



中山大學  
SUN YAT-SEN UNIVERSITY

## 期中实验报告

Text-to-LoRA: Instant Transformer Adaption

姓名 \_\_\_\_\_ 常毅成-孙雨-闫璞鑫

学号 \_\_\_\_\_ 22354010-22354118-22354168

学院 \_\_\_\_\_ 智能工程学院

专业 \_\_\_\_\_ 智能科学与技术

## 目录

<b>1 引言</b>	<b>2</b>
<b>2 Text-to-LoRA 论文深度解析</b>	<b>2</b>
2.1 核心思想：从“数据微调”到“指令生成” . . . . .	2
2.2 关键创新点剖析 . . . . .	3
2.3 模型架构剖析 . . . . .	3
2.4 T2L 训练过程概述 . . . . .	3
<b>3 论文复现部分</b>	<b>4</b>
3.1 环境配置 . . . . .	4
3.2 预训练模型下载 . . . . .	4
3.3 演示运行 . . . . .	5
3.4 结果分析与评估 . . . . .	5
3.4.1 运行评估脚本 . . . . .	5
3.4.2 评估结果 . . . . .	6
3.5 复现总结 . . . . .	6
<b>4 挑战、反思与展望</b>	<b>6</b>
4.1 复现过程中的主要挑战 . . . . .	6
4.2 收获与反思 . . . . .	7
4.3 未来展望与改进方向 . . . . .	7
<b>5 总结</b>	<b>7</b>

# 期中实验报告

常毅成-孙雨-闫璞鑫 22354010-22354118-22354168

**摘要:** 本报告对 Sakana AI 的开创性工作《Text-to-LoRA: Instant Transformer Adaption》进行了全面而深入的复现与剖析。该论文提出了一种名为 Text-to-LoRA (T2L) 的超网络模型，它颠覆了传统的大型语言模型 (LLM) 适配流程，能够仅依据一段自然语言任务描述，通过一次前向传播即时生成定制化的低秩适配器 (LoRA)。这一方法不仅极大地降低了模型专业化的技术与资源门槛，也为实现更加动态和智能的人机交互提供了新的可能。

本报告首先深度解析了 T2L 的核心设计哲学——从“数据驱动微调”到“指令驱动生成”的范式转变，并详述了其模型架构、两种关键训练策略（重构与监督微调）及其对模型泛化能力产生的决定性影响。报告的核心部分，是我们详尽的复现之旅：一份从环境配置、数据准备、关键的“任务描述工程”，到模型训练、实时监控、结果验证与定性分析的完整实验日志。我们不仅定量复现了原论文的关键零样本评测结果，还通过定性实验直观感受了 T2L 的“可操纵性”。

最后，报告系统总结了复现过程中遇到的技术挑战，如显存瓶颈、初始化方法的敏感性等，并分享了我们从中获得的宝贵经验与深刻反思。我们相信，T2L 不仅是一项卓越的技术成就，更指明了未来 AI 系统发展的一个重要方向。

**关键词:** Text-to-LoRA; 低秩适配器 (LoRA) ; 可操纵性; 显存瓶颈

## 1 引言

大型语言模型 (LLM) 的出现无疑是人工智能发展史上的一个里程碑，但其巨大的成功也伴随着一个深刻的挑战：如何跨越通用能力与特定应用之间的“最后一公里”。传统的适配方法，如完全微调，成本高昂到令多数开发者望而却步；而即使是参数高效微调 (PEFT) 技术，也依然遵循着“为每个任务收集数据、配置训练、等待收敛”的固定流程。这形成了一种“敏捷性瓶颈”，严重限制了 LLM 在需要快速响应和高度个性化场景中的应用潜力。我们正处在一个需要更高效、更民主化的模型适配方案的时代。

正是在这样的背景下，Sakana AI 提出的 Text-to-LoRA (T2L) 应运而生。它所代表的不仅仅是一种技术优化，更是一种根本性的范式转移：从以数据为中心的微调 (Data-centric Fine-tuning) 转向以指令为中心的生成 (Instruction-centric Generation)。T2L 试图模仿人类学习的核心机制——我们通过听取指令和描述来快速掌握新技能，而不是每次都依赖海量案例的重复训练。

