

Lightmem: 高效知识压缩与参数效率研究

LoRA 微调中信息压缩与秩大小的探索性分析

许永烨

2025 年 7 月 11 日

项目交付报告

1. 实验思路与目标
2. 实验设置与过程
3. 主要结果与分析
 - 压缩质量分析 (实验 A)
 - 参数效率分析 (实验 B 及正交实验)
4. 鼓励探索：从结论到创新
5. 交付成果与结论

1. 实验思路与目标

核心问题

在参数高效微调（如 LoRA）日益普及的背景下，我们如何通过优化输入数据的“信息形态”，来提升模型的学习效率和最终性能？

研究方向一：信息密度

- 文本被压缩后，模型还能学到多少知识？
- 是否存在一个“最佳压缩点”，能平衡信息保真度与学习效率？
- 不同的压缩方法（摘要 vs. 抽取）对模型学习有何不同影响？

研究方向二：参数效率

- 一个固定容量的 LoRA 适配器，其知识存储是否存在“上限”？
- LoRA 的秩（rank）大小，与它能有效学习的知识量和信息密度之间，是否存在一种最佳的匹配关系？

本实验旨在通过一系列受控实验，系统性地回答上述核心问题。

1. 知识压缩极限探究：

- 验证原始文本经不同压缩比处理后，LoRA 微调效果的**衰减边界**。

2. 参数存储容量分析：

- 在固定 LoRA 秩条件下，量化测试模型 (Qwen2-3B) 可有效存储的**结构化知识上限**。

2. 实验设置与过程

实验设置：数据集与模型

数据集

我们选取了两个在任务类型和数据形态上具有显著差异的数据集，以进行交叉验证：

- **非结构化编辑数据 (WikiUpdate)**: 代表离散、事实性的知识。
- **对话数据集 (LongMemEval)**: 代表连续、上下文相关的知识。

模型与微调方案

- **基座模型**: Qwen2.5-3B-Instruct
- **微调方法**: QLoRA (4-bit a-symmetric) 在实验中动态设置秩为 $r \in \{4, 8, 16\}$
- **微调平台**: 单卡 H800

实验过程：一个三阶段的系统性流程

1. 数据预处理与压缩:

- 使用 qwen-plus 模型 API，对源文本进行 5 个级别的摘要式压缩和 1 个级别的抽取式压缩。
- 使用 Qwen2.5-3B 的分词器，计算每个压缩版本的 Token 缩减率 (TRR)。

2. 模型微调 (核心执行阶段):

- **实验 A ($r=8$):** 训练 14 个模型，覆盖两个数据集的所有压缩级别。
- **实验 B ($r=8$):** 训练 10 个模型，覆盖 WikiUpdate 从 50 到 500 条数据的不同知识量。
- **正交实验 ($r=4, 16$):** 训练 12 个模型，完成“高效精简版”的正交实验。

3. 自动化评估:

- **WikiUpdate:** 使用 LLM-as-a-Judge (qwen-plus) 进行 1-5 分制自动评分。
- **LongMemEval:** 调用官方 evaluate_qa.py 脚本，自动计算 Accuracy。

核心指标：Token 缩减率 (TRR)

定义与目的

Token 缩减率 (Token Reduction Rate, TRR) 是一个用于定量衡量文本压缩程度的标准化指标。它回答了核心问题：“与原文相比，压缩后的文本在长度上缩减了多少？”

计算公式

其计算公式由以下表达式定义：

$$TRR = 1 - \frac{\text{Tokens}_{\text{compressed}}}{\text{Tokens}_{\text{original}}}$$

- $\text{Tokens}_{\text{original}}$ ：原始文本的 Token 总数。
- $\text{Tokens}_{\text{compressed}}$ ：压缩后文本的 Token 总数。
- 所有 Token 计数均使用基座模型 Qwen2.5-3B-Instruct 的分词器完成。

示例解读

当 `summ_l3_medium` 的 TRR 为 64.47% 时，意味着经过中度摘要压缩后，输入给模型的上下文 Token 数量，只有原始文本的 $1 - 0.6447 \approx 35.5\%$ 。

