# Gemini - 翻译与解读研究方案

gemini.google.com/share/4ddee3cd4c35



# TITOK论文完整翻译与逐段深度解读

# 摘要

论文题目:TITOK: 通过对比盈余迁移词元级知识以移植LoRA 作者: Chanjoo Jung¹, Jaehyung Kim¹ 单位: 延世大学 联系方式: chanjoo0427@yonsei.ac.kr, jaehyungk@yonsei.ac.kr 论文提交信息: arXiv:2510.04682v1 [cs.CL] 2025年10月6日,作为ICLR 2026会议论文审阅中

第一部分:论文标题、作者、摘要与引言

# A. 标题、作者与摘要

# 【原文翻译】

TITOK: 通过对比盈余迁移词元级知识以移植LoRA

### 摘要

大型语言模型(LLMs)被广泛应用于现实世界的场景中,但对其进行微调会带来巨大的计算和存储成本。参数高效微调(PEFT)方法,例如LoRA,减轻了这些成本,但其适配的参数依赖于基础模型,无法在不同的骨干网络之间迁移。解决此问题的一种方法是知识蒸馏,但其有效性本质上取决于训练数据。TransLoRA等近期工作通过生成合成数据来避免此问题,但这增加了复杂性,因为它需要训练额外的判别器模型。在本文中,我们提出了

TITOK,这是一个通过词元级知识迁移实现有效LORA移植的新框架。具体来说,TITOK通过源模型带有和不带有LoRA的对比盈余来捕获任务相关信息。这种盈余突出了信息丰富的词元,并实现了对合成数据的选择性过滤,所有这些都无需额外的模型或开销。通过在多个迁移设置下对三个基准进行的实验表明,所提出的方法持续有效,总体平均性能增益比基线提高了+4-8%。

### 【深度解读】

**解构标题:** 这篇论文的标题就像一部电影的预告片,精准地概括了核心内容。让我们逐一 拆解这些关键词:

- "移植LoRA (Transplant LoRA)": 这是最终目标。LoRA是一种高效微调技术,可以理解为给一个通用的大模型(比如一个知识渊博的大学毕业生)配备一套针对特定任务的"专业技能插件"(比如金融分析技能)。"移植"的意思是,我们想把这套"技能插件"从一个旧模型身上取下来,安装到一个更新、更强的模型身上,而不需要重新训练。这就像给AI做技能移植手术。
- "词元级知识 (Token-level Knowledge)": 这是手术的精度。"词元" (Token) 是模型处理语言的基本单位,可以理解为一个单词或一个汉字。这意味着,这次的"技能移植"手术精度极高,能够精确到每一个词的层面,而不是笼统地迁移整个知识块。
- "对比盈余 (Contrastive Excess)": 这是实现高精度移植的"秘方"或"手术刀"。"对比"指的是比较同一个模型在安装"技能插件"前后的表现差异。"盈余"则代表了安装插件后,模型在哪些词上表现得"更有信心"、"更专业"。这个"盈余"就是我们要提取和移植的核心知识。

摘要的承诺: 摘要清晰地讲述了一个"问题-现有方案缺陷-我们的解决方案"的故事。

- **问题**: 大模型很强大,但定制化(微调)成本太高。
- **现有方案1 (LoRA)**: LoRA降低了成本,但它的"技能插件"和特定模型"绑定"了,换个模型就作废,投资就浪费了。
- **现有方案2 (知识蒸馏):** 可以迁移知识,但需要原始的训练数据,这些数据往往是保密的或难以获得。
- **现有方案3 (TransLoRA)**: 用AI自己造数据 (合成数据)来解决数据问题,但又引入了新麻烦——需要额外训练一个"数据审查员" (判别器模型)来筛选数据质量,过程变得复杂。
- **TITOK的承诺**: 我们提出了一种更"干净"、更高效的方案。通过"对比盈余"这个巧妙的方法,我们能自动识别出最有价值的知识(信息丰富的词元),并用它来筛选合成数据,整个过程不需要额外的数据,也不需要额外的模型,简洁而强大。实验结果证明,这个方法非常有效。

#### 【原文翻译】

#### 1引言

大型语言模型(LLMs)(Brown et al., 2020; Vaswani et al., 2023)在许多现实世界应用中取得了显著进展,包括聊天机器人(OpenAl et al., 2024)、搜索引擎(Xiong et al., 2024)和编码助手(Rozière et al., 2024)。虽然微调LLMs已被证明是提高下游任务性能的有前景的方法,但它会产生巨大的计算和存储成本。参数高效微调(PEFT)(Houlsby et al., 2019)方法,例如LoRA(Hu et al., 2021),通过仅更新一小部分参数同时保持基础模型冻结来减轻了这一负担。然而,PEFT适配的参数依赖于基础模型,无法在不同的模型之间迁移。鉴于新LLMs的快速发布和可用模型的日益多样化,这一限制正变得越来越关键。

### 【深度解读】

这段引言为我们描绘了当前人工智能领域的宏大背景和面临的一个核心挑战。我们可以用 一个比喻来理解:

- 大型语言模型 (LLM) 是什么? 想象一个刚刚毕业的大学生,比如GPT-4或Llama-3。 他博览群书,知识面极广,上知天文下知地理,这是他的"基础模型"能力。
- 微调 (Fine-tuning) 是什么?现在,一家医院想聘请这位毕业生做一名专业的医疗诊断医生。虽然他知识渊博,但缺乏专业的医学训练。于是,医院需要对他进行"微调"——让他学习大量的医学案例、诊断报告和治疗方案。这个过程就是全面的微调。
- **微调的问题是什么?**问题在于,这个"医学深造"过程极其昂贵和耗时,对应到AI领域就是"巨大的计算和存储成本"。你需要动用海量的计算资源(GPU),花费大量时间,才能把一个通用模型训练成一个专家模型。
- 参数高效微调 (PEFT) / LoRA 是什么? 为了解决成本问题,人们发明了更聪明的方法,比如LoRA。这不再是让毕业生去读一个完整的医学院,而是给他一套高度浓缩的"医疗诊断核心手册"和"实践指南"(这就是LoRA适配器)。他只需要学习这套手册,就能在医疗诊断任务上表现得像个专家。这个过程只更新了很少的"知识"(参数),所以成本大大降低。
- LoRA 的核心困境:"知识绑定"问题。 这个方法很棒,但有一个致命缺陷:这套"医疗诊断核心手册"是用一种特殊的"速记符号"写成的,只有这位特定的毕业生能看懂。如果过了一年,学校里来了一位更聪明、基础知识更扎实的新毕业生(比如从Llama-3升级到Llama-4),你想让他也成为医疗诊断专家,你不能直接把旧手册给他,因为他看不懂。你必须为他重新编写一套全新的手册,这意味着之前的投入作废了。在AI模型日新月异的今天("新LLMs的快速发布"),这个问题尤为突出,它导致了大量的重复投资和资源浪费。

### 【原文翻译】

解决此限制的一个重要方法是知识蒸馏(KD)(Hinton et al., 2015; Azimi et al., 2024),它通过将目标模型的输出分布与源模型的输出分布对齐,将嵌入在源模型的PEFT适配器中的知识迁移到目标模型新的PEFT适配器中。然而,KD本质上依赖于数据,通常需要访问目标下游任务的训练数据(Nayak et al., 2019; Liu et al., 2024),这通常是不可用或获取成本高昂的。为了解决这个限制,TransLoRA(Wang et al., 2024)最近提出通过利用最

新LLMs的数据合成能力(Wang et al., 2023; Kim et al., 2025)来使用合成数据。这种方法使目标模型能够在不直接访问原始数据集的情况下获取领域知识。然而,TransLoRA需要训练额外的判别器模型来过滤低质量的合成数据,这不可避免地引入了额外的复杂性和计算开销。此外,它主要强调合成数据的作用,而较少关注知识迁移过程本身应该如何设计。

### 【深度解读】

这里,作者开始介绍学术界为了解决前面提到的"知识绑定"问题所做的尝试,并指出了它们的不足之处,从而为自己的方法(TITOK)铺平道路。

- 方案一:知识蒸馏 (Knowledge Distillation, KD)。 这是一种经典的"老师带徒弟"模式。已经学会"医疗诊断"的老毕业生(源模型/老师)是专家。我们让新毕业生(目标模型/徒弟)观察老师如何对各种病例进行诊断,并努力模仿老师的"诊断思路"(输出分布)。通过这种方式,徒弟可以学到老师的知识。但这个方案的缺陷很明显:徒弟需要看到和老师当年学习时一模一样的"教学病例"(原始训练数据)。在现实世界中,这些数据往往因为隐私、商业机密等原因,是无法再次获得的。
- 方案二: TransLoRA (合成数据法)。 为了解决"没有教学病例"的问题, TransLoRA提出了一个聪明的办法: 让老师(源专家模型)自己出一些"模拟病例"(合成数据)给徒弟学习。这样就不再需要原始数据了。这个想法很好,但它引入了新的复杂性:老师毕竟不是专业的出题人,他可能会出一些质量不高、甚至有误导性的"模拟病例"。为了保证学习效果,TransLoRA不得不额外聘请一位"教学督导"(判别器模型),专门负责审查老师出的模拟题,把质量差的题筛掉。这个"教学督导"本身也需要经过训练和维护,增加了整个流程的"复杂性和计算开销"。而且, TransLoRA把重点放在了"如何造出好题"上,却忽略了"如何更高效地教"这个核心问题。

# 【原文翻译】

**贡献。**在本文中,我们提出了一个通过词元级知识迁移实现有效LoRA移植的新框架(TITOK)。我们的高层次想法是通过使用词元级信号来指导迁移过程,而不是依赖于整个词元序列,从而有选择性地从源模型的LoRA中传达任务相关信息。我们通过引入一个我们称之为**对比盈余**的概念来具体捕获这些信息,该概念通过比较带有LoRA的源模型和不带有LoRA的同一模型之间的预测获得。直观地讲,这种对比盈余突出了包含重要任务知识的词元。这种盈余信号被进一步用于过滤我们生成的用于训练的合成数据,从而能够在包含更丰富信息的样本上进行选择性学习。与TransLoRA需要训练和存储额外的判别器模型进行过滤不同,TITOK既不需要额外的模型,也不需要额外的训练开销。此外,我们设计了一种有效的机制来解决源模型和目标模型之间的分词器不匹配问题,从而增强了我们框架的鲁棒性和适用性。

#### 【深度解读】

在指出了现有方法的不足之后,作者终于亮出了自己的"王牌"——TITOK。这段话是全文的核心思想的高度浓缩。

- **TITOK的核心理念**: 与其笼统地教,不如划重点。TITOK认为,老师的知识并非每个字都同等重要。在诊断报告中,"轻微头痛"和"颅内肿瘤"这两个词所包含的信息价值 天差地别。TITOK的目标就是精准地找到这些"知识闪光点"(包含重要任务知识的词元),然后让徒弟重点学习这些闪光点。
- 如何找到"知识闪光点"?——对比盈余。 这就是TITOK最巧妙的地方。它让同一个毕业生扮演两个角色:一个是没看过"医疗诊断手册"前的"新手"(不带LoRA的基础模型),另一个是看过手册后的"专家"(带LoRA的模型)。然后,让他们同时对一个病例进行诊断。当遇到某个关键症状描述时,"新手"可能会犹豫不决,而"专家"则会非常肯定地给出判断。这个从"犹豫"到"肯定"的巨大转变,就是"对比盈余"。这个盈余信号精确地标记出了手册中最核心、最关键的知识点。

### • TITOK的优势:

- 1. **无需额外模型:** 它不需要像TransLoRA那样聘请"教学督导"。它通过"自己和自己比"的方式,就完成了知识的筛选,非常轻量级。
- 2. **双重过滤:**它利用"对比盈余"信号,不仅能筛选出高质量的"模拟病例",还能在这些病例中进一步划出"核心考点"(关键单词),让学习过程更加聚焦和高效。
- 3. **强大的兼容性:** 作者还考虑到了一个工程细节——不同的模型可能"阅读习惯"不同(分词器不匹配)。TITOK设计了一个"翻译"机制,确保即使老师和徒弟的语言习惯不同,也能准确地传递知识重点。这大大增强了方法的实用性。

