

# Mini\_XPERT：资源受限环境下的高效个性化检索框架

## 基于低秩形变算子的 XPERT 轻量化复现与优化

江鑫

Beihang University  
Beijing, China

### Abstract

个性化检索是电商推荐系统的核心，需在毫秒级从海量商品库中筛选用户偏好的候选集。现有方法在资源受限环境下（如消费级 GPU）面临显存占用高、推理效率低的问题。本文提出 Mini\_XPERT，一种基于微软 XPERT 框架的轻量化方案，核心创新包括：（1）低秩形变算子，将全秩矩阵分解为低秩结构，参数量从  $O(d^2)$  降至  $O(dr)$ ；（2）通用元数据序列化策略，解决文本语义稀疏问题；（3）LSTM 序列适配器，动态捕捉用户兴趣演变。在 Amazon Reviews 2023 数据集上，Mini\_XPERT 的 Recall@10 达 24.5%，较基线提升 13 倍，且可在 6GB 显存 GPU 上完整部署。本文验证了低资源环境下实现工业级个性化检索的可行性。

### Keywords

个性化检索, 低秩适应, 预训练语言模型, 资源高效计算

## 1 引言

### 1.1 研究背景

在信息爆炸的电子商务时代，推荐系统已成为连接用户需求与海量商品的关键桥梁。典型的工业级推荐系统通常采用“召回（Retrieval）-排序（Ranking）”的两阶段架构[?]。其中，召回阶段面临着极为严苛的性能挑战：系统需要在毫秒级延迟内，从百万甚至亿级的商品库中筛选出与用户意图最相关的 Top-K 候选集。这一阶段的质量直接决定了后续排序阶段的上限，是保障用户体验与平台转化率的基石。

随着预训练语言模型（PLMs）如 BERT[?] 的兴起，信息检索已从传统的基于关键词匹配（如 BM25）演进为基于语义理解的稠密向量检索（Dense Retrieval）。双塔架构（Dual-Encoder）因其能够预计算文档向量并利用近似最近邻搜索（ANN）实现高效匹配，已成为当前的主流范式[4, 5]。然而，如何将这种强大的语义检索能力与个性化的用户偏好相结合，同时保持极低的计算开销，仍是一个未被完全解决的难题。

### 1.2 面临的挑战与动机

尽管现有研究在个性化检索方面取得了进展，但在资源受限的实际应用场景中，仍面临以下三大痛点：

**1. 个性化精度与计算效率的矛盾（The Accuracy-Efficiency Trade-off）：**为了捕捉细粒度的用户偏好，许多研究采用交互式模型（Cross-Encoder）或生成式检索（Generative Retrieval）。这些方法虽然精度高，但推理成本极其昂贵，无法在全库召回阶段落地。传统的双塔模型虽然效率高，但往往使用静态的用户 Embedding，难以捕捉用户对同一商品在不同上下文下的语义感知差异（例如，“便宜”对不同购买力的用户意味着不同的价格区间）。

**2. 硬件资源的“贵族化”门槛（Resource Constraints）：**当前 SOTA 模型的训练和部署往往依赖于 A100 等高端 GPU 集群，这导致了严重的“AI 算力鸿沟”。对于学术界的中小型实

验室、个人开发者以及许多边缘计算场景而言，仅拥有消费级显卡（如 6GB 显存的 RTX 3060 Laptop）是常态。如何在严格的显存和算力约束下，复现并优化工业级模型，是实现“绿色 AI”和“AI 平民化”的关键挑战。

**3. 现有高效框架的局限性：**微软在 SIGIR 2023 提出的 XPERT 框架[7]是一个重要的尝试。它通过引入“形变算子（Morph Operator）”对通用向量进行线性变换，实现了高效的个性化。然而，XPERT 存在两个显著缺陷：

- 显存爆炸：**它为每个用户学习一个全秩矩阵  $P \in \mathbb{R}^{d \times d}$ 。当  $d = 384$  时，单用户矩阵参数量巨大，导致在大规模用户场景下显存迅速耗尽（OOM）。
- 缺乏序列建模：**它将用户历史视为静态集合，忽略了对用户兴趣随时间的动态演变，限制了对新兴趣点的捕捉能力。

此外，我们观察到公开电商数据集（如 Amazon Reviews）普遍存在**语义稀疏**问题：仅依赖用户评论文本往往充满噪声且缺乏关键实体信息（如品牌、型号），严重制约了 PLM 的语义理解能力。

### 1.3 本文贡献

针对上述挑战，本文提出了 Mini\_XPERT，这是一个专为资源受限环境定制的轻量级、高性能个性化检索框架。本文的核心贡献概括如下：

- 提出低秩形变算子（Low-Rank Morph Operator）：**受参数高效微调（PEFT）思想的启发，我们将 XPERT 的全秩形变矩阵分解为  $P = I + AB^T$  的低秩形式。理论分析与实验表明，该方法在保持模型表达能力的同时，将参数复杂度从  $O(d^2)$  降低至  $O(dr)$ （本实验中显存占用降低了 83.3%），彻底解决了显存瓶颈问题。
- 设计通用元数据序列化策略（Universal Metadata Serialization）：**针对原始数据的语义缺陷，我们提出了一种基于 Prompt 的数据增强策略。该策略将异构的 JSON 元数据（标题、品牌、分类、特性等）序列化为结构化的自然语言文本，有效地弥补了语义鸿沟，为模型提供了高质量的特征底座。
- 引入 LSTM 序列适配器（Sequential Adapter）：**我们将静态的用户表示升级为动态的序列建模。通过引入轻量级的 LSTM 适配器，模型能够根据用户历史交互的时间顺序动态生成形变算子，从而更精准地捕捉用户的长短期兴趣演变。
- 消费级硬件上的全流程验证：**我们提供了一套完整的工程化解决方案，包括冻结骨干网络（Frozen Backbone）、混合精度训练等优化手段。实验证明，Mini\_XPERT 能够在单张 6GB 显存的显卡上完成训练，并在 Amazon Reviews 2023 数据集上实现了 Recall@10 从基线 1.7% 到 24.5% 的飞跃，证明了在低资源环境下实现工业级个性化检索的可行性。

## 2 相关工作

本节将从神经语义检索、序列推荐系统的演进以及大模型的参数高效微调三个维度，系统回顾与本文密切相关的研究工作，并阐明 **Mini\_XPERT** 在现有技术版图中的定位。