3. 主要结果与分析

结果与发现 1: WikiUpdate 性能衰减拐点

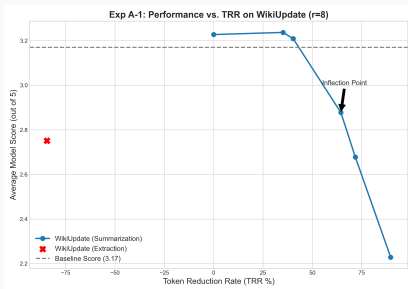


图 1: WikiUpdate 数据集上，模型性能随压缩比（TRR）的变化趋势（ $r=8$ ）。

- **有益的去噪:** 轻度压缩 (TRR 35%) 使性能达到峰值 (3.2367)，超越了原始文本。
- **衰减拐点:** 当 TRR 超过 64% 时，性能出现首次断崖式下跌，信息损失开始加剧。
- **结构破坏的代价:** 抽取式 (ext) 的 TRR 为负，其性能 (2.7512) 远逊于摘要式。不仅在于变相延长了文本，还在于文本结构的破坏。
 - LLM 无法高效地重构上下文
 - 使用 LoRA 微调效果不如全参微调（参数不易学习）
 - LLM as Judge 更易因上下文割裂被误判为“不一致”

结果与发现 2: LongMemEval 性能的 “Aha Moment”

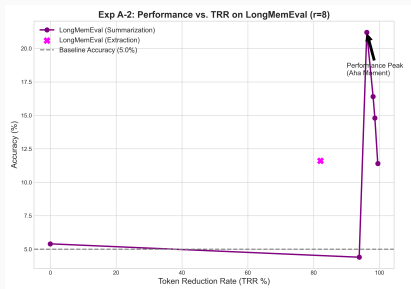


图 2: LongMemEval 数据集上, 模型性能随压缩比 (TRR) 的变化趋势 ($r=8$)。

- **原始文本的“毒性”**: 直接微调超长、嘈杂的原始对话历史, 效果极差 (5.4%)。
- **性能涌现**: 适度的摘要预处理 (summ_l2, TRR 96%) 带来了 “Aha Moment”, 性能从 5% 跃升至 21.2%。
- **规律验证**: 性能曲线呈现 “先升后降” 的模式, 与 WikiUpdate 的发现相互印证。

数据集之间的差异

- **事实数据集 WikiUpdate:** 基座模型在无微调的情况下仍然达到了一个较好的表现，甚至可以说和微调后的模型相差无几——预训练 LLM 经过反复的知识注入和微调，已经能在事实问题上达到了良好的表现。而 LoRA 只是强化这个零样本的问答能力。
 - 事实是使用基座模型和微调后的模型手动测试几个问题，得到的输出都是正确且类似的。
- **对话数据集 LongMemEval:** 对话数据集的问题相当于针对上下文相关知识进行问答。失去上下文的基座模型必然达到较低的基础准确率，而在经过适度摘要的文本微调后才出现了较大的提升。对于长程、嘈杂的对话历史，输入信息的“质量”（信噪比）远比“数量”更重要。

结果与发现 3：知识存储的“天花板效应”

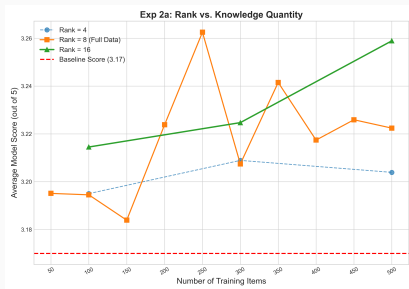


图 3：不同 LoRA 秩下，模型性能随训练知识数量的变化。

- **未现饱和曲线**: 模型性能在训练了少量数据后，便迅速达到了一个围绕基座模型性能小幅波动的“天花板”。
- **瓶颈非容量**: 这证明对于这种“开卷考试”式的任务，限制模型性能的并非 LoRA 的知识存储容量，而是其已经接近饱和的基础阅读理解能力。后续增加再多的同类数据，并不能提升它底层的“阅读理解”这个核心能力，因此分数便进入了围绕其能力“天花板”波动的平台期。