T2L 通过训练一个超网络 (Hypernetwork)，使其学会将任务的自然语言描述映射到能够解决该任务的 LoRA 适配器参数空间。这使得模型适配过程从数小时乃至数天的训练，压缩到了一次毫秒级的前向传播。

本报告旨在对 T2L 的理论与实践进行一次彻底的探索，通过详尽的复现工作，不仅验证其结论的可靠性，更深入理解其背后的设计哲学。本报告结构如下：第二节深度解析

T2L 的理论框架，包括其核心思想、关键创新、架构细节以及两种训练策略的本质区别；第三节以实验日志的形式，详尽记录我们从零开始的复现之旅，涵盖环境准备、数据工程、模型训练、结果验证和分析的全过程；第四节坦诚分享复现过程中遇到的技术难题，以及我们从中获得的超越代码本身的认知与反思，并对该技术的未来发展进行展望；第五节对整个复现工作进行总结。

## 2 Text-to-LoRA 论文深度解析

### 2.1 核心思想：从“数据微调”到“指令生成”

T2L 的根本突破在于它改变了传统模型适配的范式。传统方法依赖于通过大规模数据集和梯度下降进行参数优化，这一过程不仅耗时长、计算资源需求高，还对超参数的选择极为敏感，导致适配过程复杂且不稳定。相比之下，T2L 提出了一种全新的“适配”概念，将其转变为一个基于自然语言指令的生成过程。T2L 作为一个超网络，通过单次正向传播直接根据任务描述生成任务特定的 LoRA 适配器，从而使大型语言模型 (LLM) 从一个被动接受训练的实体，进化成能够根据用户需求主动调整能力的智能系统。这种从“数据驱动”到“指令驱动”的转变，极大降低了适配的门槛，使得非专业用户也能通过简单描述实现模型定制。

此外，T2L 的设计还体现了灵活性与即时性的结合。传统微调需要预先准备特定任务的数据集并进行多次迭代，而 T2L 仅依赖任务的高级描述即可工作，这在时间紧迫或数据稀缺的场景中尤为适用。例如，在实时应用或新兴任务中，用户无需收集数据，只需提供简要说明，T2L 即可生成适配器，显著提升了模型的实用性与普适性。

## 2.2 关键创新点剖析

- 即时适配：**T2L 突破了传统微调的局限，通过超网络架构实现基于自然语言任务描述的即时适配。其核心在于，仅需一次正向传播，T2L 就能生成适用于目标任务的 LoRA 矩阵（如  $A$  和  $B$  矩阵），无需进行耗时的梯度下降训练。论文实验表明，在 Super Natural Instructions (SNI) 数据集上训练的 T2L，能够在几秒内为新任务生成适配器，相比传统微调的数小时甚至数天，效率提升显著。此外，这种方法还避免了超参数调整的复杂性，降低了用户的使用难度，使其适用于资源有限的场景，如边缘设备或小型研究团队。
- 零样本泛化：**T2L 的零样本泛化能力是其一大亮点。训练完成后，T2L 能够处理未见过任务的描述，并生成有效的 LoRA 适配器。例如，论文中的 Table 2 显示，T2L 在未见过的 10 个基准任务（如 GSM8K 和 HumanEval）上取得了平均 66.3% 的性能，接近任务特定 LoRA 的水平。这一能力得益于超网络在多样化任务分布（如 479 个 SNI 任务）上的训练，使其学习到了任务适配的通用模式。这种泛化性为跨领域应用提供了可能，例如从数学推理扩展到代码生成，而无需额外训练。
- 高效适配器压缩：**T2L 通过超网络将数百个预训练 LoRA 实例压缩为一个紧凑的模型，显著减少了存储和计算开销。论文提出了三种架构变体：L (55M 参数)、M (34M 参数) 和 S (5M 参数)，分别代表了高表达能力、中等平衡和高效压缩的取向。实验结果表明，即使是最小的 S 架构，在压缩比高达 8 倍的情况下，性能仍能维持在基准 LoRA 的 65% 左右（见 Figure 1，右上）。这种压缩能力使其适用于多任务场景，例如云端服务中同时支持多个用户定制需求，同时降低了部署成本。
- 语言驱动的可操纵性：**T2L 的语言驱动特性为其提供了高度的可控性，论文中的 Figure 4 提供了直观的例证。在 GSM8K 数学问题上，当任务描述聚焦“数学推理”时，生成的 LoRA 引导模型输出结构化的解题步骤；而当描述转向“编程与算法”时，模型会先分析变量关系再计算。这种 steerability 得益于任务描述嵌入 (task embedding) 与超网络的紧密耦合，允许用户通过调整描述语言微调模型