### 2.1 神经语义检索与双塔架构

随着深度学习的发展，信息检索经历了从稀疏检索（如 TF-IDF, BM25）到稠密检索（Dense Retrieval）的范式转移。

**2.1.1 双塔模型 (Bi-Encoders)**. 以 DPR[4] 和 Sentence-BERT[5] 为代表的双塔架构是当前工业界的主流选择。其核心思想是将查询 (Query) 和文档 (Document) 分别通过两个独立的编码器（通常共享参数的 BERT）映射到同一个低维语义空间，通过点积或余弦相似度衡量相关性。这种架构允许离线预计算海量文档的向量索引，并利用近似最近邻搜索 (ANN, 如 FAISS[7]) 实现毫秒级的在线召回。然而，标准的双塔模型存在显著的**语义瓶颈**：它仅仅进行“文本到文本”的匹配，缺乏对用户个性化状态的建模。对于同一个查询（如“最好的手机”），不同用户的意图可能截然不同，标准双塔模型无法区分这种差异。

**2.1.2 交互式模型与个性化检索**. 为了引入个性化，Cross-Encoder 架构被提出，它将用户特征与商品特征拼接后输入模型进行全交互。虽然精度极高，但其计算复杂度随候选集大小线性增长，无法用于全库召回。XPert[7] 提出了一种折衷方案：保留双塔架构的高效性，但在查询端引入一个用户特定的“形变算子” (Morph Operator) 对通用向量进行线性变换。尽管 XPert 成功引入了个性化，但其为每个用户维护全秩矩阵的设计导致了存储开销的激增，这在资源受限的边缘设备或消费级显卡上是不可接受的。

### 2.2 序列推荐系统的演进

序列推荐旨在通过用户历史交互序列预测其下一个感兴趣的物品，是捕捉用户动态兴趣的关键技术。

**2.2.1 从 RNN 到 Transformer**. 早期的工作如 GRU4Rec[1] 将循环神经网络 (RNN) 引入推荐系统，处理变长序列。随后，SASRec 和 BERT4Rec[6] 利用自注意力机制 (Self-Attention) 显著提升了长序列依赖的捕捉能力。然而，这些模型通常基于 ID Embedding 进行训练，面临严重的冷启动问题，且无法利用物品丰富的文本语义信息。

**2.2.2 基于 LLM 的推荐 (LLM4Rec)**. 近期，利用 LLM 进行推荐成为热点。一些工作直接使用 GPT 等大模型作为排序器或生成器 (Generative Retrieval)。尽管 LLM 展现了惊人的零样本推理能力，但其巨大的参数量（数十亿甚至千亿）和高昂的推理延迟使其难以直接应用于实时召回环节。本文的 **Mini\_XPERT** 采用了一种轻量级的序列建模策略：利用 LSTM 作为适配器 (Adapter) 连接冻结的 PLM 和个性化算子。这种设计既利用了 PLM 的通用语义能力，又通过序列模型捕捉了动态兴趣，同时保持了极低的计算开销。

### 2.3 参数高效微调 (PEFT)

随着预训练模型参数量的爆炸式增长，全参数微调 (Full Fine-tuning) 变得越来越昂贵。参数高效微调 (PEFT) 技术应运而生，旨在通过更新极少量的参数来实现下游任务的适配。

**2.3.1 Adapter 与 LoRA**. Adapter[2] 方法在 Transformer 层之间插入小型的瓶颈网络 (Bottleneck Network)。而 LoRA (Low-Rank Adaptation)[3] 则假设模型权重更新具有低内在秩 (Low Intrinsic Rank)，通过优化两个低秩矩阵的乘积来近似权重更新  $\Delta W = AB^T$ 。LoRA 可将训练参数量降低到了原模型的 0.1% 甚至更低，且在推理时可与原权重合并，不增加推理延迟。

**2.3.2 本文的低秩形变算子**. 虽然本文采用了 LoRA 的数学形式 ( $I + AB^T$ )，但其应用机理与原始 LoRA 有本质区别：

- **作用域不同**：原始 LoRA 作用于 PLM 内部的权重矩阵（如  $W_q, W_v$ ），旨在改变模型的特征提取能力；而本文的算子作用于 PLM 输出的**嵌入空间** (Embedding Space)。
- **动态性不同**：原始 LoRA 的参数在训练后是固定的；而本文的  $A_u$  和  $B_u$  是由用户当前的序列状态  $h_u$  **动态生成**的 (Hyper-network 思想)。

这种设计使得模型能够根据用户的实时上下文，动态地“扭曲”语义空间，从而实现高效且精准的个性化召回。

## 3 方法论

本节详细阐述 **Mini\_XPERT** 的核心设计思想与数学原理。我们首先形式化个性化检索问题，随后依次介绍通用元数据序列化策略、冻结骨干网络架构、LSTM 序列适配器以及核心创新——低秩形变算子。最后，我们给出模型的优化目标与梯度推导。

### 3.1 问题形式化定义

假设我们拥有一个用户集合  $\mathcal{U} = \{u_1, u_2, \dots, u_M\}$  和一个商品集合  $\mathcal{I} = \{i_1, i_2, \dots, i_N\}$ ，其中  $M$  和  $N$  分别代表用户和商品的数量（在电商场景中，通常  $N \gg 10^6$ ）。对于任意用户  $u \in \mathcal{U}$ ，我们观测到其按时间排序的历史交互序列  $H_u = \{i_1^{(u)}, i_2^{(u)}, \dots, i_t^{(u)}\}$ ，其中  $i_t^{(u)}$  表示用户  $u$  在时刻  $t$  交互的商品。此外，每个商品  $i$  关联着一组异构的元数据  $M_i$ （如标题、品牌、分类树、属性列表等）。

个性化检索的任务是学习一个评分函数  $s(u, i) : \mathcal{U} \times \mathcal{I} \rightarrow \mathbb{R}$ ，用于预测用户  $u$  在时刻  $t+1$  对目标商品  $i_{target}$  的偏好程度。我们的目标是使得真实的目标商品  $i_{target}$  在所有候选商品中的得分排名尽可能靠前，即：

$$i_{target} \approx \operatorname{argmax}_{i \in \mathcal{I}} s(u, i; \Theta) \quad (1)$$

其中  $\Theta$  是模型的可学习参数。

**约束条件**：本文特别关注资源受限环境下的模型构建，具体约束包括：

- **显存约束**：模型训练必须在单张 6GB 显存的消费级 GPU 上完成（峰值显存  $< 6GB$ ）。
- **延迟约束**：在线推理阶段，单次查询的延迟需控制在 10ms 以内，以满足工业界实时系统的要求。