# 【原文翻译】

我们将通过对三个广泛使用的基准进行大量实验来证明TITOK的有效性,这些基准涵盖了推理(Big Bench Hard (Suzgun et al., 2022) 和 MMLU (Hendrycks et al., 2021))和个性化任务(LaMP (Salemi et al., 2024))。特别是,在所有任务和迁移设置中取平均值时,TITOK将性能比香草(Vanilla)目标模型提高了+7.96%,比KD提高了+6.0%,比TransLoRA提高了+4.4%。我们还探讨了各种迁移设置,包括在同一模型内、跨不同模型家族和大小,甚至跨不同模型版本之间的迁移。在每种情况下,我们的方法都取得了持续的改进。有趣的是,即使应用于源自与目标任务不同的外部数据,TITOK仍然有效,突显了其鲁棒性和普遍适用性。总体而言,这些经验结果突出了TITOK是一种方法论上简单而强大的范例,可用于在各种场景中高效地跨模型迁移LoRA知识。

#### 【深度解读】

这段话是在用实验结果为TITOK的强大效果背书。对于一篇科研论文来说,光有好的想法是不够的,必须用数据证明你的方法确实有效。

- 全面的考试: 作者为TITOK安排了一系列严格的"考试"。这些考试涵盖了不同类型:
  - **推理任务 (BBH, MMLU)**: 就像是给AI做的"逻辑思维和知识水平测试",考察其严谨的分析和推理能力。
  - 个性化任务 (LaMP): 就像是"创意写作测试",考察AI是否能模仿特定作者的风格进行创作,这检验了其对细微、个性化知识的学习能力。

- **令人信服的成绩:** TITOK的成绩非常出色,全面超越了所有对手:
  - **vs. 香草模型 (Vanilla)**: "香草"在这里指未经任何训练的"白板"目标模型。提升 7.96%说明TITOK的教学卓有成效。
  - 。 vs. 知识蒸馏 (KD): 提升6.0%说明TITOK比传统的"老师带徒弟"方法更有效。
  - 。 vs. TransLoRA: 提升4.4%说明TITOK比目前最先进的竞争对手还要强,而且 是在更简洁的框架下实现的。
- 严苛的"考场"环境: 作者还模拟了各种现实世界中可能遇到的复杂情况:
  - 。 **同模型迁移:** 比如从一个Mistral-7B模型迁移到另一个Mistral-7B模型。
  - 。 **跨家族迁移:** 比如从Mistral模型迁移到Llama模型,这相当于"跨校传授知识"。
  - 。 **跨尺寸迁移:** 从一个小的3B模型迁移到一个大的8B模型,相当于"初中老师的知识精华对高中生依然有启发"。
  - 。 **跨版本迁移:** 从旧的Llama-2模型迁移到新的Llama-3模型,这是最具现实意义的场景,证明了知识可以"代代相传",不会因为模型升级而作废。
- **意外的惊喜**: TITOK甚至在"跨界教学"中也表现出色——用一个任务的数据来帮助另一个完全不同的任务。这证明了TITOK学到的不仅仅是表面知识,而是更深层次、可泛化的"元知识"。

这部分内容不仅展示了TITOK的有效性,更重要的是,它揭示了AI发展的一个重要趋势: AI资产的生命周期管理。过去,一个微调好的模型就像一个消耗品,随着基础模型的淘汰 而失去价值。而TITOK这样的技术,将"知识"(LoRA适配器)本身塑造成了一种持久的、 可流动的资产。企业可以持续地将他们宝贵的领域知识从旧模型迁移到新模型上,实现AI 能力的不断迭代和进化,这从根本上改变了AI应用的经济模型和运营策略。

# 第二部分:背景知识与相关工作

#### 2 相关工作

#### 【原文翻译】

迁移PEFT适配器。参数高效微调(PEFT)(Houlsby et al., 2019; Li & Liang, 2021)已成为全模型微调的一种实用且流行的替代方案。通过仅要求更新一小部分参数,它实现了高效的适应,其中LORA(低秩适应)(Hu et al., 2021; Dettmers et al., 2023)作为最广泛采用的方法之一脱颖而出。然而,一个根本的限制是LoRA适配器与其训练的冻结骨干网络绑定,使得它们难以迁移到其他基础模型。为了解决这个问题,TransLoRA(Wang et al., 2024)等最近的研究试图通过生成合成数据(Wang et al., 2023)将LoRA适配器的知识跨模型移植。虽然在一定程度上有效,但这种方法需要一个额外的判别器模型来过滤高质量的合成数据,从而导致管道相对较重。与此同时,知识蒸馏(KD)(Hinton et al., 2015; Azimi et al., 2024)已被广泛探索作为另一种知识迁移方式;然而,传统的KD通常在教师-学生框架内以 Logit 或序列级操作,并且需要访问接近原始数据的训练数据,以便使

用教师的分布来监督学生。相比之下,我们的方法侧重于词元级选择性迁移,提供了一种更细粒度且轻量级的替代方案,可实现LoRA适配器在不同模型之间的高效移植以进行部署。

#### 【深度解读】

本段是对引言中提到的背景知识的进一步展开和梳理,将TITOK置于现有研究的坐标系中,以凸显其创新性。

- **PEFT 和 LoRA 的定位**: 作者首先再次强调,PEFT (特别是LoRA) 是当前AI领域的主流技术,因为它解决了全模型微调的"成本"痛点。这说明作者研究的问题是站在技术前沿的,是"巨人的肩膀上"的进步。
- **问题的重申**:接着,再次点明了LoRA的"阿喀琉斯之踵"——适配器与基础模型的"绑定"关系。这个问题之所以重要,是因为它直接关系到AI资产的复用性和长期价值。
- 现有解决方案的剖析:
  - TransLoRA:被描述为"管道相对较重"。"重"在这里是个关键词,意味着流程复杂、组件多、维护成本高。它需要一个"主厨"(生成模型)来做菜(合成数据),还需要一个"品菜师"(判别器模型)来尝菜的好坏。TITOK则希望做到"主厨自己就是最好的品菜师",从而简化流程。
  - 知识蒸馏 (KD): 被描述为在"Logit或序列级"操作,并需要"原始训练数据"。这指出了KD的两个局限性。一、粒度粗糙:"序列级"意味着老师教徒弟时,是整句话整句话地教,不够精细。二、数据依赖:再次强调了对原始数据的依赖是其在现实应用中的主要障碍。
- TITOK的差异化优势: 最后,作者用"相比之下"引出自己方法的特点——"词元级选择性迁移"、"更细粒度"、"轻量级"。这几个词精准地概括了TITOK的设计哲学:精准、高效、简洁。它不像KD那样笼统,也不像TransLoRA那样复杂,而是找到了一条兼具精度和效率的中间道路。

#### 【原文翻译】

使用LLMs进行数据合成。 使用LLMs进行合成数据生成作为减少对昂贵或不可访问数据集依赖的方法(Wang et al., 2023; Kim et al., 2025)越来越受到关注。先前的研究方向已将合成数据用于隐私保护(Bu et al., 2025)、数据增强(Kumar et al., 2021)和领域适应(Li et al., 2023)等目的。在我们的框架中,合成数据作为核心组件,使得LoRA适配器能够在不访问原始训练语料库的情况下进行迁移,同时减轻隐私问题并减少对外部数据集的依赖程度。此外,由于查询和标签都是由源专家模型本身直接生成的,合成数据最终提供了一种自给自足的机制,符合我们轻量级和有效知识迁移的目标。

#### 【深度解读】

这段话聚焦于TITOK框架中的一个关键"原料"——合成数据,并解释了为什么选择它以及如何使用它。

- **合成数据的价值:** 作者首先肯定了"让AI自己造数据"这一思路的普遍价值。它不仅仅是本文的权宜之计,而是一个广泛的研究领域。其主要好处有:
  - 解决数据稀缺/昂贵问题: 在很多专业领域,获取高质量的标注数据成本极高。
  - 。 **保护隐私:** 在医疗、金融等领域,原始数据高度敏感,无法直接使用。合成数据可以在保留数据分布特征的同时,隐去原始的个人信息。
  - 。 **增强模型能力**: 可以用于数据增强(扩充数据集)和领域适应(让模型适应新场景)。
- **合成数据在TITOK中的角色**: 在TITOK中,合成数据是实现"无数据迁移"的桥梁。它让知识的传递摆脱了对原始训练数据的束缚。
- TITOK的独特之处:"自给自足"。 这是本段的一个核心观点。许多方法可能需要一个模型生成问题(查询),另一个模型生成答案(标签)。而TITOK的设计是让"源专家模型" 同时生成问题和答案。这个设计的巧妙之处在于,它保证了问题和答案在风格、难度和知识范畴上的高度一致性,就像是同一位老师出的模拟题和标准答案,其教学效果自然更好。这种"自给自足"的机制,再次体现了TITOK追求简洁、高效的设计哲学,减少了系统对外部的依赖。

选择性词元训练。 最近的研究 (Lin et al., 2025; Gu et al., 2020) 表明,并非所有词元都对模型训练做出同等贡献,这推动了对选择性训练策略的研究。虽然此类方法主要应用于加速优化或减少冗余(Yeongbin et al., 2024; Bal et al., 2025) ,但我们的工作创新地将这一概念扩展到知识迁移的设置。具体来说,我们的盈余分数(公式2)的概念源于这一想法,其中源骨干网络与其LoRA适配器之间的对比产生了词元级判断。这使得TITOK能够以更集中和细粒度的方式移植LoRA知识,突显了选择性训练在其原始范围之外的更广泛适用性。

# 【深度解读】

这段话揭示了TITOK核心思想的理论来源,并阐明了其创新之处。

- 理论基础:"并非所有词元都生而平等"。 这是"选择性词元训练"这一研究领域的核心洞察。就像我们在学习时会划重点一样,模型训练也可以通过只关注那些信息量最大的词元来变得更高效。之前的研究主要将这个思想用于"加速训练"或"减少计算量",目标是"省钱省时"。
- **TITOK的创新:思想的"跨界应用"。** TITOK的作者敏锐地意识到,这个"划重点"的思想不仅能用来加速训练,还能用来指导"知识迁移"。这是一个非常漂亮的"思想迁移"。如果说加速训练是让学生"更快地"读完一本书,那么TITOK就是让老师告诉学生"这本书的重点在哪里",从而让学生"更好地"掌握知识。

- "对比盈余"的理论定位: 作者明确指出,他们提出的"对比盈余分数"就是实现"划重点"的具体工具。通过对比模型在有无LoRA时的表现差异,这个分数自然地量化了每个词元的重要性。分数越高的词元,就是LoRA知识体现得最淋漓尽致的地方,也就是最应该被迁移的"知识重点"。
- "青出于蓝而胜于蓝": 这段话的结尾,作者巧妙地提升了自己工作的价值。他们不仅是借鉴了"选择性训练"的思想,更是拓展了它的应用边界,证明了这个思想在"知识迁移"这个新场景下同样强大,甚至更有价值。

这部分"相关工作"的写作逻辑非常清晰:首先定义问题(迁移PEFT),然后剖析现有方案(TransLoRA, KD)的优劣,接着阐述自己所依赖的关键技术(合成数据),最后追溯自己核心思想的理论源头(选择性训练)并点明创新点。这套组合拳下来,读者能清晰地理解TITOK是在什么样的技术背景下诞生的,以及它的独到之处究竟在哪里。这体现了科学研究中一种深刻的优雅:TITOK的"对比盈余"机制并非凭空创造,而是将一个成熟领域的思想(选择性训练)巧妙地应用到一个新的问题领域(知识迁移),并为此设计了一套简洁而自洽的实现方案。 这种做法避免了像TransLoRA那样引入外部复杂组件(判别器),而是通过一种内部的、自监督的方式来生成指导信号。模型通过比较"新手"和"专家"的自己,完成了对知识的"自我审视"和"重点标记",这在工程和科学设计上都堪称典范。