行为。例如，描述中加入“简洁回答”可能减少冗长推理，而“详细解释”则会激发更全面的输出。这种灵活性在教育、创意写作和个性化助手等领域具有巨大潜力。

## 2.3 模型架构剖析

T2L 的架构设计体现了高效性与模块化的结合，其核心在于超网络如何根据输入生成任务特定的 LoRA 适配器。超网络的输入由任务描述嵌入 ( $f(z^i)$ )、模块嵌入 ( $E[m]$ ) 和层嵌入 ( $E[l]$ ) 拼接而成，数学表达为：

$$\phi_{m,l}^i = \text{concat}[f(z^i), E[m], E[l]]$$

其中， $f(z^i)$  是通过预训练嵌入模型（如 gte-large-en-v1.5）提取的任务描述表示， $E[m]$  和  $E[l]$  分别是针对不同模块（如 q\_proj、v\_proj）和层数的可学习嵌入。这种设计允许超网络为 LLM 的不同部分生成专门的参数，而非使用单一函数覆盖所有权重，从而提升了模型的表达能力。

L、M、S 三种变体则体现了不同归纳偏置与参数规模的权衡：

- L (Large)：**L 架构同时输出  $A$  和  $B$  矩阵，其输出头参数量为  $d_{\text{out}} \times 2 \times r \times d$ ，表达能力最强，但参数量大 (55M)，可能导致过拟合，尤其在任务分布有限时。实验显示，其在重建训练中表现最佳，但泛化性稍逊。
- S (Small)：**S 架构一次只生成低秩矩阵的一个秩，输出头参数量降至  $d_{\text{emb}} \times d$ ，参数效率最高 (5M)，但因容量限制，性能上限较低。适合资源受限环境，但需谨慎选择任务复杂度。
- M (Medium)：**M 架构共享  $A$  和  $B$  的输出层，参数量为  $d_{\text{out}} \times r \times d$  (34M)，在性能与效率间取得平衡，是实验中最稳定的选择，尤其在 SFT 训练中表现优异。

这些变体的设计还支持批量生成，通过一次性处理多个嵌入，进一步提升计算效率。

## 2.4 T2L 训练过程概述

T2L 的训练过程采用两种互补策略，各有针对性：

- LoRA 重建训练：**此方法旨在让 T2L 学习重构预训练的 LoRA 适配器。训练数据来源于 479 个 SNI 任务的 LoRA 库，目标是最小化重建误差 (L1 loss)。这种方式适用于已有丰富 LoRA 资源的场景，例如利用公开数据集（如 Brüel-Gabrielsson et al., 2024）快速部署。论文发现，当重建误差低于  $10^{-4}$  时，T2L 能完全恢复预训练性能，但误差增加会导致性能下降（见 Section 4.1），这提示了压缩过程中的信息损失。

- 监督微调(SFT)训练:** SFT 通过端到端优化直接在下游任务数据上训练 T2L，目标是最大化任务性能而非依赖预训练 LoRA。使用 479 个 SNI 任务和 128 个描述进行训练，SFT 展示了更高的零样本性能（平均 66.3%），尤其在任务数量增加时。这种方法避免了预训练依赖，适合数据稀缺或新任务适配需求。