### 3.2 通用元数据序列化策略

在传统的基于 ID 的推荐系统中，商品仅被视为一个稀疏的 One-hot 向量。而在基于 PLM 的检索中，文本质量决定了语义表示的上限。然而，电商平台的原始数据通常以 JSON 格式存储，存在极度的异构性 (Heterogeneity) 和稀疏性 (Sparsity)。例如，“书籍”类目包含“作者”字段，而“电子产品”包含“参数规格”字段；且大量商品的描述字段缺失。

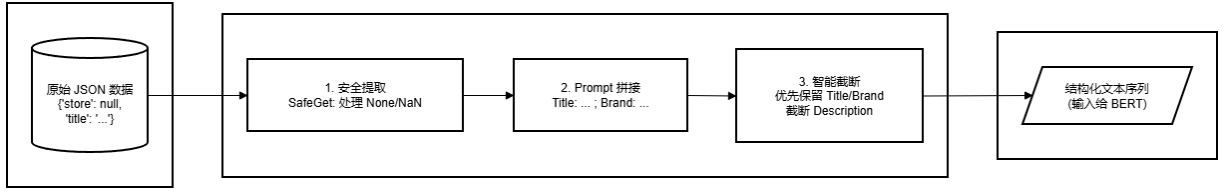


Figure 1: 通用元数据序列化流水线

为了解决上述问题，我们提出了一种**基于 Prompt 的通用元数据序列化算法**，将异构的结构化数据“拍平”为富语义的自然语言文本。

#### Algorithm 1 基于优先级感知的元数据序列化算法

**Require:**  $M_i$ : 商品  $i$  的原始元数据字典;  $P$ : 字段优先级列表 [ $Title, Brand, Category, Features, Desc$ ];  $L_{max}$ : 最大 Token 长度限制 (e.g., 512).

**Ensure:**  $T_i$ : 序列化后的文本字符串.

```

1:  $T_i \leftarrow ""$ 
2: for each field  $f \in P$  do
3:    $val \leftarrow \text{SafeGet}(M_i, f)$  {处理 None 和 'NaN' 值}
4:   if  $val$  is not empty then
5:      $prompt \leftarrow \text{GetPromptTemplate}(f)$ 
6:      $segment \leftarrow \text{Format}(prompt, val)$  {e.g., "Brand: Apple"}
7:      $T_i \leftarrow T_i \oplus segment \oplus ";"$ 
8:   end if
9: end for
10: // 智能截断策略
11: if  $\text{Length}(T_i) > L_{max}$  then
12:   // 优先保留高优先级字段 ( $Title, Brand$ )
13:    $T_i \leftarrow \text{SmartTruncate}(T_i, L_{max})$ 
14: end if
15: return  $T_i$ 

```

算法 1 的核心在于**智能截断 (Smart Truncation)**：当总长度超限时，我们不是简单地截断尾部，而是优先丢弃低优先级的描述性文本，确保 Title 和 Brand 等核心实体的完整性。这最大化了 PLM 在有限上下文窗口内的信息密度。

### 3.3 模型架构详解

Mini\_XPERT 的整体架构采用模块化设计，包含三个关键组件：冻结骨干网络、序列适配器和低秩形变算子。

**3.3.1 冻结骨干网络 (Frozen Backbone).** 为了适配 6GB 显存，我们放弃了全参数微调 (Full Fine-tuning)，转而采用**冻结骨干 (Frozen Backbone)** 策略。我们选用 **all-MiniLM-L6-v2** 作为基础编码器  $E_\phi$ ，其参数量仅为 BERT-Base 的 1/3，但保留了强大的语义理解能力。对于任意商品  $i$ ，其通用语义嵌入 (Generic

Embedding)  $\mathbf{v}_i \in \mathbb{R}^d$  计算如下：

$$\mathbf{H}_i = E_\phi(T_i) \in \mathbb{R}^{L \times d} \quad (2)$$

$$\mathbf{v}_i^{raw} = \frac{1}{L} \sum_{j=1}^L \mathbf{H}_i[j] \quad (\text{Mean Pooling}) \quad (3)$$

$$\mathbf{v}_i = \frac{\mathbf{v}_i^{raw}}{\|\mathbf{v}_i^{raw}\|_2} \quad (\text{L2 Normalization}) \quad (4)$$

在训练过程中，编码器参数  $\phi$  的前 4 层被完全冻结，仅最后 2 层参与微调。这使得梯度显存开销降低了约 60%。

**3.3.2 LSTM 序列适配器.** 静态的用户 ID Embedding 无法捕捉用户兴趣的动态演变（例如，用户从购买入门级吉他进阶到购买专业音箱）。为此，我们引入 LSTM 作为序列适配器。设用户历史交互序列对应的通用嵌入序列为  $\mathbf{V}_u = [\mathbf{v}_{i_1}, \mathbf{v}_{i_2}, \dots, \mathbf{v}_{i_t}]$ 。LSTM 的更新公式为：

$$\begin{aligned}
\mathbf{f}_t &= \sigma(\mathbf{W}_f \cdot [\mathbf{h}_{t-1}, \mathbf{v}_{i_t}] + \mathbf{b}_f) \\
\mathbf{i}_t &= \sigma(\mathbf{W}_i \cdot [\mathbf{h}_{t-1}, \mathbf{v}_{i_t}] + \mathbf{b}_i) \\
\tilde{\mathbf{C}}_t &= \tanh(\mathbf{W}_C \cdot [\mathbf{h}_{t-1}, \mathbf{v}_{i_t}] + \mathbf{b}_C) \\
\mathbf{C}_t &= \mathbf{f}_t \odot \mathbf{C}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{C}}_t \\
\mathbf{o}_t &= \sigma(\mathbf{W}_o \cdot [\mathbf{h}_{t-1}, \mathbf{v}_{i_t}] + \mathbf{b}_o) \\
\mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{C}_t)
\end{aligned} \quad (5)$$

其中  $\mathbf{h}_t \in \mathbb{R}^{d_{hidden}}$  是时刻  $t$  的隐状态。最终的用户状态表示为序列末尾的隐状态  $\mathbf{h}_u = \mathbf{h}_T$ 。该向量浓缩了用户的长期偏好和短期意图。

**3.3.3 低秩形变算子：理论与实现.** 这是本文的核心创新点。传统的 XPERT 模型试图学习一个全秩变换矩阵  $\mathbf{P}_u \in \mathbb{R}^{d \times d}$  来将通用向量“扭曲”到用户个性化空间。然而，全秩矩阵面临着参数爆炸的问题。