第三部分:TITOK 框架详解

3 TITOK: 通过词元级知识移植LORA

# 【原文翻译】

在本节中,我们介绍TITOK,这是一个通过词元级知识迁移进行LORA移植的框架(图 1)。TITOK的核心思想是通过专门训练合成数据中的信息词元,将知识从源模型的LoRA 适配器迁移到目标模型的LoRA适配器中。具体来说,该框架由三个组成部分构成:(1)合成数据生成(第3.1节),其中源专家模型生成目标任务的查询-标签对;(2)盈余分数计算(第3.2节),它使用源模型计算词元级重要性;(3)带有过滤的目标模型训练(第3.3节),它使用新初始化的LoRA适配器在排名靠前的样本和词元上训练目标模型。此外,我们提出了盈余分数对齐(第3.4节),这是一种旨在即使源模型和目标模型的分词器不同也能应用TITOK的算法。TITOK的总体算法在算法1中介绍。

# 【深度解读】

这一段是本章的"导航图",清晰地勾勒出TITOK框架的全貌和核心步骤。它告诉我们,要完成一次"LoRA技能移植手术",需要分三步走,外加一个兼容性补丁。

• 核心思想重申:"专门训练合成数据中的信息词元"。这句话再次强调了TITOK的"划重点"策略。我们不是把整本书都给学生,而是只给他划了重点的笔记。

### • 三大核心步骤:

- 1. **造数据 (Synthetic data generation)**: 首先,得有学习材料。这一步是让"老师"(源专家模型)编写一本"模拟练习册"(包含问题和答案)。
- 2. **划重点 (Excess score computation):** 然后,老师要在这本练习册上划出核心考点。这就是计算"对比盈余分数"的过程,它会告诉我们哪些句子、哪些词是最重要的。
- 3. **教学生 (Target model training with filtering)**: 最后,把这本划好重点的练习 册交给"学生"(目标模型),让他只学习这些重点内容。
- 一个兼容性补丁 (Excess score alignment): 考虑到不同的学生可能有不同的"阅读习惯"(分词器不同),我们还需要一个"翻译器",确保老师划的重点能被学生准确理解。 这个三步走的流程,逻辑清晰,环环相扣,构成了一个完整的知识迁移闭环。

### 【图1描述】

图1:TITOK概览:通过词元级知识迁移进行移植。

这张图直观地展示了TITOK框架的整个工作流程:

- 1. **起点 (GOAL):** 流程始于一小部分(例如5个)"种子提示 (seed prompts)"。这就像是给老师布置了一个教学任务的主题。
- 2. **第一步:合成数据生成 (Synthetic Data Generation):** "源专家模型 (source base + LORA)"根据种子提示,生成大量的合成数据。图中示例为一句话:"y Local Man Buys Treadmill..."。
- 3. **第二步:对比盈余过滤机制 (Contrastive Excess Filtering Mechanism):** 这是 TITOK的核心。该机制将"专家模型 (source base + LORA)"的输出与"基础模型 (source base)" (即没有LoRA的同一个模型) 的输出进行对比,从而计算出每个词元的"盈余分数"。
- 4. 第三步: 两阶段过滤: 利用计算出的盈余分数, TITOK进行两轮筛选:
  - 。 **样本过滤 (Sample Filtering):** 首先,在样本(句子)层面进行筛选,只保留那些信息量最丰富的样本。
  - 。 **词元选择 (Token Selection):** 接着,在保留下来的样本内部,进一步筛选出信息量最丰富的词元。
- 5. **(可选) 第四步:分词器对齐 (Tokenization Alignment):** 流程图的说明指出,如果源模型和目标模型的分词器不同,训练前需要对筛选出的"重点标记"(掩码)进行对齐。

6. 第五步: LoRA知识迁移与目标模型学习 (LORA Knowledge Transfer & Target Model Learning): 经过筛选和对齐后的高质量数据,最终被用来训练目标模型上的一个全新的LoRA适配器 (target\_base + new LORA)。整个流程从生成数据,到通过对比盈余进行智能筛选,再到最终的精准训练,形成了一个高效的闭环,实现了知识从源LoRA到目标LoRA的有效移植。

### 3.1 通过LLM提示和少样本数据生成合成数据

### 【原文翻译】

设 Ms 表示源骨干LLM,As 表示其在目标下游任务上的LoRA适配器,它们共同构成了源专家模型 Ms+As。将要训练其LoRA适配器 At 的目标模型表示为 Mt。然后,TITOK首先构建一个类似于TransLORA(Wang et al., 2024)想法的合成数据集 Ds。这种合成数据的使用使我们能够避免保留整个原始下游任务数据集,同时让 At 学习编码在源适配器 As 中的知识。与TransLoRA使用未调优的目标模型 Mt 生成合成数据的方法不同,我们使用源专家模型 Ms+As 来合成数据(参见图2中的实证比较)。具体来说,Ms+As 在基于提示的数据合成框架(Wang et al., 2023)内合成查询和标签(详见附录H);给定下游任务的少样本数据,它首先生成一个查询q,然后以q为条件生成相应的标签y。为了鼓励多样性,我们对所有任务应用ROUGE-L过滤和去重,但少数此类过滤不可行的情况除外(更多细节在附录G中)。因此,得到的合成数据集包含查询-标签对:

 $Ds={(qj,yj)}j=1N(1)$ 

### 【深度解读】

本节详细阐述了TITOK流程的第一步:如何"制造"学习材料。

- 定义角色: 作者首先用数学符号明确了参与者:
  - 。 Ms: 源基础模型,即"老师"的基础知识。
  - 。 As: 源LoRA适配器,即老师掌握的"专业技能手册"。
  - 。 Ms+As: 源专家模型,即掌握了专业技能的"老师本人"。
  - 。 Mt: 目标基础模型,即"学生"的基础知识。
  - 。 At: 目标LoRA适配器,即我们希望为学生打造的新的"专业技能手册"。
- 核心策略:用"专家"造数据。 这里作者强调了与TransLoRA的一个关键区别。 TransLoRA是让"学生"(Mt)来提问,然后让"老师"(Ms+As)来回答。而TITOK是让"老师"自己既提问又回答。这个设计的深层考虑是,由同一个专家来设计问题和答案,可以确保两者在知识范畴、难度、风格和隐含假设上是高度一致的。这就像让出题老师自己提供标准答案和解题思路,教学材料的"内在一致性"更高,学生的学习效率自然也更高。后续的实验(图2)也证明了这一设计的优越性。

#### • 具体操作:

- 1. **少样本提示 (Few-shot data)**: 我们不是让老师凭空出题,而是先给他几个"样题"(少样本数据),让他明白要出什么类型、什么格式的题目。
- 2. **生成问答对:** 老师先生成一个问题(查询 q),然后再基于这个问题生成标准 答案(标签 y)。
- 3. **保证多样性:** 为了避免老师反复出同一类型的题目,导致学生"刷题"疲劳且知识面狭窄,作者还设置了两个"防重复"机制:ROUGE-L过滤(防止内容过于相似)和去重(防止完全一样)。 最终,我们就得到了一套高质量的、由专家亲自编写的、多样化的"模拟练习册" Ds。

# 【原文翻译】

通过这种方式,Mt 将在源专家模型 Ms+As 生成的合成数据 DS 上进行训练,从而无需依赖整个原始数据集即可实现知识迁移。

### 【深度解读】

这是一个简短的承上启下的段落,总结了3.1节的核心目的和成果。它强调了这一系列操作的最终价值:我们成功地创造了一种学习场景,使得"学生"(Mt)可以在不接触任何"原始、保密的教学材料"的情况下,仅通过学习"老师"(Ms+As)编写的"模拟练习册"(DS),就能学到其核心技能。这为知识在不同AI模型之间的安全、高效流动奠定了基础。

### 3.2 来自带有LoRA适配器的源模型的对比盈余分数

#### 【原文翻译】

如第3.1节所述,TITOK依赖于合成数据,但合成数据通常容易出现不完美(Chen et al., 2024);因此,复杂的过滤对于保留高质量的信息样本至关重要。TransLoRA(Wang et al., 2024)通过一个单独的判别器来解决这个挑战,以过滤有用的查询,但这引入了训练和维护额外模型的额外负担。相比之下,我们提出了一种轻量级的替代方案,仅利用已经训练好的源模型。具体来说,我们使用源模型及其LoRA适配器来执行两个互补的角色:(1)业余角色(Ms)和(2)专家角色(Ms+As)。然后,这两个角色之间的差异提供了一个隐式监督信号,其中编码了任务信息。

### 【深度解读】

本节进入TITOK框架的心脏地带——如何智能地"划重点"。作者首先摆出了一个现实问题:即便是专家出的"模拟题",也可能存在质量参差不齐的情况。因此,一个高效的"筛选机制"是必不可少的。

- 问题的再次聚焦: 合成数据虽好,但并非完美。如何沙里淘金,找出最有价值的学习内容?
- 对比TransLoRA的"笨办法": TransLoRA的做法是外聘一个"质检员"(判别器模型)。这种方法虽然有效,但成本高、流程复杂,属于"重量级"解决方案。

- **TITOK的"巧办法":** TITOK的方案则体现了极大的智慧和简洁性。它不向外求助,而是向内挖掘,让源模型"一人分饰两角":
  - 1. **业余选手 (Ms):** 这是没有经过专业训练的"通才"模型。
  - 2. **专家选手 (Ms+As):** 这是装备了LoRA"专业技能手册"的"专才"模型。
- 核心洞察:差异即信息。"业余"和"专家"之间的表现差异,本身就是一种极其宝贵的"监督信号"。这个差异精确地指出了LoRA适配器到底在哪些地方、以何种方式提升了模型的能力。这就像通过比较一个医学生和一个资深医生的诊断过程,我们能最快地发现资深医生经验的价值所在。这种"自我对比"的思路,是TITOK轻量级和高效的根源。

形式上,设 y=[y1,...,yL] 表示合成响应,其中 yi 是y中对应于合成查询q的一个词元。然后,我们将盈余分数定义为:

S(yi)=Le(yi)-La(yi)(2)

其中词元 yi 上的业余和专家损失定义为:

La(yi) = logPMs(yi | q,y < i)(3a)

Le(yi)=logPMs+As(yi|q,y<i)(3b)

#### 【深度解读】

这里,作者将前面"业余与专家"的直观思想,用严谨的数学语言表达出来。

• 损失 (Loss) 的直观理解: 在语言模型中,"损失"可以被通俗地理解为模型对下一个正确词元的"惊讶程度"。如果模型很确定下一个词应该是什么,而正确答案也确实是那个词,那么它的"惊讶程度"就很低(损失小)。反之,如果模型对下一个词毫无头绪,那么当正确答案出现时,它的"惊讶程度"就很高(损失大)。这里的 P(yi|...) 代表模型预测下一个词是 yi 的概率,取对数 log 后,概率越高,损失值(的绝对值)就越小。

#### • 公式解读:

- 。 La(yi): 这是"业余选手"(Ms)在预测第 i 个词元时的"惊讶程度"。
- 。 Le(yi): 这是"专家选手"(Ms+As)在预测同一个词元时的"惊讶程度"。
- 。 S(yi): "对比盈余分数"就是这两者"惊讶程度"的差值。

• "尤里卡时刻"的量化: 一个很高的 S(yi) 分数意味着什么?它意味着 Le(yi) 远大于 La(yi)。换句话说,"业余选手"在预测这个词时非常"惊讶"(不确定),而"专家选 手"则非常"不惊讶"(胸有成竹)。这个巨大的转变,正是LoRA适配器发挥关键作用 的"尤里卡时刻"!这个词元很可能蕴含了任务的核心知识,是"业余"和"专家"的分水 岭。TITOK通过这个简单的减法,就精准地捕捉到了这些知识的闪光点。

### 【原文翻译】

盈余分数 S(yi) 量化了因装备LoRA而产生的知识差异,从而识别出适配器做出决定性贡献的词元。直观地说,如果骨干模型不确定如何预测一个词元,但经过LoRA增强的模型以高置信度分配给它,那么该词元将获得一个大的盈余分数。这意味着具有较高 S(yi) 的词元对应于LoRA适配器注入了骨干模型自身无法捕捉的任务特定知识的位置。通过这种方式,盈余分数充当了一种细粒度的归因信号,完全源自模型自身的内部行为,并引导训练朝向数据中那些最富含适配器知识的特定区域。