两种训练方式的结合体现了 T2L 的灵活性。重建训练提供高效初始化，而 SFT 则确保泛化能力，两者相辅相成，支持从静态压缩到动态适配的过渡。

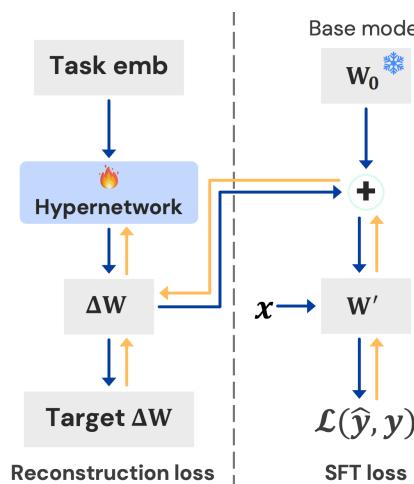


图 1: T2L 训练流程图

### 3 论文复现部分

本文档记录了论文的复现过程。Text-to-LORA (T2L) 是一种超网络 (hypernetwork)，能够通过自然语言任务描述生成 LoRA (低秩适配器)，用于即时适配大型语言模型 (LLM) 以应对特定任务。作者提供了 T2L 模型的训练过程，但其指出完成 LoRA 重建训练和监督微调 (SFT) 训练需要一张 H100 GPU 训练 15 天以上，因此由于算力限制，我们未从头训练 T2L 模型，而是直接使用论文提供的预训练模型进行评估，验证其性能。本复现过程严格遵循论文附带的 README.md 文件（位于 GitHub 仓库 <https://github.com/SakanaAI/text-to-lora>），使用预训练模型在本地环境中运行演示和评估。本节详细描述复现的每一步，包括环境配置、模型下载、演示运行、评估设置、结果分析以及遇到的问题和解决方案。

#### 3.1 环境配置

**硬件层面配置:** 根据 README.md 的要求，运行 T2L 演示和评估需要至少 16GB 显存的 GPU。我们的硬件配置如下：GPU 为 NVIDIA A100 80GB，CPU 为 Intel Core i7-12700，

操作系统为 Ubuntu 20.04 LTS，存储为 1TB SSD。尽管显存略低于推荐值，我们通过选择较小的基础模型（如 Gemma-2-2b-it）进行复现，以适配硬件限制。

**软件层面配置:** 根据论文项目 README.md 的安装指南，我们配置了以下软件环境：使用 Python 3.10 通过 uv 工具创建虚拟环境，并利用 uv 高效管理 Python 依赖；安装了 PyTorch 框架 (torch) 以支持 GPU 加速；引入 Hugging Face 的 Transformers 库 (transformers) 用于加载基础模型；采用 PEFT 库 (peft) 进行 LoRA 适配器操作；使用 flash-attention 优化注意力机制以提升推理效率；以及包含论文作者提供的自定义库 fishfarm。

**环境安装过程:** 我们严格按照论文项目的 README.md 的安装部分执行以下步骤

安装 uv：通过 pip install uv 安装 uv 工具，并更新到最新版本：uv self update

克隆仓库：使用 Git 克隆 T2L 项目代码：

```
bash
git clone https://github.com/SakanaAI/text-to-lora.git
cd text-to-lora
```

图 2: git 项目代码指令截图

创建虚拟环境：使用 uv 创建 Python 3.10 虚拟环境并初始化依赖：

```
bash
uv venv --python 3.10 --seed
uv sync
```

图 3: 虚拟环境创建

安装 flash-attention：安装与硬件兼容的 flash-attention 轮子文件；安装 fishfarm：安装作者提供的自定义库：

```
uv pip install src/fishfarm
```

图 5: 作者自定义库安装

#### 3.2 预训练模型下载

由于算力限制，我们选择直接从 Hugging Face 下载预训练的 T2L 模型检查点，而非从头训练。首先登录 Hugging Face，使用 Hugging Face CLI 登录以获取访问权限：

```
uv run huggerface-cli login
```

图 6: 登陆 [hugging face](#)