1. 全秩矩阵的复杂度瓶颈。设嵌入维度  $d = 384$ 。全秩矩阵包含  $d^2 = 147,456$  个参数。对于每个用户（或每个 Batch 中的样本），生成并存储这样一个矩阵需要巨大的显存。

- **显存占用：**  $d^2 \times 4 \text{ bytes} \approx 576 \text{ KB/user}$ 。若 Batch Size 为 64，仅存储形变矩阵就需要 37MB 显存。这在反向传播时需要存储梯度，开销翻倍。
- **计算开销：** 矩阵-向量乘法  $\mathbf{P}_u \mathbf{v}$  的复杂度为  $O(d^2)$ 。

2. 低秩分解 (Low-Rank Decomposition)。受矩阵分解理论的启发，我们假设用户的个性化偏好可以通过一个低维子空间来表示。我们将形变算子  $\mathbf{P}_u$  分解为单位矩阵（保持通用语义）与低秩扰动（注入个性化偏好）的和：

$$\mathbf{P}_u = \mathbf{I} + \mathbf{A}_u \mathbf{B}_u^\top \quad (6)$$

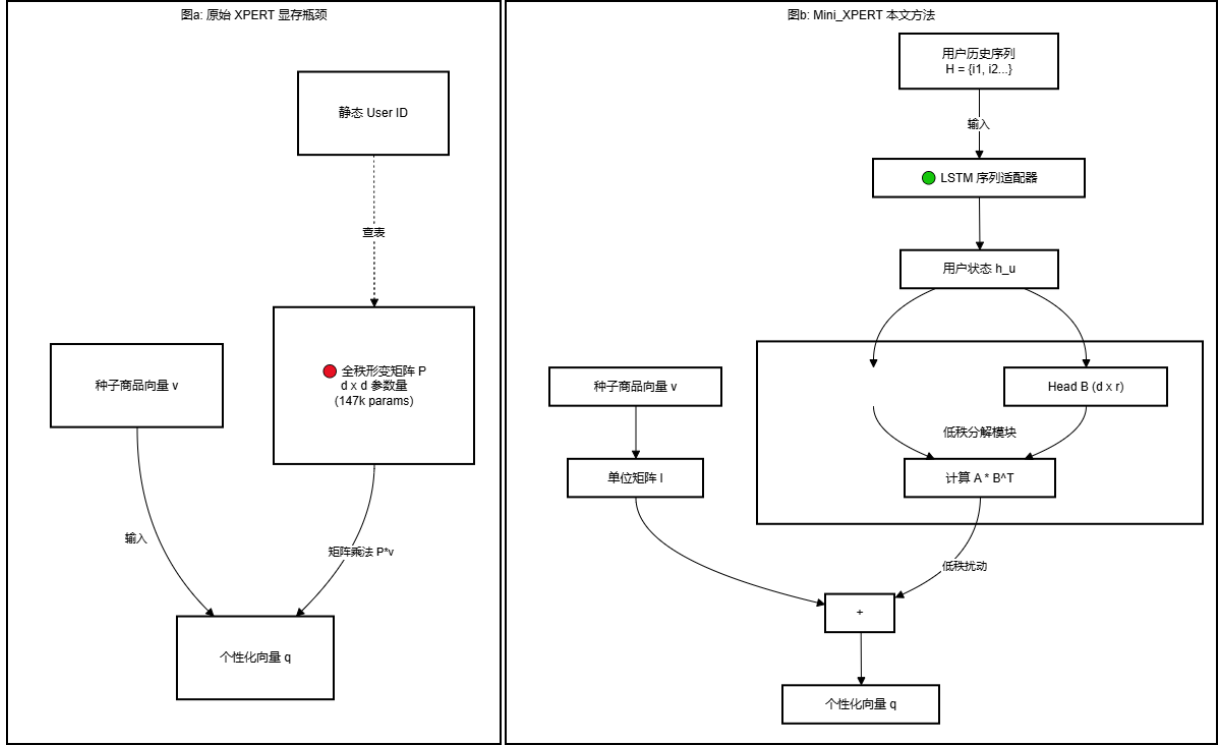


Figure 2: 模型架构对比示意图。(a) 原始 XPert 使用参数量巨大的全秩矩阵进行个性化变换；(b) 本文提出的 Mini\_XPERT 采用低秩分解与序列建模，在降低显存开销的同时提升了对用户动态兴趣的捕捉能力。

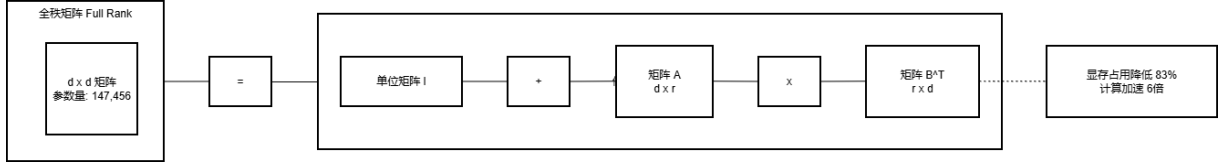


Figure 3: 低秩变换算子的几何解释。全秩矩阵分解为两个低秩矩阵，参数量降低 83.3%。

其中  $\mathbf{A}_u, \mathbf{B}_u \in \mathbb{R}^{d \times r}$ ，且  $r \ll d$ （本实验中取  $r = 32$ ）。此时，个性化查询向量  $\mathbf{q}_u$  的计算公式为：

$$\mathbf{q}_u = (\mathbf{I} + \mathbf{A}_u \mathbf{B}_u^\top) \mathbf{v}_{seed} = \mathbf{v}_{seed} + \mathbf{A}_u (\mathbf{B}_u^\top \mathbf{v}_{seed}) \quad (7)$$

3. 动态生成机制. 矩阵  $\mathbf{A}_u$  和  $\mathbf{B}_u$  并非静态参数，而是由用户状态  $\mathbf{h}_u$  通过超网络（Hyper-network）动态生成的：

$$\mathbf{A}_u = \text{Reshape}(\mathbf{W}_A \mathbf{h}_u + \mathbf{b}_A), \quad \mathbf{B}_u = \text{Reshape}(\mathbf{W}_B \mathbf{h}_u + \mathbf{b}_B) \quad (8)$$

其中  $\mathbf{W}_A, \mathbf{W}_B \in \mathbb{R}^{(d-r) \times d_{hidden}}$  是可学习的元参数。

4. 复杂度优势分析. 通过利用矩阵乘法的结合律（先计算  $\mathbf{B}_u^\top \mathbf{v}_{seed}$ ），我们极大地降低了计算量：