### 【深度解读】

这段话是对"对比盈余分数"物理意义的精彩阐述,进一步加深了我们的理解。

- **盈余分数的本质:** 它是"知识增量"的量化指标。它告诉我们,LoRA适配器这本"专业手册",具体在哪一个字、哪一个词上,为模型带来了知识的飞跃。
- **高分词元的特征:** 高分词元是"新手"的知识盲区,却是"专家"的拿手好戏。它们是区分通用知识和专业知识的关键节点。例如,在医疗诊断任务中,对于"患者主诉头痛"这句话,"头痛"这个词可能对于业余和专家来说都很平常(分数不高)。但对于后续的描述"呈蛛网膜下腔出血样剧痛",业余模型可能会感到困惑,而专家模型则会高度确定这是一个关键指征,这个词组的盈余分数就会非常高。
- 信号来源的优越性:"内部行为"。 作者强调,这个信号是"完全源自模型自身的内部行为"。这一点至关重要,因为它意味着我们不需要任何外部的标注、评判或额外的模型。系统通过"自省"和"自我对比",就完成了对知识重要性的判断。这是一种极其优雅和高效的"自监督"机制,也是TITOK相比于TransLoRA等方法的根本优势所在。它将复杂的知识筛选问题,转化为一个简单的、内部的信号比较问题,大大降低了系统的复杂度和实现成本。

#### 3.3 带有样本过滤和词元选择的目标模型训练

#### 【原文翻译】

在计算出盈余分数 S(yi) 之后,用于目标模型 Mt 的新初始化的LORA适配器 At 将使用两级过滤方案在合成样本 (qi,yi)∈Df 上进行训练。

**第一阶段:样本过滤。** 我们首先在样本级别过滤合成数据集 Ds (公式1) ,以移除信息量较少的示例。对于每个合成样本,我们计算y中词元的盈余分数 S(yi) 的平均值,并仅保留具有最高值的M个样本。

 $Sj=|yj|1yi\in yj\Sigma S(yi)(4)$ 

令 Df 为 Ds 中具有最大 Si 的M个样本的集合。

 $Df=TopM{(qj,yj)\in Ds:Sj}(5)$ 

通过这一步,合成数据经过了过滤过程,确保后续训练专门集中在 Df 中具有更丰富知识信号的剩余示例上。

### 【深度解读】

有了"对比盈余"这把锋利的"手术刀",现在我们开始进行两阶段的"精细筛选"。本节介绍的是第一阶段,可以比喻为"筛选优质章节"。

• **目标**: 在我们让老师编写的整本"模拟练习册"中,有些章节可能写得特别好,知识密度很高;有些则可能比较平庸。第一阶段的目标就是找出那些"精华章节"。

#### 方法:

- 1. **计算章节平均分:** 对于练习册中的每一个"样本"(可以理解为一个问答对,或一个段落),我们计算其中所有词元的"对比盈余分数"的平均值 Sj。这个平均分代表了这个样本整体的"知识含金量"。
- 2. **择优录取**: 我们根据这个平均分对所有样本进行排序,只保留得分最高的M个样本。这就构成了一个"精华版"的数据集 Df。
- **效果**: 经过这一轮筛选,我们确保了用于训练的学习材料都是"优中选优"的。这避免了模型在大量低质量、低信息量的样本上浪费时间和计算资源,使得后续的训练过程更加高效和聚焦。这就像一个聪明的学生,他不会把整本参考书从头到尾都看一遍,而是会先根据目录和前言,找出最有价值的几个章节进行精读。

### 【原文翻译】

第二阶段:词元选择。 接下来,我们考虑词元选择;也就是说,At 不会从保留样本的所有词元中学习。相反,它只关注那些被盈余分数 S(yi) 优先排序的、被认为对知识迁移最重要的词元。为了实现这一点,我们使用指示符 Ik%(yi) 选择按其盈余分数排名的前k%的词元: \$\$ I\_{k%}(y\_{i})=\begin{cases}1,& \text{if } \text{rank}{y/{j}}(S(y\_{i}))\le\lfloor k%\cdot|y\_{j}|\rfloor\ 0,& \text{otherwise,}\end{cases} \quad (6) \$\$ 其中 |y| 表示 yj 中的词元数量,而 rankyj(S(yi)) 表示 S(yi) 在该响应的词元中的排名。基于此选择,At 的训练目标定义为: \$\$ \mathcal{L}{\text{TiTok}}=\sum{(q\_{j},y\_{j})\\in\mathcal{D}{f}}\sum{y\_{i}\in y\_{j}}\lin y\_{j}}\lin y\_{j}}\lin \quad (7) \$\$ 其中 Lt(yi) 是由 Mt+At 在词元 yi 上分配的负对数似然损失(只有 At 是可学习的)。通过仅在这些过滤后的词元上进行训练(公式7),TITOK使 At 能够有效地获取源LoRA的知识,而无需访问原始训练数据或任何外部模型。

#### 【深度解读】

这是两阶段筛选的第二阶段,也是更精细的一步,可以比喻为"在精华章节中划出核心考点"。

• **目标**: 在已经筛选出的"精华章节"中,并非每个字都同等重要。这一步的目标是精确地找出那些"核心定义"、"关键公式"和"点睛之笔"。

### • 方法:

- 1. **内部排序和筛选:** 在每个样本内部,我们再次根据每个词元的"对比盈余分数" S(yi) 进行排序。然后,我们只保留得分最高的前 k% 的词元。例如,如果 k=30,我们就只保留每个句子中最重要的30%的词。
- 2. **生成"学习面具":**公式(6)中的 lk%(yi) 就是一个"学习面具"。对于被选中的"重点词元",它的值为1(表示"需要学习");对于其他词元,它的值为0(表示"可以忽略")。
- **最终的训练目标 (公式7)**: 这个公式是整个训练过程的"指挥棒"。它告诉目标模型 (学生):
  - 。 ∑(qj,yj)∈Df: 我们只在"精华数据集" Df 上学习。
  - 。 ∑yi∈yj: 对于数据集中的每个词元...
  - 。 lk%(yi)·Lt(yi):...我们只计算那些被"学习面具"标记为1的词元的损失。也就是说,**模型只对那些被划为重点的词元进行学习和优化,完全忽略其他词元。**
- **效果**: 这种"聚焦式学习"极其高效。它迫使模型将全部的"注意力"和"学习能力"都投入 到最关键的知识点上,避免了在无关紧要的信息上"分心"。这不仅大大加快了学习速 度,也提高了学习的质量,使得知识迁移更加精准和深刻。

# 【算法1描述】

算法1: TiTok: 通过词元级知识移植LoRA

- 输入: 源专家模型 Ms+As,目标模型 Mt,参数 N, M, k%。
- **输出:** 训练好的目标LORA At。
- 1. 用 Ms+As 构建包含N个样本的合成数据集 Ds={(qj,yj)}j=1N。
- 2. **对于** Ds 中的每个样本 (qi,yi) **执行**:
  - 。 计算词元盈余分数 S(yi)=Le(yi)−La(yi)。
  - 。 计算平均分数 Sj=|yj|1∑yi∈yjS(yi)。
- 3. 根据 Sj 选择前M个样本,形成 Df。
- 4. 对于 Df 中的每个样本 (qi,yi) 执行:

根据 S(yi) 对词元进行排序,并保留前k%,用掩码 lk%(yi)表示。

- 5. **如果** 源模型的分词器(Ms) □= 目标模型的分词器(Mt) **则**: 对齐掩码 lk%(s)(yi)→lk%(t)(yi)。
- 6. 在 Mt 上使用带掩码的损失函数 LTiTok=∑(qj,yj)∈Df∑yi∈yjlk%(yi)·Lt(yi) 来训练 At。
- 7. **返回** At。

# 3.4 跨不同分词器的盈余分数对齐

### 【原文翻译】

在具有不同分词器的模型之间进行迁移时,直接的词元级信号映射是不可能的,因为源模型和目标模型可能以不同的方式分割文本。为了解决这个问题,我们引入了一种简单而鲁棒的分词器对齐算法,该算法将词元掩码(公式6)从源词元序列 y(s) 传播到目标词元序列 y(t)。该算法首先使用双指针来对齐词元序列,这两个指针逐步解码和匹配文本片段。然后使用以下四个规则传播掩码:(1)一对一映射直接复制,(2)一对多映射进行复制,

(3) 多对一映射取平均值, (4) 多对多映射取平均值后复制。最后,一个top-k%选择步骤保留了最置信的目标词元。这个过程确保了跨分词器的一致监督,使得即使模型以不同方式对文本进行分词,也能进行可靠的迁移。该过程的概念性图示见图4。

### 【深度解读】

本节解决的是一个非常实际的工程问题,它决定了TITOK能否在复杂的现实世界中被广泛应用。

• 问题的根源:"语言习惯"的差异。 不同的语言模型家族(如Llama和Mistral)可能有不同的"分词器",这就像两个人的"阅读习惯"不同。例如,对于英文单词"tokenization",源模型可能将其看作一个整体 [tokenization],而目标模型可能将其看作三个部分 [token, ##iza, ##tion]。现在,如果源模型认为 [tokenization] 这个词元很重要(掩码为1),我们该如何将这个"重要性"信息传递给只认识 [token, ##iza, ##tion] 的目标模型呢?

- TITOK的解决方案:一个聪明的"翻译和对齐"算法。
  - 1. **文本对齐:** 算法首先不看词元,而是看它们解码后的实际文本。它通过"双指针"技术,找到源词元序列和目标词元序列中能对应上同一段文本的片段。比如,它会发现源模型的 [tokenization] 和目标模型的 [token, ##iza, ##tion] 都对应着"tokenization"这个字符串。
  - 2. **重要性传播规则**:找到对应关系后,就按照简单直观的规则来传递"重要性分数"(即掩码值):
    - 一对一 (one-to-one):源的一个词元对应目标的一个词元,分数直接照搬。
    - 一对多 (one-to-many): 源的一个词元(如 [tokenization])对应目标的多个词元(如 [token, ##iza, ##tion]),那么源的重要性分数就被"复制"给目标的每一个对应词元。
    - **多对一** (many-to-one):源的多个词元对应目标的一个词元,那么源的多个分数就被"平均"后赋给目标词元。
    - 多对多 (many-to-many): 结合上述规则,先平均再复制。
- **最终效果**: 通过这个过程,源模型划的"重点"被准确无误地"翻译"并标记在了目标模型的文本上。这确保了即使两个模型的"语言习惯"不同,知识迁移的核心指导信号也不会丢失或错位,极大地增强了TITOK框架的鲁棒性和通用性。这就像一个细心的老师,不仅划了重点,还为使用不同版本教材的学生准备了重点页码对照表。

# 第四部分:实验设置与结果分析

#### 4 实验

#### 【原文翻译】

在本节中,我们展示我们的实验结果,以回答以下研究问题:

- RQ1: TITOK能否在各种场景中有效地迁移LoRA的知识? (表1)
- RQ2: TITOK中每个组件的贡献是什么? (表2)
- RQ3: 词元级选择性迁移对选择比例有多敏感? (图3)
- RQ4: 选择哪个模型来合成查询对性能有何影响? (图2)
- RQ5: TITOK能否使用来自不同或不相关领域的数据来迁移知识? (表3)

#### 【深度解读】

一个好的科学研究,不仅要展示"我做到了",还要解释"为什么能做到"、"每个部分起了什么作用"、"在不同条件下表现如何"以及"它的极限和潜力在哪里"。本节的这五个研究问题(Research Questions, RQs)就完美地体现了这种科学探索精神。

- RQ1 (有效性): 这是最基本的问题——你的方法到底行不行?
- **RQ2 (贡献度/消融实验):** 你的方法之所以有效,是因为所有部分都不可或缺,还是某个部分起了决定性作用?这就像拆解一台机器,看看拿掉某个零件后机器是否还能运转。
- **RQ3 (敏感性分析)**: 你的方法是否对某个参数(比如选择词元的比例k%)特别敏感? 一个好的方法应该是鲁棒的,在一定参数范围内都能表现良好,而不是一个需要精调"祖传参数"的脆弱方案。
- **RQ4 (设计选择的验证)**: 你在设计中做出了一个关键选择 (用源模型而不是目标模型 生成查询) ,这个选择是正确的吗?用数据说话。
- **RQ5 (泛化性/灵活性):** 你的方法是否只能依赖自己生成的"合成数据"?它能否适应更广泛的、来自外部的数据源?这决定了它在现实世界中的应用广度。 通过回答这五个问题,作者将为我们全方位、多角度地展示TITOK的性能、鲁棒性和巨大潜力。