结果与发现 4：“有效压缩比-秩”的最佳匹配

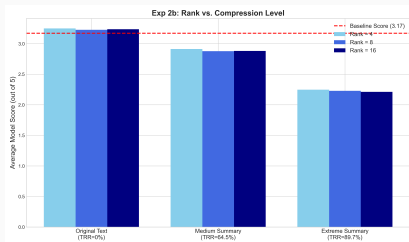


图 4：不同 LoRA 秩下，模型在不同压缩级别上下文上的性能。

经验性匹配建议

对于我们的 3B 基座模型和上下文内问答任务，一个更小的 LoRA 秩 ($r=4$) 在几乎所有压缩比下，都表现出了**最佳或接近最佳**的性能，是性价比最高的选择。这挑战了“秩越大越好”的普遍直觉。

4. 鼓励探索：从结论到创新

探索方向一：混合压缩策略 (Hybrid Compression)

数据洞见：摘要与抽取的优劣势

我们的实验明确表明，单纯的抽取式 (ext) 压缩因破坏自然语言结构而性能不佳 (WikiUpdate 得分 2.75)，远逊于摘要式 (summ_l3 得分 2.88)。这证明了保持文本的可读性对于模型学习至关重要。

方案设想：两阶段混合压缩 (Two-Stage Hybrid Compression)

“信息保真度”和“文本可读性”的更优策略应运而生：

1. 第一步 (核心事实抽取):

- 使用一个强大的 LLM，对原始文本进行一次预处理，只抽取其中最核心、最不可或缺的实体、关系和关键事实。

2. 第二步 (约束性摘要生成):

- 以第一步抽出的核心信息为“必须包含的约束条件”，再让 LLM 对原始文本生成一段通顺、连贯的摘要。
- **预期优势:** 该策略生成的上下文，既能确保关键信息 100% 不丢失，又保持了自然语言的流畅性，有望在高压缩比下取得远超任何单一方法的性能。

探索方向二：动态秩分配方案 (Dynamic Rank Allocation)

数据洞见：固定秩的低效性

我们的正交实验结果是提出此方案的最直接证据。数据显示，对于我们 3B 的基座模型和“开卷”任务，一个更小的 LoRA 秩 ($r=4$) 在几乎所有压缩比下，都表现出了最佳或接近最佳的性能。这证明了“一刀切”地使用固定秩（如 $r=8$ ）是一种资源效率低下的做法。

方案设想：基于信息密度的自适应微调系统

1. 第一步 (训练前信息密度评估):

- 在微调前，系统首先对准备好的上下文 (context) 进行一次快速的、轻量级的评估。我们可以直接使用已经定义的 TRR (Token 缩减率) 作为这个评估指标。

2. 第二步 (基于规则的动态秩分配):

- 根据我们从正交实验中得到的“匹配建议”，系统会自动为本次训练选择一个最合适的 rank 值。
- 预期优势:** 在保证性能的同时，极大地节约了计算资源，将 LoRA 调优从一个需要人工反复尝试的“炼丹”过程，变成了一个能自适应调节的智能化流程。

5. 交付成果与结论

1. 压缩质量分析报告

我们成功地量化了不同压缩方法对模型性能的影响，并通过图表清晰地定位了两个数据集上的性能衰减“拐点”与“峰值点”。

2. 参数效率实验报告

我们通过正交实验，验证了秩与知识存储量之间的非线性关系，并基于实验数据，提出了一个反直觉但极具价值的匹配建议：对于离散知识问答任务，一个更小的 LoRA 秩 ($r=4$) 往往是最高效的选择。

3. 探索性方案

基于坚实的实验数据，我们为“混合压缩策略”和“动态秩分配方案”这两个前沿方向，提出了具体、可行的设计思路。

感谢聆听 Q & A