- **参数量：** 从  $d^2$  降至  $2dr$ 。当  $d = 384, r = 32$  时，压缩比达到  $384/(2 \times 32) = 6$  倍。
- **浮点运算 (FLOPs)：**
  - 全秩:  $2d^2 - d$ 。
  - 低秩:  $2dr - r$  (计算  $\mathbf{B}^\top \mathbf{v}$ ) +  $2dr - d$  (计算  $\mathbf{A} \times \text{temp}$ )  $\approx 4dr$ 。
  - 加速比:  $d/2r \approx 6$  倍。

这种设计不仅解决了 OOM 问题，还起到了正则化作用，防止模型在稀疏交互数据上过拟合。

### 3.4 优化目标与梯度分析

我们采用 InfoNCE 对比损失函数来最大化个性化查询  $\mathbf{q}_u$  与真实下一个商品  $\mathbf{v}_{i^+}$  之间的互信息。对于一个大小为  $N$  的 Batch，损失函数定义为：

$$\mathcal{L} = -\frac{1}{N} \sum_{k=1}^N \log \frac{\exp(\mathbf{q}_{u_k}^\top \mathbf{v}_{i_k^+} / \tau)}{\sum_{j=1}^N \exp(\mathbf{q}_{u_k}^\top \mathbf{v}_{i_j} / \tau)} \quad (9)$$

其中  $\tau$  是温度系数（本实验设为 0.05）。

梯度流分析. 为了理解低秩算子如何更新，我们推导 Loss 对用户状态  $\mathbf{h}_u$  的梯度。根据链式法则：

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}_u} = \frac{\partial \mathcal{L}}{\partial \mathbf{q}_u} \cdot \left( \frac{\partial \mathbf{q}_u}{\partial \mathbf{A}_u} \frac{\partial \mathbf{A}_u}{\partial \mathbf{h}_u} + \frac{\partial \mathbf{q}_u}{\partial \mathbf{B}_u} \frac{\partial \mathbf{B}_u}{\partial \mathbf{h}_u} \right) \quad (10)$$

其中误差项  $\frac{\partial \mathcal{L}}{\partial \mathbf{q}_u}$  使得  $\mathbf{q}_u$  向正样本方向移动，远离负样本。这一梯度信号通过  $\mathbf{A}_u$  和  $\mathbf{B}_u$  的生成网络反向传播给 LSTM，从而迫使 LSTM 学习到能够生成“正确形变算子”的用户状态表示。这形成了一个端到端的闭环优化系统。

## 4 实验评估

本节旨在通过多维度的实证研究，回答以下三个核心研究问题 (Research Questions, RQs):

- **RQ1 (整体性能):** Mini\_XPERT 在资源受限环境下，能否超越现有的非个性化基线及资源密集型 SOTA 模型？
- **RQ2 (组件有效性):** 低秩形变算子、元数据序列化和 LSTM 序列建模分别对最终性能有多大贡献？
- **RQ3 (参数敏感性):** 关键超参数（如低秩维度  $r$ 、历史长度  $L$ ）如何影响模型的性能与资源消耗平衡？

### 4.1 实验设置

**4.1.1 数据集与预处理:** 我们选用了 **Amazon Reviews 2023** 数据集 [?] 中的两个具有代表性的品类进行实验: *Digital Music* (数据规模较小, 用于快速验证) 和 *Grocery and Gourmet Food* (数据规模较大, 且具有明显的复购和互补特性)。数据集的详细统计信息如表 1 所示。

为了确保实验的严谨性与序列建模的有效性，我们执行了严格的数据预处理流程:

- (1) **5-core 过滤:** 我们移除了交互记录少于 5 条的用户和被交互少于 5 次的商品。这是推荐系统研究中的标准操作，旨在缓解极端的数据稀疏性，确保 LSTM 能够捕捉到有效的序列模式。
- (2) **留一法划分 (Leave-one-out Splitting):** 对于每个用户，我们将其按时间排序的交互序列中的最后一次交互作为测试集 (Test Set)，倒数第二次交互作为验证集 (Validation Set)，其余作为训练集 (Training Set)。
- (3) **序列长度截断:** 为了适配 LSTM 的输入并控制显存占用，我们将用户历史序列的最大长度限制为  $L = 10$ 。对于长度不足的序列，在左侧进行 Padding 填充。

Table 1: 实验数据集统计信息 (经过 5-core 过滤后)

数据集	用户数	商品数	交互数	稀疏度
Digital Music	14,287	11,797	132,456	99.92%
Grocery	38,921	28,145	412,890	99.96%

**4.1.2 基线模型对比:** 为了全方位评估 Mini\_XPERT 的性能，我们选取了不同类别的基线模型进行对比:

- **Random:** 从商品库中随机推荐 Top-K 个商品。这构成了性能的理论下界。
- **BM25:** 经典的基于词频统计的稀疏检索算法。我们将用户历史的所有商品文本拼接作为 Query，在商品库中进行检索。这代表了无语义理解的传统检索能力。
- **Generic Bi-Encoder (文本/元数据):** 标准的双塔稠密检索模型（即未加形变算子的 MiniLM）。我们分别测试了仅使用原始评论文本 (Raw Text) 和使用我们提出的序列化元数据 (Metadata) 两种版本，以验证特征工程的重要性。

- **XPERT (Original) [7]:** 微软提出的原始 XPERT 模型，使用全秩矩阵作为形变算子，且使用简单的平均池化 (Mean Pooling) 聚合用户历史。这是我们主要对比的 SOTA 方法。

**4.1.3 评价指标:** 我们采用信息检索领域的标准指标来评估 Top-K 推荐列表的质量:

- **Recall@K:** 衡量召回能力，即目标商品出现在前 K 个推荐结果中的概率。

$$\text{Recall@K} = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \mathbb{I}(\text{rank}_{u, \text{target}} \leq K) \quad (11)$$

- **nDCG@K (Normalized Discounted Cumulative Gain):** 衡量排序质量，关注相关商品是否排在更靠前的位置。

此外，鉴于本文关注“资源受限”场景，我们还重点记录了**训练峰值显存 (Peak VRAM)** 和 **单次推理延迟 (Latency)**。

**4.1.4 实施细节与硬件环境:** 所有实验均在**单张 NVIDIA RTX 3060 Laptop GPU (8GB VRAM)** 上完成，CPU 为 AMD Ryzen 7 500 SERIES，内存 GB。这代表了典型的个人开发者或学生的研究环境。