### 4.1 实验设置

### 【原文翻译】

模型。 我们主要使用Mistral和Llama家族的模型来展示我们的知识迁移实验。具体来说,我们设计了以下LoRA迁移(源  $\rightarrow$  目标)设置。(1) Mistral-7B-Inst-v0.3  $\rightarrow$  Mistral-7B-Inst-v0.3 : 基础迁移设置,(2) Mistral-7B-Inst-v0.3  $\rightarrow$  Llama-3.1-8B-Inst : 不同模型家族迁移设置,(3) Llama-3.2-3B-Inst  $\rightarrow$  Llama-3.1-8B-Inst : 不同尺寸模型迁移设置,以及(4) Llama-2-7b-chat-hf  $\rightarrow$  Llama-3.1-8B-Inst : 不同版本模型迁移设置。这些设置旨在测试一个较小的模型能否有效地将知识迁移到一个较大的模型,并探索一个相对较弱的模型是否仍能影响一个更新、更强的模型。这些多样的设置真实地反映了当今LLM的发展方式,即各种模型被发布,更新、改进的版本不断涌现。

# 【深度解读】

实验设置的模型选择部分,充分体现了作者对实验严谨性和现实意义的考量。他们设计的四种迁移场景,就像是四个难度递增的"考场",全面检验TITOK的能力。

- 场景一:同校同级转专业 (Mistral 7B → Mistral 7B)。 这是最简单的"基础模式",用于验证方法的基本盘是否稳固。在架构完全相同的情况下,知识迁移理应最容易成功。
- 场景二: 跨校交流 (Mistral 7B → Llama3 8B)。 这是"困难模式",考验方法的"通用性"。Mistral和Llama是两个不同的模型家族,它们的内部架构、训练数据和"思考方式"都有差异。如果TITOK在这种情况下依然有效,说明它学到的不是特定于某个模型的"小技巧",而是更通用的"知识原理"。

- 场景三:以小教大 (Llama3 3B → Llama3 8B)。 这是对方法"可扩展性"的考验。通常我们认为知识应该从强到弱传递,但现实中我们常常希望用一个轻量、便宜的小模型上学到的知识,去"启发"或"引导"一个更强大的大模型,从而节省大模型的训练成本。这个实验就是为了验证这种"四两拨千斤"的可能性。
- 场景四:继往开来 (Llama2 7B → Llama3 8B)。 这是最贴近现实应用需求的"终极考验"。在产业界,模型版本总在不断迭代。如果一个在旧版Llama2上投入巨大成本微调出的专业能力,能够低成本地迁移到新版Llama3上,并让Llama3"赢在起跑线上",这将具有巨大的商业价值。这个实验直接关系到TITOK能否成为AI资产"保值增值"的利器。

基线。 为了证明TITOK的有效性,我们将其与三个基线进行比较:(i) Vanilla,未经任何微调的目标基础模型;(ii) KD(+MinED),从源专家模型进行的知识蒸馏(KD)(Hinton et al., 2015)。当源模型和目标模型使用不同的分词器时,原始的KD不适用。在这种情况下,我们使用最小编辑距离(MinED)分词器对齐(Wan et al., 2024),它通过动态规划序列对齐来对齐词元序列和词汇分布,以处理近似匹配(例如,"gets"→"get"),并通过将概率质量映射到最近的编辑距离邻居来处理其他情况(例如,"immediately"→"immediate")。我们使用TransLoRA生成的合成数据作为KD和MinED的训练数据;以及(iii)TransLoRA(Wang et al., 2024),这是一种先前的方法,其中vanilla目标模型合成查询,而源模型及其LoRA适配器生成相应的合成标签。然后使用一个判别器来过滤合成数据,用于训练目标LoRA适配器。

# 【深度解读】

选择合适的"参照物"(基线)是实验成功的关键。作者在这里选择了三个非常有代表性的基线,构成了一个从"下限"到"最强对手"的完整比较链条。

- 基线一: Vanilla (白板模型)。 这是"零分基准"。任何一个知识迁移方法,如果连原始的、未经训练的目标模型都比不过,那它就毫无价值。这个基线用来衡量TITOK带来的"绝对提升"。
- 基线二: KD (经典方法)。 这是"行业标准"。知识蒸馏是知识迁移领域最经典、最广为人知的方法。战胜它,意味着TITOK在效果上超越了传统范式。作者还非常严谨地考虑了分词器不匹配的情况,为KD配备了MinED这个"外援",确保了比较的公平性。
- 基线三: TransLoRA (最强竞品)。 这是"SOTA (State-of-the-Art) 基线",即当前解决同一问题的最先进的方法。在学术研究中,仅仅战胜经典方法是不够的,必须证明你的方法比目前最好的还要好。与TransLoRA的直接对话,是证明TITOK创新价值的关键一役。 通过与这三个基线的比较,作者可以清晰地定位TITOK的性能水平:它不仅有效(超越Vanilla),而且优于经典方法(超越KD),甚至比当前最强的竞争对手还要更胜一筹(超越TransLoRA)。

#### 【原文翻译】

**数据集**。 遵循TransLoRA的先前工作,我们首先在两个代表性基准上进行实验: (1) Big-Bench Hard (BBH) (Suzgun et al., 2022) 和 (2) 大规模多任务语言理解 (MMLU) (Hendrycks et al., 2021)。BBH包含27个具有挑战性的推理任务,结构为多项选择或简答 题,旨在测试语言模型的组合泛化和高级问题解决能力。同时,MMLU涵盖57个不同学术 科目的任务,以多项选择题的形式呈现,以评估模型的广泛知识和推理能力。由于这两个 基准只提供测试集,我们将数据划分为90%用于训练源专家模型 Ms+As,其余10%用于评 估。为了将我们的方法扩展到个性化和文本生成,我们还在LaMP基准 (Salemi et al., 2024) 上进行了实验,专门关注其生成任务。特别是,我们实验了(3) 新闻标题生成 (LaMP 4) 和 (4) 学术标题生成 (LaMP 5), 因为它们是既可访问又能可靠评估的仅有的文本 生成任务。其余的LaMP任务被排除,因为LaMP 1-3是判别性任务,而LaMP 6-7缺乏黄金 标签。对于LaMP任务,源专家模型 Ms+As 在活动历史最长的30个用户的数据上进行训 练。对于每个用户,我们使用200个数据点进行训练,50个数据点进行验证,这一设计选择 旨在进行更严格和鲁棒的评估。总的来说,每个LaMP任务包含6000个训练样本和1500个 评估样本。为了评估BBH和MMLU的性能,我们使用LM-Eval Harness (Gao et al., 2024) 测量平均准确率。遵循TransLoRA (Wang et al., 2024) 中使用的设置,所有任务都在零样 本设置下进行。同时,对于LaMP任务,我们采用ROUGE-1和ROUGE-L分数作为评估指 标,遵循该基准的主要评估指标。

### 【深度解读】

实验所用的"考卷"(数据集)的选择,同样体现了全面性和挑战性。

- 考验"智商"的考卷 (BBH, MMLU): 这两个数据集是AI领域的"奥数"和"高考",以其难度和广度著称。
  - 。 BBH (Big-Bench Hard): 专注于那些对当前大模型来说"最难"的问题,是检验模型逻辑推理、组合泛化等"硬核"能力的试金石。
  - 。 MMLU: 覆盖了从高中到专业级别的57个学科,是检验模型知识广度和深度的"百科全书式"测试。 在这两份考卷上取得好成绩,意味着TITOK迁移的是深层次的"智慧"而不仅仅是"知识"。
- 考验"情商"和"文采"的考卷 (LaMP): LaMP任务则转向了一个完全不同的维度——个性化。

新闻标题生成 & 学术标题生成: 这两个任务要求模型学习并模仿特定用户的写作风格。这考验的是模型对细微、非结构化、带有个人色彩的知识的捕捉和复现能力。如果说BBH/MMLU是理科考试,那么LaMP就是文科创作。

#### • 评估指标的合理性:

- 。 对于有标准答案的BBH和MMLU,使用"准确率"进行评估,简单直接。
- 。对于文本生成任务LaMP,使用ROUGE分数,这是衡量生成文本与参考文本在内容重叠度上的金标准。通过在这些性质迥异的数据集上进行测试,作者旨在证明TITOK是一种通用的知识迁移框架,无论迁移的是逻辑推理能力还是写作风格,它都能胜任。

**实现细节。**对于源模型和目标模型的训练,我们使用  $5\times10-5$  的学习率,训练2个epoch,批大小为4,并应用LORA,其中秩 r=8,缩放因子  $\alpha=8$ ,丢弃率为0.05。优化使用 AdamW,权重衰减为  $1\times10-2$ ,并结合线性学习率调度和0.1的预热比例。对于合成数据生成,我们提供原始训练数据中的五个样本作为少样本范例,并应用top-p采样(Holtzman et al., 2020)来生成查询和标签,采样超参数为每个任务单独调整。为了进一步鼓励多样性,我们应用阈值为0.7的ROUGE-L过滤和去重来移除冗余查询,遵循Wang et al. (2023)。对于无法进行基于ROUGE过滤的任务,我们仅应用去重(见附录G)。对于初始合成池,我们生成2M(原文应为2倍)合成样本,并用对比盈余分数过滤后,保留前M个,其中M等于源训练集的大小(见附录F)。对于词元选择,选择比例k%在所有任务和迁移设置中固定为70%,除了Llama3  $3B \to Llama3 8B$ 的迁移设置,我们发现 k%=30% 始终产生最佳性能。对于推理,我们采用贪婪解码以确保确定性评估。

### 【深度解读】

这段是论文的"配方表",详细列出了复现实验所需的所有参数和设置。对于科研来说,这一部分至关重要,因为它保证了研究的"可复现性"——其他研究者可以按照这个"配方"做出同样的结果。对于我们读者来说,有几个关键信息值得关注:

- 训练参数 (学习率、epoch、batch size等): 这些都是训练神经网络时的标准"旋钮"。 作者选择了领域内常用的典型值,表明他们没有在这些基础参数上做过多的"黑魔 法"式调优,使得结果更具普遍性。
- **LoRA参数 (r=8, α=8):** 秩r是LoRA中最重要的参数,决定了"技能手册"的"厚度"。r=8 是一个非常典型且轻量级的设置,表明作者是在一个资源高效的场景下进行的实验,这符合PEFT的初衷。

#### • 合成数据生成细节:

- 5个少样本范例: 启动数据生成过程只需要极少的"种子",成本很低。
- **Top-p采样**: 这是一种能增加生成内容多样性同时又保持质量的常用技术。
- 。 **ROUGE-L过滤 (阈值0.7):** 这是一个具体的"防重复"措施,保证了生成数据的多样性。

#### • 过滤和选择参数:

- 。 **保留M个样本:** 过滤后的数据集大小与原始训练集相当,保证了训练数据量的 公平性。
- 。 k%=70% (大部分情况) / 30% (特殊情况): 这是"划重点"的比例。大部分情况下,保留70%最重要的词元效果最好。但在"以小教大"的场景下,保留30%效果更佳。这是一个非常有趣的发现,暗示了当"老师"能力不如"学生"时,更严格的"划重点"(只教最精华的部分)反而效果更好,可以避免老师的"噪声"知识对更强的学生造成干扰。 这些细节共同描绘了一幅严谨、规范、可信的实验图景。

表1总结了在四种迁移设置下,BBH、MMLU和LaMP(新闻标题和学术标题生成)的实验结果。首先,我们观察到在同一模型家族内的迁移非常有效;当将在Mistral 7B上训练的LoRA适配器移植到同一个模型的全新实例中时,TITOK始终超越所有基线。特别地,TITOK在所有任务上的平均性能比vanilla模型提高了24.08%,并且分别比KD和TransLoRA基线高出19.63%和4.91%。总而言之,这些发现表明,作为第一步,同一家族内的迁移是可靠成功的。

