dense Retrieval. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (Washington DC, USA) (KDD '22). Association for Computing Machinery, New York, NY, USA, 3178–3186. <https://doi.org/10.1145/3534678.3539137>
- [13] Ningning Li, Qunwei Li, Xichen Ding, Shaohu Chen, and Wenliang Zhong. 2022. Prototypical Contrastive Learning and Adaptive Interest Selection for Candidate Generation in Recommendations. In *Proceedings of the 31st ACM International Conference on Information and Knowledge Management* (Atlanta, GA, USA) (CIKM '22). Association for Computing Machinery, New York, NY, USA, 4183–4187. <https://doi.org/10.1145/3511808.3557674>
- [14] Jiongnan Liu, Zhicheng Dou, Qiannan Zhu, and Ji-Rong Wen. 2022. A Category-Aware Multi-Interest Model for Personalized Product Search. In *Proceedings of the ACM Web Conference 2022* (Virtual Event, Lyon, France) (WWW '22). Association for Computing Machinery, New York, NY, USA, 360–368. <https://doi.org/10.1145/3485447.3511964>
- [15] Jiaqi Ma, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H. Chi. 2018. Modeling Task Relationships in Multi-task Learning with Multi-gate Mixture-of-Experts. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19–23, 2018*. ACM, 1930–1939.
- [16] Xiang-Rong Sheng, Liqin Zhao, Guorui Zhou, Xinyao Ding, Binding Dai, Qiang Luo, Siran Yang, Jingshan Lv, Chi Zhang, Hongbo Deng, and Xiaoqiang Zhu. 2021. One Model to Serve All: Star Topology Adaptive Recommender for Multi-Domain CTR Prediction. In *CIKM '21: The 30th ACM International Conference on Information and Knowledge Management, Virtual Event, Queensland, Australia, November 1 - 5, 2021*, Gianluca Demartini, Guido Zuccon, J. Shane Culpepper, Zi Huang, and Hanghang Tong (Eds.). ACM, 4104–4113.
- [17] Zihua Si, Xueran Han, Xiao Zhang, Jun Xu, Yue Yin, Yang Song, and Ji-Rong Wen. 2022. A Model-Agnostic Causal Learning Framework for Recommendation Using Search Data. In *Proceedings of the ACM Web Conference 2022* (Virtual Event, Lyon, France) (WWW '22). Association for Computing Machinery, New York, NY, USA, 224–233. <https://doi.org/10.1145/3485447.3511951>
- [18] Zihua Si, Zhongxiang Sun, Xiao Zhang, Jun Xu, Xiaoxue Zang, Yang Song, Kun Gai, and Ji-Rong Wen. 2023. When Search Meets Recommendation: Learning Disentangled Search Representation for Recommendation. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Taipei, Taiwan) (SIGIR '23). Association for Computing Machinery, New York, NY, USA, 1313–1323. <https://doi.org/10.1145/3539618.3591786>
- [19] Hongyan Tang, Junling Liu, Ming Zhao, and Xudong Gong. 2020. Progressive Layered Extraction (PLE): A Novel Multi-Task Learning (MTL) Model for Personalized Recommendations. In *RecSys 2020: Fourteenth ACM Conference on Recommender Systems, Virtual Event, Brazil, September 22–26, 2020*. ACM, 269–278.
- [20] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. *CoRR abs/2302.13971* (2023).
- [21] Jun Xu, Xiangnan He, and Hang Li. 2018. Deep Learning for Matching in Search and Recommendation. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval* (Ann Arbor, MI, USA) (SIGIR '18). Association for Computing Machinery, New York, NY, USA, 1365–1368. <https://doi.org/10.1145/3209978.3210181>
- [22] Jing Yao, Zhicheng Dou, Ruobing Xie, Yanxiong Lu, Zhipeng Wang, and Ji-Rong Wen. 2021. USER: A Unified Information Search and Recommendation Model Based on Integrated Behavior Sequence. In *Proceedings of the 30th ACM International Conference on Information and Knowledge Management* (Virtual Event, Queensland, Australia) (CIKM '21). Association for Computing Machinery, New York, NY, USA, 2373–2382. <https://doi.org/10.1145/3459637.3482489>
- [23] Hamed Zamani and W. Bruce Croft. 2018. Joint Modeling and Optimization of Search and Recommendation. In *Proceedings of the First Biennial Conference on Design of Experimental Search & Information Retrieval Systems, Bertinoro, Italy, August 28–31, 2018* (CEUR Workshop Proceedings, Vol. 2167), Omar Alonso and Gianmaria Silvello (Eds.). CEUR-WS.org, 36–41.
- [24] Hamed Zamani and W. Bruce Croft. 2020. Learning a Joint Search and Recommendation Model from User-Item Interactions. In *WSDM '20: The Thirteenth ACM International Conference on Web Search and Data Mining, Houston, TX, USA, February 3–7, 2020*, James Caverlee, Xia (Ben) Hu, Mounia Lalmas, and Wei Wang (Eds.). ACM, 717–725.
- [25] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b: An open bilingual pre-trained model. *arXiv preprint arXiv:2210.02414* (2022).
- [26] Kai Zhao, Yukun Zheng, Tao Zhuang, Xiang Li, and Xiaoyi Zeng. 2022. Joint Learning of E-Commerce Search and Recommendation with a Unified Graph Neural Network. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining* (Virtual Event, AZ, USA) (WSDM '22). Association for Computing Machinery, New York, NY, USA, 1461–1469. <https://doi.org/10.1145/3488560.3498414>
- [27] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A Survey of Large Language Models. *CoRR abs/2303.18223* (2023).