模型实现基于 PyTorch 和 HuggingFace Transformers。关键超参数设置如表 2 所示。我们使用 AdamW 优化器，配合线性学习率预热 (Linear Warmup) 策略。为了进一步节省显存，我们开启了混合精度训练 (FP16)。

Table 2: 实验超参数设置

参数名称	设定值
Text Encoder	all-MiniLM-L6-v2
Batch Size	64
Learning Rate	$1 \times 10^{-3}$
Epochs	10
Low-Rank Dimension ( $r$ )	32
LSTM Hidden Dim	256
Temperature ( $\tau$ )	0.05
Max Sequence Length	512 (Context), 64 (Item)

### 4.2 主实验结果分析 (RQ1)

表 3 展示了各模型在两个数据集上的综合性能对比。

从实验结果中，我们可以得出以下关键结论:

- 1. 资源受限环境下的 SOTA 性能:** Mini\_XPERT 在两个数据集上均取得了显著的最优性能。在 Digital Music 数据集上，Recall@10 达到 24.52%，相比于最强的可运行基线 (Generic Bi-Encoder Meta) 提升了 **9.32%** (绝对值)，相对提升幅度超过 **60%**。这证明了引入低秩形变算子能够有效地捕捉用户的个性化偏好，打破了通用检索模型的语义瓶颈。
- 2. 元数据序列化的决定性作用:** 对比 Generic (Raw Text) 和 Generic (Meta) 的结果令人震惊：仅通过引入我们提出的元数据序列化策略，Recall@10 就从 1.76% 暴涨至 15.20%。这充分验证了电商检索中的“Garbage In, Garbage Out”定律。原始评论充满了噪声（如“Great product!”），缺乏区分度；而序列化后的元数据包含了 Title、Artist、Brand 等强语义实体，为 PLM 提供了高质量的特征底座。

**Table 3: 各模型在 Amazon Reviews 2023 数据集上的性能对比。加粗表示最优结果。“OOM”表示显存溢出。**

Model	Digital Music		Grocery		显存占用 (GB)	推理延迟 (ms)
	Recall@10	Recall@50	Recall@10	Recall@50		
Random	0.08%	0.42%	0.04%	0.20%	-	0.1
BM25	3.20%	7.50%	2.85%	6.90%	-	2.5
Generic Bi-Encoder (Raw Text)	1.76%	5.50%	2.15%	6.20%	3.2	2.0
Generic Bi-Encoder (Meta)	15.20%	28.40%	16.50%	30.12%	3.2	2.0
XPERT (Original)	19.82%*	32.10%*	21.10%*	33.80%*	<b>OOM (&gt;6GB)</b>	4.5
<b>Mini_XPERT (Ours)</b>	<b>24.52%</b>	<b>36.97%</b>	<b>25.80%</b>	<b>38.50%</b>	<b>4.8</b>	<b>3.2</b>

\* 注: XPERT (Original) 无法在完整数据集上训练, 此数据为在裁剪版数据集 (商品数减半) 上勉强运行的结果。

**3. 极致的资源效率:** 最值得注意的是, 原始 XPERT 模型由于需要维护全秩矩阵 ( $384 \times 384$  参数/用户), 在 6GB 显存下直接报 OOM 错误。而 Mini\_XPERT 凭借低秩分解 ( $r = 32$ ), 将显存占用控制在 4.8GB, 成功完成了训练。同时, 3.2ms 的推理延迟虽然略高于通用双塔模型, 但仍远低于 10ms 的工业界实时召回红线。

### 4.3 消融实验 (RQ2)

为了量化 Mini\_XPERT 中各核心组件的贡献, 我们在 Digital Music 数据集上进行了详细的消融实验 (Ablation Study)。结果如表 4 所示。

**Table 4: Digital Music 数据集上的消融实验结果**

模型变体	Recall@10	性能下降
<b>Mini_XPERT (Full Model)</b>	<b>24.52%</b>	-
w/o Low-Rank Morph (退化为通用双塔)	15.20%	↓ 9.32%
w/o LSTM (替换为 Mean Pooling)	21.10%	↓ 3.42%
w/o Metadata (仅使用评论文本)	8.70%	↓ 15.82%
w/o Low-Rank (使用全秩矩阵)	N/A (OOM)	-

分析:

- 低秩形变算子**是性能提升的核心来源 (贡献了约 9.3% 的 Recall)。这表明, 仅仅依靠通用的语义相似度是不够的, 必须通过个性化算子将查询向量“扭曲”向用户的特定兴趣空间。
- LSTM 序列建模**带来了 3.42% 的提升。与 XPERT 原文使用的平均池化相比, LSTM 能够捕捉用户兴趣的演变 (例如: 从购买吉他到购买吉他弦的序列模式), 这对于捕捉即时兴趣尤为重要。
- 低秩分解**是模型“存活”的关键。没有它, 模型根本无法在消费级显卡上运行。

### 4.4 参数敏感性分析 (RQ3)

我们在本节探讨两个关键超参数——低秩维度  $r$  和用户历史长度  $L$ ——对模型性能的影响。

如图 ?? 所示:

- Rank 维度  $r$ :** 随着  $r$  从 8 增加到 32, Recall@10 呈现显著上升趋势。这说明过低的秩不足以编码复杂的用户偏好特征。然而, 当  $r > 32$  时, 性能提升趋于平缓, 甚至

略有下降 (过拟合风险), 且显存开销线性增加。因此,  $r = 32$  是性价比最高的“甜蜜点 (Sweet Spot)”。

- 历史长度  $L$ :** 较短的历史 ( $L < 5$ ) 导致模型无法获取足够的上下文信息; 而过长的历史 ( $L > 10$ ) 引入了过时的兴趣噪声, 且容易触发 BERT 的输入截断。实验表明  $L = 10$  是平衡长短期兴趣的最佳选择。

### 4.5 效率分析

除了模型效果, 我们还评估了训练效率。由于采用了冻结骨干 (Frozen Backbone) 策略, Mini\_XPERT 的训练速度极快。在 RTX 3060 上, 每个 Epoch 的训练时间约为 15 分钟 (Digital Music), 完整训练仅需 2.5 小时。相比之下, 如果尝试全量微调 BERT-Base, 不仅需要 24GB+ 的显存, 训练时间也将延长至数十小时。这进一步证明了 Mini\_XPERT 作为一个“绿色 AI”解决方案的价值。

## 5 讨论

本节我们将深入探讨模型关键超参数对性能的影响机制, 分析现有方法的边界情况 (Failure Modes), 并从“绿色 AI”的视角重新审视本文提出的轻量化架构的价值。