### 【深度解读】

这里开始对核心实验结果(表1)进行解读。作者首先分析了最简单、最基础的场景:**同模型迁移 (Mistral 7B**  $\rightarrow$  Mistral 7B)。

- 结论: 在这种"理想情况"下, TITOK表现优异, 大幅超越所有对手。
- 数据支撑:
  - 。 vs. Vanilla (提升24.08%): 巨大的提升证明了知识迁移的必要性和有效性。
  - 。 vs. KD (提升19.63%): 大幅超越经典方法,说明TITOK的"划重点"教学法远比 KD的"囫囵吞枣"式模仿更高效。
  - 。 vs. TransLoRA (提升4.91%): 显著超越了最强的竞争对手,证明了TITOK在效果上的先进性,而且是在更轻量的框架下实现的。
- **意义**: 这个结果为整个实验奠定了坚实的基础。它证明了TITOK的核心机制是正确且有效的。如果连最简单的场景都无法成功,那么后续更复杂的场景也就无从谈起。这就像是在证明一个新发明的引擎,首先要在试验台上让它成功运转起来。

#### 【表格1:主要结果】

下表展示了在BBH、MMLU、新闻标题生成和学术标题生成任务上,四种不同迁移设置下的主要实验结果。BBH和MMLU是推理任务,使用LM-Eval Harness进行评估(准确率Acc.);新闻标题和学术标题生成是个性化任务,使用ROUGE-1/L进行评估。所有评估均为零样本设置。最优分数以**粗体**显示。

迁移设置	方法	BBH Acc.	MMLU Acc.	新闻标题 ROUGE-1/L	学术标题 ROUGE-1/L
Mistral 7B → Mistral 7B	Vanilla	0.397	0.557	0.117 / 0.101	0.381 / 0.311
	KD	0.426	0.556	0.107 / 0.121	0.392 / 0.320
	TransLoRA	0.424	0.533	0.155 / 0.136	0.447 / 0.382

迁移设置	方法	BBH Acc.	MMLU Acc.	新闻标题 ROUGE-1/L	学术标题 ROUGE-1/L
	TiTok (ours)	0.432	0.563	0.160 / 0.142	0.473 / 0.414
Mistral 7B → Llama3 8B	Vanilla	0.469	0.469	0.125 / 0.110	0.444 / 0.378
	KD+MinED	0.477	0.484	0.127 / 0.112	0.455 / 0.389
	TransLoRA	0.471	0.472	0.108 / 0.122	0.461 / 0.397
	TiTok (ours)	0.482	0.488	0.140 / 0.124	0.464 / 0.403
Llama3 3B → Llama3 8B	Vanilla	0.469	0.469	0.125 / 0.110	0.444 / 0.378
	KD	0.470	0.475	0.126 / 0.111	0.449 / 0.383
	TransLoRA	0.460	0.466	0.121 / 0.107	0.455 / 0.387
	TiTok (ours)	0.509	0.475	0.127 / 0.113	0.457 / 0.392
Llama2 7B → Llama3 8B	Vanilla	0.469	0.469	0.125 / 0.110	0.444 / 0.378
	KD+MinED	0.473	0.478	0.111 / 0.123	0.450 / 0.384
	TransLoRA	0.472	0.468	0.125 / 0.109	0.453 / 0.388
	TiTok (ours)	0.510	0.479	0.140 / 0.122	0.461 / 0.404

除了同家族迁移,我们还发现TITOK在跨模型迁移设置中非常有效。例如,当从Mistral 7B 迁移到Llama 8B时,TITOK在所有任务上的平均性能比vanilla模型提高了7.11%,同时比 KD高出4.73%,比TransLoRA高出6.24%。这突出表明TITOK不仅限于家族内知识迁移,还可以成功地跨越不同架构。在Llama3 3B到Llama3 8B的情况下,TITOK相对于vanilla、KD和TransLoRA的平均增益分别为3.45%、2.49%和4.14%。这些结果表明TITOK能有效地随模型大小扩展,使其在从轻量级模型迁移到更大型号模型时非常有用,而不会失去效率。最后,当从Llama2 7B迁移到Llama3 8B时,TITOK的平均优势分别为7.43%(对vanilla)、6.28%(对KD)和7.22%(对TransLoRA)。这表明TITOK即使在不同模型版本之间也能稳健地适应,意味着当模型在现实世界管道中升级时具有实际相关性。总的来

说,这些结果提供了强有力的证据,表明TITOK不仅在同一模型家族内的迁移中表现出色,而且在跨模型迁移中也表现出广泛的有效性,从而突显了其在各种模型规模、版本和家族中的鲁棒性。

### 【深度解读】

这段是对表1中更具挑战性的三个场景的分析,这些结果是证明TITOK强大泛化能力的关键。

- **跨架构迁移 (Mistral 7B** → **Llama3 8B)**: 这个场景的成功意义重大。它证明了TITOK 所提取和迁移的"对比盈余"是一种相对"普适"的知识表达,可以被不同架构的模型所 理解和吸收。TITOK在这里全面战胜了所有对手,尤其是大幅领先TransLoRA (6.24%),这可能意味着TransLoRA的方法(特别是其判别器)可能对其模型架构有 更强的依赖性,而TITOK的自监督信号则更加通用。
- **以小教大 (Llama3 3B** → **Llama3 8B)**: 这个结果同样令人振奋。一个3B的小模型所携带的专业知识,能够让一个更强大的8B模型在推理任务(BBH)上获得巨大提升(从0.469提升到0.509)。这验证了之前的一个猜想:专业知识的价值在一定程度上是独立于承载它的模型大小的。这为利用小模型进行低成本的"知识探索",然后将成果"注入"大模型提供了一条可行的路径。
- 版本升级 (Llama2 7B → Llama3 8B): 这是TITOK"商业价值"最直接的体现。实验结果表明,在旧版Llama2上训练的LoRA,可以通过TITOK成功地迁移到新版Llama3上,并且带来的性能提升(BBH上从0.469提升到0.510)甚至超过了其他迁移场景。这一现象背后可能的原因是,Llama3作为Llama2的升级版,其基础能力更强,能够更好地理解和利用从Llama2迁移过来的专业知识,从而产生"1+1>2"的效果。这揭示了一个深刻的道理:专业知识并不会因为基础平台的过时而贬值,相反,当它与一个更强大的新平台结合时,可能会爆发出更大的价值。 这种"非对称的知识价值"现象,是TITOK最重要的发现之一,它为企业如何管理和迭代其AI资产提供了全新的思路。企业不再需要抛弃在旧模型上的投入,而是可以将其视为宝贵的"知识遗产",在新一代模型上实现继承和发扬。

#### 4.3 补充分析

#### 【原文翻译】

消融研究。为了验证我们框架中每个组件的贡献,我们通过选择性地排除第3.3节中的样本过滤和词元选择机制来进行额外实验。我们报告了所有四种迁移设置的平均性能得分,结果呈现在表2中。当不应用样本过滤时(第1和第2行),目标模型的训练数据是从完整的合成数据集中随机抽样的。在这种设置下,仅应用词元选择(第2行)比纯粹的随机基线表现得明显更好。这证明了我们词元选择方法的有效性,并凭经验证实了对比盈余能够可靠地识别和选择信息最丰富的词元。同样,当仅应用样本过滤时(第1和第3行),结果比纯粹的随机抽样有所改善,表明选择高质量的数据至关重要。这一发现进一步强调了为合成数据整合一个有效的过滤机制是必不可少的。在我们的框架中,对比盈余通过可靠地选择具有更丰富和信息量更大信号的样本来扮演这个角色,如经验结果清晰地展示。最后,结果表明,结合两个阶段(第4行)取得了最佳性能,证实了它们在过滤高质量示例和选择信息词元以实现有效知识迁移方面的互补作用。

### 【深度解读】

消融研究(Ablation Study)是科学实验中非常重要的一环,它的目的是通过"做减法"来理解一个复杂系统中每个部分的功能。这就像拆解一台精密的手表,拿掉一个齿轮看看手表会怎么样,从而理解这个齿轮的作用。

### • 表2的解读逻辑:

- 第1行 (无过滤,无选择): 这是最差的情况,相当于学生拿到整本练习册,不分好坏、不划重点地随机看。
- 第2行 (无过滤,有选择): 相当于学生拿到整本练习册,但老师在里面划了重点词句。结果比第1行好,证明了"划重点"(词元选择)这个操作本身是有效的。
- 第3行 (有过滤,无选择): 相当于学生只拿到了练习册的"精华章节",但没有划重点。结果也比第1行好,证明了"筛选章节"(样本过滤)这个操作也是有效的。
- 。 **第4行 (有过滤,有选择):** 这是完整的TITOK方法,相当于学生拿到"精华章节", 并且里面还划了"核心考点"。结果是最好的。
- 结论: 这个实验清晰地证明了TITOK的成功并非偶然,也不是某个单一组件的功劳。 样本过滤和词元选择这两个阶段,就像是"战略筛选"和"战术打击",各有分工,且相 辅相成,共同构成了TITOK高效的知识筛选体系。缺少任何一环,效果都会打折扣。 这证明了TITOK框架设计的合理性和完整性。

# 【表格2:消融研究】

"样本过滤"使用平均词元级对比盈余分数来移除信息量不足的样本,而"词元选择"通过选择性地保留每个样本中最具信息量的词元来进一步提炼保留的数据。在没有样本过滤的情况下,数据是随机抽样的。结果报告为BBH和MMLU的平均准确率(Acc.),以及LaMP任务的ROUGE-1(R-1)和ROUGE-L(R-L)。

<b>✓</b>	1	0.483	0.501	0.142 / 0.125	0.464 / 0.403
<b>✓</b>	Х	0.470	0.500	0.139 / 0.122	0.460 / 0.397
×	✓	0.463	0.496	0.137 / 0.121	0.460 / 0.397
X	X	0.458	0.485	0.133 / 0.117	0.456 / 0.393
样本过滤	词元选择	BBH Acc.	MMLU Acc.	新闻标题 R-1/R-L	学术标题 R-1/R-L

#### 【原文翻译】

**查询生成模型的影响。**我们现在来研究查询生成模型的选择如何影响性能。实际上,合成查询可以由源模型或目标模型生成,后者是TransLoRA管道(Wang et al., 2024)最初采用的方法。图2比较了这两种选择在不同迁移设置下的情况,报告的分数是所有BBH任务的

平均值。值得注意的是,我们观察到使用源专家模型(源骨干+LORA)进行查询生成通常会产生更强的性能。一个可能的解释是,当合成查询和相应的标签都由同一个模型生成时,它们更接近其训练分布。查询和标签之间的这种对齐可能使监督更加连贯和准确,从而促进更有效的迁移。这些结果表明,保持查询和标签之间的分布对齐是提高知识迁移效果的关键因素,从而为我们选择使用源模型作为合成查询和标签生成器的设计提供了经验依据。

### 【深度解读】

这个实验验证了作者在3.1节中做出的一个关键设计选择的正确性。

#### • 两种策略的对比:

- 1. TransLoRA策略: 让"学生"(目标模型)提问,"老师"(源模型)回答。
- 2. **TITOK策略:**让"老师"自己出题,自己给答案。
- **图2的结果**: 实验数据清晰地显示,在所有迁移场景下,TITOK的策略(蓝色柱子) 都优于或等于TransLoRA的策略(橙色柱子)。
- 深层原因分析: 作者给出的解释非常深刻——"分布对齐"。当老师自己出题和作答时,题目(查询)的风格、难度、用词习惯和知识假设,与答案(标签)是完全匹配的,它们源自同一个"知识体系"。这种内在的一致性使得学习材料的质量更高,学生学起来也更顺畅。而如果让学生提问,他可能会问出一些超出老师专业范畴或者与老师思路不符的问题,导致"鸡同鸭讲",教学效果自然会打折扣。
- 设计启示: 这个发现为未来的研究提供了宝贵的经验。在设计类似的知识迁移系统时,优先考虑由同一个专家模型生成完整的"教案"(问题+答案),是保证教学质量的关键。

#### 【图2与图3描述】

**图2:查询来源的影响。** 该柱状图比较了在四种迁移设置下,使用"目标模型(vanilla)"生成查询和使用"源模型+LORA"生成查询对BBH任务平均准确率的影响。图中显示,蓝色柱子(源模型生成)在所有设置中都高于或等于橙色柱子(目标模型生成),表明使用源专家模型来合成查询通常能带来更好的性能。

**图3:不同k%下的代表性性能趋势。** 该图包含三条折线图,展示了在不同任务和迁移设置下,模型的性能如何随着"词元选择比例k%"的变化而变化。

- (a) BBH (Mistral 7B → Mistral 7B): 在同模型迁移的推理任务中,性能呈现出一条"抛物线"形状。当k%在40-70%之间时,效果最好。太低(如10%)会导致学习不足,太高(100%,即不选择)则失去了"划重点"的意义。
- **(b) BBH (Llama2 7B** → **Llama3 8B)**: 在"以弱教强"的推理任务中,趋势完全不同。 k%越低,性能越好。这表明当老师较弱时,只教最精华、最有把握的知识点(最低的k%)是最好的策略,可以最大限度地避免"噪声"干扰。

• (c) 新闻标题生成 (Llama2 7B → Llama3 8B): 在"以弱教强"的生成任务中,趋势再次 反转。k%越高,性能越好。这可能是因为对于风格模仿这类任务,即使是老师的"非重点"词汇,也能提供有价值的文体线索,所以"多多益善"。

### 【原文翻译】

词元选择比例的影响。 我们现在探讨我们框架中词元选择比例(k%)的影响。图3展示了三条代表性曲线:一个任务(BBH)在两种迁移设置下,以及另一个任务(新闻标题生成)在相同的迁移设置下进行比较。对于源和目标具有相同骨干的BBH(Mistral 7B→Mistral 7B),40-70%的中等词元选择比例产生最佳结果。非常低的比例(10-20%)导致欠拟合,而100%的比例是无效的,因为它根本没有应用过滤。有趣的是,在同一BBH任务的弱到强迁移(Llama2 7B → Llama3 8B)中,性能随着k%的降低而提高。这表明,仅选择具有最高盈余损失的词元能有效过滤掉较弱模型不确定的大部分噪声区域。通过丢弃较弱模型缺乏信心的大部分词元,TITOK可以显著减少负迁移。然而,在相同的弱到强迁移设置下,这种趋势对于新闻标题生成任务是相反的。对于这个任务,通常k%越大越好,可能是因为即使是一个较弱的模型也能提供有价值的词汇和文体线索,这对于个性化是有用的。因此,与对噪声监督敏感的推理任务(如BBH)不同,像新闻标题生成这样的个性化任务能更广泛地从源模型的输出中受益,即使源模型相对比目标模型弱。

### 【深度解读】

图3的分析揭示了一个非常深刻和细致的发现:最优的"划重点"策略,取决于"老师"、"学生"和"科目"三者之间的关系。

- 场景一:强师教强生,学推理(图3a)。 当老师和学生水平相当(Mistral 7B → Mistral 7B),学习的是逻辑严谨的推理任务时,一个"适度"的划重点范围(40-70%)是最好的。划得太少,知识点不够;划得太多,重点不突出。
- 场景二:弱师教强生,学推理(图3b)。 当老师比学生弱(Llama2 7B → Llama3 8B),学习的仍然是推理任务时,策略应该变为"极度聚焦"。老师应该只教自己最懂、最有把握的核心知识(k%越低越好)。因为对于更聪明的学生来说,老师的那些"半懂不懂"的知识(低盈余分数词元)不仅没有帮助,反而是一种"噪声"和干扰,会妨碍学生自己的思考。这是一种"宁缺毋滥"的教学智慧。
- 场景三: 弱师教强生,学写作(图3c)。 当老师比学生弱,但教的是"文科"——风格模仿类的写作任务时,策略又变了。这时候,老师说的每一句话,即使不是核心知识点,也可能携带着他独特的"口音"和"文风"(词汇和文体线索)。对于想学习这种风格的学生来说,这些信息都是有价值的。因此,在这种情况下,k%越大越好,学生应该"兼收并蓄",尽可能多地接触老师的语言材料。 这个分析充分展示了TITOK框架的灵活性和其背后深刻的洞察力。它告诉我们,知识迁移不是一个一成不变的公式,而是一个需要根据具体情况调整策略的动态过程。TITOK不仅提供了一个强大的工具,还通过实验揭示了如何更智慧地使用这个工具。

# 【原文翻译】

通过外部数据源迁移的有效性。 我们进一步研究在不倾向于使用合成数据,因此使用外部数据作为替代方案的情况下,TITOK是否能有效迁移知识。为此,我们在LaMP任务上评估了三种替代设置,用于Mistral 7B → Mistral 7B的迁移设置: (1) 使用来自随机选择的不同用户的数据, (2) 混合来自多个用户的数据,以及 (3) 跨任务迁移,其中学术标题生成的数据用于训练新闻标题生成,反之亦然(即,分布外场景)。结果呈现在表3中。值得注意的是,TITOK在这些异构的外部设置中始终优于所有基线。这些发现表明,TITOK不仅限于合成数据场景,也能在外部或用户提供的数据条件下有效适应。这突显了TITOK在多样化的现实世界应用场景中进行实际部署的灵活性,并进一步强调了其跨不同数据条件的适应性,因为它即使在使用外部数据作为合成数据的替代品时也有效。

### 【深度解读】

这个实验是TITOK实用性的"终极证明"。它回答了一个关键问题:如果我不想用AI生成的合成数据,而是想用我自己手头的、真实的数据(即使这些数据不那么"完美"),TITOK还能工作吗?

### • 三种"真实世界"的数据场景:

- 1. **用别人的数据 (other-user):** 我想学习用户A的写作风格,但我手头没有用户A的数据,只有用户B的。我能用用户B的数据作为"桥梁"来迁移知识吗?
- 2. **用混合数据 (mixed-user):** 我手头的数据是多个用户混在一起的,比较杂乱。
- 3. **用不相关任务的数据 (cross-task):** 我想学习写新闻标题,但我只有学术论文的数据。
- **表3的结果**: 结果令人惊讶。在所有这些"数据不纯"、"任务错配"的苛刻条件下, TITOK的表现依然是最好的(所有粗体数字)。
- 深刻的启示: 这个实验证明了TITOK的核心机制——"对比盈余"信号的普适性。无论 底层的数据是合成的、真实的、来自他人的还是来自其他任务的,只要存在一个"新手"(基础模型)和一个"专家"(带LoRA的模型),TITOK就能通过对比它们的差 异,找出其中蕴含的、与任务相关的"知识增量"。这个机制是独立于数据来源的。
- 巨大的实用价值: 这一特性使得TITOK的应用场景被极大地拓宽了。在许多商业应用中,企业可能拥有大量不同来源的、未标记的文本数据。TITOK提供了一种可能性,即利用这些现成的"外部数据"作为载体,来完成LoRA知识的迁移和更新,而无需为每个任务都去生成或标注新的数据集。这使得TITOK从一个"聪明的算法"变成了一个"灵活实用的工业级解决方案"。

### 【表格3:使用外部数据进行迁移】

在三种数据设置下,新闻标题和学术标题生成任务的ROUGE-1 (R-1) 和 ROUGE-L (R-L) 分数,用于Mistral 7B $\rightarrow$ Mistral 7B迁移:(1)其他用户,(2)混合用户,(3)跨任务。最优分数以**粗体**显示。

数据设置	方法	新闻标题 R-1/R-L	学术标题 R-1/R-L
其他用户	Vanilla	0.117 / 0.101	0.381 / 0.311
	KD	0.127 / 0.113	0.405 / 0.333
	TiTok (ours)	0.151 / 0.136	0.481 / 0.425
混合用户	Vanilla	0.117 / 0.101	0.381 / 0.311
	MinED	0.124 / 0.110	0.409 / 0.337
	TiTok (ours)	0.151 / 0.137	0.480 / 0.422
	Vanilla	0.117 / 0.101	0.381 / 0.311
	MinED	0.118 / 0.106	0.403 / 0.331
	TiTok (ours)	0.133 / 0.120	0.450 / 0.383

第五部分:结论、局限性与附录

# 5 结论

# 【原文翻译】

在本文中,我们提出了TITOK,一个高效的框架,通过仅在一组经过选择性挑选的高度信息化的词元上进行训练,将LoRA知识从源模型迁移到目标模型。TITOK的核心是利用词元级信号从源适配器中提炼任务相关信息,而不是依赖整个词元序列。通过将监督集中在适配器贡献最大的区域,TITOK实现了更强、更有针对性的知识迁移。凭借这种简单而有效的设计,TITOK在各种任务中表现出鲁棒性,并持续超越现有基线,使其成为高效知识迁移的实用解决方案。

### 【深度解读】

结论部分是对整篇论文工作的高度概括和升华。

- **再次强调核心思想:** "仅在一组经过选择性挑选的高度信息化的词元上进行训练"。这是TITOK的灵魂,也是它区别于其他方法的根本所在。
- 机制的本质:"利用词元级信号…提炼任务相关信息"。这指出了"对比盈余"的本质作用——它是一个信息提纯器。
- **效果的根源:**"将监督集中在适配器贡献最大的区域"。这解释了为什么TITOK效果好——因为它把学习资源用在了"刀刃"上。

• 最终评价:"简单而有效"、"鲁棒"、"实用"。这三个词精准地概括了TITOK的优点。在 计算机科学领域,一个好的解决方案往往不是最复杂的,而是最能以简单的方式解决 复杂问题的。TITOK正是这种"大道至简"设计哲学的典范。它为如何高效、低成本地 复用和传承AI知识资产,提供了一个非常漂亮且实用的答案。

### 【原文翻译】

### 局限性与未来方向

虽然TITOK已显示出鲁棒性和相较于现有方法的明显优势,但仍有机会对该框架进行改进和扩展。首先,尽管TITOK旨在最大限度地减少对原始数据的依赖,但它仍然需要少量的种子示例通过提示来生成合成数据。然而,这种依赖是适度的,因为我们避免使用完整的数据集。此外,我们的实验证实,外部数据也可以作为一种有效的替代方案,增强了该方法的灵活性。关于词元选择,TITOK目前应用一个固定的阈值来确定保留哪些词元。虽然这种简单的设计在各种任务中是有效和稳定的,但未来的工作可以探索更具适应性或数据驱动的阈值策略,以在不损害鲁棒性的前提下进一步提高效率。

### 【深度解读】

任何一项严谨的科学研究都会坦诚地指出自身的局限性,并为后续的研究者指明方向。这不仅是科学诚信的体现,也是推动学科不断进步的方式。

- **局限一:对"种子"的依赖。** TITOK虽然不需要完整的训练数据,但它的合成数据生成过程仍然需要5个左右的"种子样本"来启动。这是一个轻微的依赖。作者也提出了解决方案:实验已经证明外部数据也可以用,这在很大程度上缓解了这个问题。未来的研究可以探索完全不需要任何示例的"零样本"启动方法。
- **局限二:"一刀切"的阈值。** 目前,"划重点"的比例k%是人为设定的一个固定值(比如70%)。虽然实验表明这个简单的策略效果不错,但它可能不是最优的。一个更智能的系统应该能够根据任务的类型、模型的强弱等因素,**自动地、动态地**调整k%的大小。比如,系统自动识别出这是一个"弱师教强生,学推理"的场景,然后自动将k%调整到一个较低的值。实现这种"自适应阈值"将是未来一个非常有价值的研究方向,它将使TITOK变得更加智能和高效。

#### 【原文翻译】

#### 道德声明

我们完全依照既定的道德标准进行本研究。对于我们的实验,我们完全依赖于公开可用的数据集,如LaMP、MMLU和BBH,并严格按照其用于学术研究的既定目的使用它们。对于涉及用户数据的LaMP任务,我们的TITOK框架通过最小化数据依赖性来符合道德考量。它不存储或暴露原始用户数据,仅更新一小组特定于任务和用户的参数。这有助于在实现高效知识迁移的同时保护隐私。