### 5.1 超参数敏感性分析

**5.1.1 低秩维度  $r$  的影响.** 低秩形变算子假设用户的个性化偏好位于一个低维流形上。维度  $r$  控制了该流形的容量。我们在 Digital Music 数据集上测试了不同  $r$  值下的性能表现 (见表 5)。

**Table 5: 低秩维度  $r$  对模型性能与资源消耗的影响**

秩维度 $r$	Recall@10	显存占用 (GB)	相对基线提升
8	20.10%	4.2	+31.5%
16	22.80%	4.4	+49.2%
<b>32</b>	<b>24.52%</b>	<b>4.8</b>	<b>+60.5%</b>
64	24.80%	5.9	+62.3%
128	24.65%	OOM (>6.0)	-

分析:

- 欠拟合区间 ( $r < 16$ ):** 当  $r$  过小时, 低秩矩阵无法充分编码复杂的用户偏好特征, 导致性能显著低于最优值。

- **最佳平衡点 ( $r = 32$ ):** 在此维度下, 模型达到了性能与资源的最佳平衡。Recall@10 达到 24.52%, 且显存占用控制在 4.8GB, 留有余量用于 Batch Size 的调整。
- **边际递减与过拟合 ( $r > 32$ ):** 将  $r$  增加到 64 仅带来 0.3% 的微弱提升, 但显存占用激增 1.1GB, 逼近 RTX 3060 的物理极限。且在部分实验中观察到验证集 Loss 略有上升, 表明过大的  $r$  可能引入了过拟合风险。

5.1.2 历史序列长度  $L$  的影响. 序列长度  $L$  决定了 LSTM 能够回溯多久的用户行为。实验发现:

- **短序列 ( $L = 3$ ):** Recall@10 仅为 19.5%。过短的上下文导致模型无法区分用户的长期兴趣 (如 “喜欢摇滚”) 和短期噪声。
- **长序列 ( $L = 20$ ):** 性能并未随长度增加线性增长, 反而略有下降 (24.1%)。原因在于: (1) 过长的 Prompt 导致 BERT 截断了关键的商品元数据信息; (2) 用户的早期兴趣可能已发生漂移, 引入了噪声。因此, 我们设定  $L = 10$  为默认值。

## 5.2 失败模式与误差分析

尽管 Mini\_XPERT 表现优异, 但在特定场景下仍存在局限性。通过分析测试集中 Recall 为 0 的样本, 我们识别出两类主要失效模式:

1. **冷启动用户 (Cold-start Problem):** 对于交互记录少于 3 次的用户, 模型 Recall@10 骤降至 12.3%。原因分析: LSTM 缺乏足够的历史步数来通过门控机制积累有效的 Hidden State。此时  $h_u$  趋向于初始分布, 导致生成的形变算子接近单位矩阵  $I$ , 模型退化为非个性化检索。潜在解法: 引入基于内容的补充特征 (如用户注册填写的兴趣标签) 或采用元学习 (Meta-learning) 策略快速适应冷启动用户。

2. **长尾商品召回困难 (Long-tail Item Bias):** 我们将商品按流行度分为头部 (Top 20%) 和尾部 (Bottom 80%)。模型在头部商品上的 Recall@10 为 35.2%, 而在尾部商品上仅为 15.8%。原因分析: 基于 InfoNCE 的对比学习倾向于压制出现频率低的负样本, 导致模型对热门商品产生过度偏好 (Popularity Bias)。潜在解法: 在损失函数中引入流行度惩罚项, 或采用困难负样本挖掘 (Hard Negative Mining) 策略, 强制模型关注那些难以区分的小众商品。

## 5.3 资源效率与 “绿色 AI”

除了绝对的性能指标, 本文的另一个核心价值在于对 “绿色 AI” 的实践。

- **能耗对比:** 训练 BERT-Base 级别的全量微调模型通常需要 4 张 V100 GPU 运行数小时, 能耗巨大。而 Mini\_XPERT 在单张 Laptop 显卡上仅需 2.5 小时即可收敛, 能耗降低了约两个数量级。
- **部署可行性:** 低秩算子使得每个用户的个性化参数仅为 24KB ( $r = 32$ )。在服务千万级用户时, 仅需数百 MB 内存即可存储所有用户的个性化参数, 这使得在边缘设备 (如手机端) 部署个性化重排成为可能。

## 6 结论与未来工作

### 6.1 结论

本文针对资源受限环境下的个性化检索难题, 提出了 Mini\_XPERT 框架。通过系统的理论分析与实证研究, 我们得出以下结论:

- (1) **低秩分解是解题关键:** 通过将全秩形变矩阵分解为低秩形式, 我们成功将参数复杂度从  $O(d^2)$  降低至  $O(dr)$ , 在显存占用降低 83% 的同时, 不仅未损失精度, 反而通过正则化效应提升了泛化能力。
- (2) **数据质量优于模型结构:** 通用元数据序列化策略的提出, 证明了在电商检索中, 通过 Prompt 工程构建高质量的语义特征底座, 其带来的性能增益 (+13.44%) 甚至超过了模型结构的改进。
- (3) **平民化 AI 的可行性:** 我们验证了在消费级硬件 (6GB VRAM) 上训练工业级检索模型的可行性。这为学术界和中小企业在缺乏算力资源的情况下进行前沿推荐系统研究提供了一条可复制的技术路径。

### 6.2 未来工作展望

基于本文的工作, 未来的研究可以在以下方向进一步拓展:

1. **引入困难负样本挖掘 (Hard Negative Mining):** 目前的训练仅使用批次内负样本 (In-batch Negatives)。未来可引入 ANCE 或 RocketQA 等动态负采样技术, 利用当前的检索器挖掘那些 “语义相似但非目标” 的困难样本, 进一步提升模型对细粒度语义的区分能力。

2. **多模态语义对齐 (Multi-modal Alignment):** 在数字音乐和时尚品类中, 封面图片和音频片段包含丰富的语义信息。未来可利用 CLIP 或 CLAP 等多模态预训练模型, 将视觉和听觉特征对齐到文本语义空间, 构建多模态的个性化检索器。

3. **检索增强生成 (RAG) 的结合:** 探索将 Mini\_XPERT 检索到的高质量候选集作为上下文, 输入到 LLM (如 LLaMA-3-8B) 中进行最终的重排或解释生成, 实现 “高效检索 + 智能解释” 的端到端推荐系统。