#### 可复现性声明

我们在第4节中提供了我们实现的全面描述,包括管道配置、超参数、模型、数据集和评估指标。我们的实现和实验的源代码将在发表后在一个代码库中公开提供。

# 参考文献

[此处省略参考文献列表的逐条翻译,保留格式]

#### 附录 A 分词器对齐算法

### 【原文翻译】

当源模型和目标模型使用不同的分词器时,直接的词元级迁移是不可能的。为了解决这个问题,我们实现了一个双指针校准程序。该算法维护两个指针,一个用于源词元,一个用于目标词元。在每一步,源指针前进一个词元,累积一个解码后的片段,而目标指针则逐步扩展自己的片段,直到规范化的文本匹配为止。一旦找到匹配,相应的跨度就被记录为一个对齐,目标指针跳到该位置。对齐后,我们应用掩码规则将源二进制掩码分数(指示一个词元应该被保留还是丢弃的值)传播到目标词元。具体来说:

1. 一对一: 二进制掩码分数被直接复制。

2. 一对多: 分数被复制到所有对齐的目标词元。

3. 多对一: 多个源词元的平均分数被分配给目标词元。

4. **多对多:** 源词元的平均分数被分配给相应的对齐目标词元。

### 【深度解读】

附录A是对正文3.4节中提到的分词器对齐算法的更详细的技术说明。它揭示了"翻译器"工作的具体细节。核心思想是:**不比较词元本身,而是比较词元解码后的实际文本内容。** 

• **双指针校准:** 想象有两个人在分别阅读两本不同出版社的《哈姆雷特》。虽然排版和断句可能不同,但内容是一样的。双指针算法就像是这两个人同时朗读,当他们读到同一句话的结尾时,就完成了一次"对齐"。

- 掩码传播规则: 对齐之后,如何传递"重点标记"?
  - 。如果老师的教材里"to be or not to be"这句被划了重点,而学生的教材里这句话 恰好也是一个整体,那就直接把重点标记复制过去(一对一)。
  - 。 如果老师的教材里"to be or not to be"是一个整体,而学生的教材里是"to be"和"or not to be"两部分,那就把重点标记复制给这两个部分(一对多)。
  - 。 反之,如果老师的教材是两部分,学生的是一部分,那就把两个部分的重点标记平均一下,赋给学生的那一部分(多对一)。 这个简单而有效的过程,确保了知识的重点能够在不同"教材版本"之间准确无误地传递。

### 【图4描述】

图4:TITOK的分词器对齐算法概览。 该图展示了当源模型和目标模型使用不同分词器时,算法如何处理。源模型分配的二进制掩码分数在对齐的文本片段内被平均,并传播到目标词元,产生小数分数,这些分数指导目标模型LoRA适配器训练的top-k%词元选择。例如,图中源模型将句子 "But seek first his kingdom and his righteousness" 分词,并给出了一个二进制的重点掩码``。目标模型的分词结果不同。对齐算法发现:

- 源的 his (掩码1) 对应目标的 his (掩码1), 为一对一, 分数直接复制。
- 源的 first (掩码1) 对应目标的 first (掩码1), 为一对一。
- 源的 kingdom (掩码1) 和 and (掩码1) 对应目标的 kingdom and (掩码1), 为多对一,
  分数平均 (1+1)/2=1。
- 源的 righteousness (掩码1) 对应目标的 right 和 eousness (此处图示有误,应为类似的多部分),为**一对多**,分数复制。
- 一个**多对多**的例子显示,源的两个词元(总分2/3=0.667)对应目标的两个词元,每个都获得0.667的分数。 最终,目标模型获得一个带小数的"重要性分数"序列,然后根据这些分数选择top-k%的词元进行训练。

#### 附录 B 数据集详情

#### 【原文翻译】

B.1 BIG-BENCH HARD (BBH) Big-Bench Hard (BBH) (Suzgun et al., 2022) 是一个为评估模型在挑战性推理问题上的性能而设计的严格基准,包括多步逻辑推理、符号操纵和常识推理。任务格式为多项选择题。由于BBH最初是一个只有测试集的基准,我们将其90%的数据用于训练源专家模型,并保留10%用于评估。27个BBH任务分类见表4。

【表格4:27个Big-Bench Hard (BBH) 任务的分类】

逻 boolean expressions, causal judgement, date understanding,

辑 disambiguation\_qa, dyck\_languages, formal\_fallacies,

logical\_deduction\_three\_objects, logical\_deduction\_five\_objects,

logical\_deduction\_seven\_objects, temporal\_sequences,

tracking\_shuffled\_objects\_three\_objects, tracking\_shuffled\_objects\_five\_objects, tracking\_shuffled\_objects\_seven\_objects, web\_of\_lies

语 hyperbaton, ruin\_names, salient\_translation\_error\_detection, snarks, word\_sorting

学

玾

geometric\_shapes, multistep\_arithmetic\_two, object\_counting,

学/ reasoning\_about\_colored\_objects

符号

数

movie recommendation, navigate, penguins in a table, sports understanding

用/

知 识

**B.2 大规模多任务语言理解 (MMLU)** 大规模多任务语言理解 (MMLU) (Hendrycks et al., 2021) 是一个用于评估模型在广泛知识密集型任务上性能的综合基准。该基准由多项选择 题组成,与BBH类似,我们也对原始的仅测试集数据应用了90%/10%的划分。所有57个子

任务的分类见表5。

【表格5:57个MMLU任务的分类】

类别 任务

abstract\_algebra, anatomy, astronomy, college\_biology, college\_chemistry, college\_computer\_science, college\_mathematics, college\_medicine, college\_physics, computer\_security, conceptual\_physics, electrical\_engineering, elementary\_mathematics, formal\_logic, high\_school\_biology, high\_school\_chemistry, high\_school\_computer\_science, high\_school\_mathematics, high\_school\_physics, high\_school\_statistics, machine\_learning, medical\_genetics, nutrition, professional\_medicine, professional\_psychology, virology, clinical\_knowledge, human\_aging, human\_sexuality

人文学 business\_ethics, formal\_fallacies, jurisprudence, logical\_fallacies, philosophy, prehistory, world\_religions, moral\_disputes, moral\_scenarios,

professional\_law, professional\_accounting, high\_school\_world\_history, high\_school\_us\_history, high\_school\_european\_history, global\_facts, security studies

社会科学

econometrics, high\_school\_macroeconomics, high\_school\_microeconomics, management, marketing, public\_relations, sociology, us\_foreign\_policy,

(10)

high\_school\_geography, high\_school\_government\_and\_politics

其他 miscellaneous, global\_facts (2)

**B.3 LAMP 任务** 我们利用LaMP基准来评估TITOK在个性化设置中是否也有效。在LaMP任务中,我们专注于两个适合、可访问且可在我们设置中评估的文本生成任务:

- **新闻标题生成 (LaMP 4)。** 给定一个由先前撰写的标题组成的作者简介,模型被要求为一个新的新闻文章生成一个标题。该任务评估模型是否能调整其输出以反映作者在新闻写作中的特有风格。
- **学术标题生成 (LaMP 5)**。 使用一个由学术出版物先前标题构建的作者简介,模型为一个新的给定摘要生成一个标题。该任务评估模型捕捉和再现学术写作独特惯例的能力。

# 附录 C 基线详情

### 【原文翻译】

- **C.1 VANILLA** Vanilla基线对应于标准的、未经任何额外训练或从源模型进行知识迁移的目标基础模型。此设置反映了目标模型在原始初始化状态下的性能,并作为评估迁移方法的下限。直观地,实现超越此基线的性能提供了明确的证据,表明目标模型已成功获取并内化了从源模型迁移的知识。
- C.2 知识蒸馏 (KD) (+MINED) 知识蒸馏 (KD) 基线训练学生模型模仿教师模型的输出分布。具体来说,损失是交叉熵目标和教师与学生分布之间KL散度的加权和。在我们的设置中,KD实验使用TransLoRA过滤后的合成数据集进行。对于源和目标模型使用不匹配分词器的情况,我们应用最小编辑距离 (MinED) 对齐方法。MinED通过最小化字符级编辑来匹配不同词汇表中的词元。
- **C.3 TRANSLORA** TransLoRA基线通过生成合成数据来迁移LORA知识。在这种方法中,vanilla目标模型合成查询,源模型及其LORA适配器提供相应的标签。随后,一个以源模型为基础单独训练的判别器被应用于过滤用于微调目标模型LoRA适配器的合成数据。为了一致性和公平比较,我们在应用TransLoRA基线时采用了与我们实验相同的超参数设置。

### 附录 E 与额外基线的比较

# 【原文翻译】

在这里,我们引入一个利用少样本监督的额外基线。实际上,我们的合成数据生成过程是由一小组提示作为种子的,这可以被视为少样本数据。为求周全,我们因此建立一个仅在这五个种子提示上训练的基线,称为KD (5-shot),以便我们对TITOK的评估也考虑到最小的少样本监督。这个额外基线的结果呈现在表6中。在所有迁移设置中,仅在5-shot样本上进行KD只比vanilla模型提供了边际改进... 相比之下,TITOK在此基线上实现了显著的改进... 战胜仅有5个少样本的KD表明,TITOK的成功不仅仅是少样本效应的结果。相反,五个种子提示仅作为起点,我们的框架将其扩展为更丰富、更有效的合成训练信号。这一发现证实了TITOK的真正优势在于它如何有效利用有限的监督,而不在于少样本数据本身。

### 【表格6:与使用五个种子提示的少样本KD基线的比较】

迁移设置	方法	BBH Acc.	MMLU Acc.	新闻标题 R- 1/R-L	学术标题 R- 1/R-L
Mistral 7B → Mistral 7B	Vanilla	0.397	0.557	0.117 / 0.101	0.381 / 0.311
	KD (5- shot)	0.402	0.558	0.104 / 0.118	0.383 / 0.312
	TiTok (ours)	0.432	0.563	0.160 / 0.142	0.473 / 0.414
	***				

# 附录 F-H:数据生成与提示模板

#### 【原文翻译摘要】

- **F节** 详细说明了合成数据池的大小。首先生成源训练集大小两倍的样本,然后通过对比盈余分数过滤,保留与源训练集大小相同的M个样本(BBH M=250, MMLU M=90, LaMP M=200)。
- **G节** 列出了在生成合成数据时,由于任务特性(如格式高度受限)而**没有**应用 ROUGE-L过滤的任务列表(见表7)。
- H节 提供了用于生成合成数据的详细提示(Prompt)模板。为BBH的每个子任务都设计了特定的指令和格式要求(以boolean\_expressions为例,见表8)。为MMLU设计了通用模板(表9)和针对历史类任务的专门模板(表10)。为LaMP的新闻标题(表11)和学术标题(表12)任务也提供了相应的模板。

### 【表格7:未应用ROUGE-L过滤的任务】

类别 任务

**BBH** bbh.boolean.expressions, bbh\_date\_understanding, bbh\_disambiguation\_qa, bbh\_geometric\_shapes,...

**MMLU** high\_school\_world\_history, high\_school\_us\_history, high\_school\_european\_history

【表格8:BBH (boolean\_expressions) 任务的合成查询生成提示】

System 你是一个boolean\_expressions任务的专家生成器。生成以 is结尾的布尔表达式求值任务 - 仅限单行。 关键格式要求: - 遵循示例中显示的确切格式结构。 - 绝对关键: 每个响应只生成一行,一个任务。 - 禁止多行,禁止换行符(\n),禁止多个表达式。 - 必须以 is结尾。 - 无效输出示例: True or False is\n False and True is - 有效输出示例: True or False and (True or False) is 生成多样的内容但保持完全相同的格式结构。只输出任务输入,不要输出解决方案。

User 按照这些确切的格式示例生成新的boolean\_expressions任务: 示例1: [布尔表达式以'is'结尾] ... 示例5: [布尔表达式以'is'结尾] **关键:** 只生成一行,完全像上面的示例一样。 按照确切的格式生成一个新任务:

[其余附录中的提示模板表格内容繁多,此处省略以保持报告核心部分的清晰度,但在实际翻译中会完整呈现。]