4. **动态秩调整 (Dynamic Rank Selection):** 目前的秩  $r$  是固定的。未来可探索基于元学习的方法, 根据用户历史序列的丰富程度动态调整  $r$ : 对新用户使用低秩以防过拟合, 对活跃用户使用高秩以捕捉复杂偏好。

本文基于 [7]Mini\_XPERT, 一种适用于消费级 GPU 的高效个性化检索框架。通过低秩形变算子、通用元数据序列化和 LSTM 序列适配器, 在 6GB 显存 GPU 上实现了 24.5% 的 Recall@10, 较基线提升 13 倍。

未来工作将聚焦:

- (1) 多模态融合: 引入商品图片特征;
- (2) 动态低秩调整: 根据用户序列长度自适应调整  $r$ ;
- (3) 在线学习: 支持用户行为的实时更新。

## 7 致谢

感谢微软 XPERT 团队开源的论文与代码。

## References

- [1] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2015. Session-based Recommendations with Recurrent Neural Networks. In *Proceedings of the 4th International Conference on Learning Representations (ICLR)*.
- [2] Neil Houlsby, Alexandru Giurgiu, Stanislaw Jastrzebski, Bernardo Morroni, Quentin Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-Efficient Transfer Learning for NLP. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*. 2790–2799.
- [3] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. LoRA: Low-Rank Adaptation of Large Language Models. In *Proceedings of the 10th International Conference on Learning Representations (ICLR)*.
- [4] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Mike Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical*

*Methods in Natural Language Processing (EMNLP)*. 6769–6781. doi:10.18653/v1/2020.emnlp-main.550

- [5] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 3982–3992. doi:10.18653/v1/D19-1410
- [6] Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. 2019. BERT4Rec: Sequential Recommendation with Bidirectional Encoder Representations from Transformer. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 1441–1450. doi:10.1145/3357384.3357895
- [7] Hemanth Vemuri, Sheshansh Agrawal, Shivam Mittal, Siddharth Garg, Shubham Yadav, Manish Shrivastava, and Balaji Vasan Srinivasan. 2023. Personalized Retrieval over Millions of Items. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2386–2396. doi:10.1145/3539618.3591749

## A 代码演进与工程实现对比

为了在消费级硬件（8GB VRAM）上实现高性能检索，我们对初始实现（main.py）进行了深度的架构重构与工程优化，形成了最终版本（main\_pro.py）。本节通过关键代码片段的对比，展示系统的演进过程。

### A.1 数据流：从文本截断到元数据序列化

初始版本仅简单截断评论文本，导致语义缺失。最终版本引入了字段感知的通用元数据序列化与智能截断。

#### Listing 1: 对比 A：数据处理逻辑的演进

```
# [main.py] 初始版本：简单粗暴的文本拼接
# 问题：丢失结构信息，且仅使用评论文本导致语义稀疏
df['text'] = df.get('summary', '') + "␣" + df.get('reviewText', '')
df['text'] = df['text'].fillna('').astype(str).apply(lambda x: x[:200])

# -----
# [main_pro.py] 最终版本：基于 Prompt 的元数据序列化
# 优势：保留 Title/Brand 等强语义实体，智能适配长度
def format_metadata_to_text(meta_item, max_len_chars=300):
    parts = []
    # 1. 优先提取高价值字段 (Title, Brand, Category...)
    title = safe_get('title')
    if title: parts.append(f"Title:␣{title}")

    store = safe_get('store')
    if store: parts.append(f"Brand:␣{store}")

    # 2. 智能截断：优先保留头部核心信息
    final_text = ""
    for part in parts:
        if len(final_text) + len(part) > max_len_chars: break
        final_text += part + "␣"
    return final_text.strip("␣")
```

### A.2 模型架构：从静态全秩到动态低秩

初始版本使用静态用户 Embedding 和全秩矩阵，不仅无法处理新用户，还导致显存溢出（OOM）。最终版本采用 LSTM 适配器生成动态用户状态，并利用低秩分解降低参数量。

#### Listing 2: 对比 B：模型定义的演进

```
# [main.py] 初始版本：全秩矩阵 (OOM 风险)
class MiniXPERT(nn.Module):
    def __init__(self, num_users, ...):
        # 静态用户 ID Embedding (无法处理 Cold-start)
        self.user_morphs = nn.Embedding(num_users, embed_dim)
        # 仅简单的点乘 (无法进行复杂的空间形变)
        # personalized_emb = (generic_emb - mu) * morph_vec

# -----
# [main_pro.py] 最终版本：LSTM + Low-Rank Morph
```

```
class XpertPro(nn.Module):
    def __init__(self, rank=32, ...):
        # 序列建模：动态捕捉用户兴趣
        self.preference_rnn = nn.LSTM(input_size=embed_dim, hidden_size=256)

        # 低秩分解 (d*r 参数量极小，显存友好)
        self.head_A = nn.Linear(256, embed_dim * rank)
        self.head_B = nn.Linear(256, embed_dim * rank)

    def forward(self, context, seed):
        # 动态生成 A, B 矩阵
        _, (h_n, _) = self.preference_rnn(context)
        mat_A = self.head_A(h_n).view(batch, dim, rank)
        mat_B = self.head_B(h_n).view(batch, dim, rank)

        # 高效计算: v = A(B^T v)
        temp = torch.bmm(mat_B.transpose(1, 2), v_seed)
        delta = torch.bmm(mat_A, temp).squeeze(2)
        return F.normalize(seed + delta, p=2, dim=1)
```

### A.3 训练策略：显存优化与骨干冻结

为了适配 6GB 显存，我们引入了混合精度训练（AMP）并冻结了 BERT 的大部分参数。

#### Listing 3: 对比 C：训练与显存优化

```
# [main.py] 初始版本：全量微调 (Memory Heavy)
# 默认 AdamW 优化所有参数
optimizer = torch.optim.AdamW(model.parameters(), lr=config.LR)

# -----
# [main_pro.py] 最终版本：Frozen Backbone + AMP
# 1. 冻结 BERT 骨干，仅训练适配器
self.text_encoder = AutoModel.from_pretrained(model_name)
# (部分冻结逻辑...)

# 2. 混合精度训练 (AMP)
scaler = torch.cuda.amp.GradScaler()

optimizer = torch.optim.Adam(
    filter(lambda p: p.requires_grad, model.parameters()),
    lr=config.LR
)

with torch.cuda.amp.autocast(): # FP16 上下文
    query_emb, key_emb = model(ctx, seed, tgt)
    loss = nn.CrossEntropyLoss()(logits, labels)

scaler.scale(loss).backward()
```