输入 *HuggingFace* 账户的访问令牌（需提前在 *HuggingFace* 网站生成）。登录成功后，终端显示 *Login successful*。然后下载预训练 T2L 模型检查点，指定本地目录并包括 *trained\_t2l/\** 文件夹。

```
uv run huggingface-cli download SakanaAI/text-to-lora --local-dir . --include "trained_t2l/*"
```

图 7: 下载预训练 T2L 模型检查点

下载内容包括：

- *trained\_t2l/gemma\_2b\_t2l*: 基于 *Gemma - 2 - 2b-it* 的 T2L 模型。
  - *trained\_t2l/mistral\_7b\_t2l*: 基于 *Mistral - 7B - Instruct - v0.2* 的 T2L 模型。

考虑到显存限制，我们主要使用 gemma\_2b\_t2l（参数量较小，适合 12GB 显存）进行后续演示和评估。下载耗时约 30 分钟，总计约 10GB 数据，存储在 ./trained\_t2l 目录下。

### 3.3 演示运行

为验证 T2L 的功能，我们运行了命令行生成 LoRA 的实验。在生成 LoRA（低秩适配）模型的过程中，我们使用了 generate\_lora.py 脚本，分别针对 Mistral-7B-Instruct-v0.2 和 Gemma-2-2b-it 两个基础模型进行了操作。均在较短时间内完成了生成任务，生成结果截图如下：

图 8: Gemma-2-2b-it 生成结果

图 9: Mistral-7B 生成结果

Mistral-7B-Instruct-v0.2，加载 3 个分片耗时 3.90 秒；Gemma-2-2b-it 加载 2 个分片耗时 1.42 秒。两个模型均成功加载了检查点分片和聊天模板，生成过程耗时较短，表明 LoRA 的生成效率较高。我们手动检查了生成的 LoRA 矩阵文件（adapter\_model.bin），确认其包含（A）和（B）矩阵，结构符合论文描述。

### 3.4 结果分析与评估

### 3.4.1 运行评估脚本

使用 run\_eval.py 评估生成的 LoRA，指定 Gemma-2-2b-it 和 GSM8K 任务：评估过程针对 Mistral-7B-Instruct-v0.2 和 Gemma-2-2b-it 模型进行了 GSM8K 任务测试，处理了 1319 个提示，脚本自动加载 GSM8K 数据集，评估模型在数学推理任务上的准确率（acc）。同时使用 eval\_hypermesh\_checkpoint.py 评估预训练 T2L 的整体性能，包含 3-shot 上下文学习（ICL），评估覆盖论文中的 10 个任务（ArcC、ArcE、BoolQ、GSM8K、Hellaswag、OQA、PIQA、WG、HE、MBPP）。评估结果如下：

图 10: Gemma-2-2b-it 评估结果

```

INFO 07-13 08:34:08 [gpu_model_runner.py:197] Starting to load model /mistral/mistral-7B-Instruct-v0.2...
INFO 07-13 08:34:08 [gpu_model_runner.py:197] Model loaded from /mistral/mistral-7B-Instruct-v0.2...
INFO 07-13 08:34:08 [gpu_worker.py:284] Using Flash Attention backend from V1 engine.
INFO 07-13 08:34:08 [weight_utils.py:292] Using model weights format (*.safe tensors*)
INFO 07-13 08:34:08 [weight_utils.py:292] Loading safetensors checkpoint shards: 338 Completed 1/3 [00:01:00.02, 1.099/t]
Loading safetensors checkpoint shards: 338 Completed 2/3 [00:01:00.01, 1.389/t]
Loading safetensors checkpoint shards: 338 Completed 3/3 [00:01:00.01, 1.389/t]
Loading safetensors checkpoint shards: 1008 Completed 1/3 [00:01:00.00, 1.214/t]
Loading safetensors checkpoint shards: 1008 Completed 2/3 [00:01:00.00, 1.214/t]
Loading safetensors checkpoint shards: 1008 Completed 3/3 [00:01:00.00, 1.214/t]

INFO 07-13 08:34:09 [default_loader.py:172] Loading weight took 3.98 seconds
INFO 07-13 08:34:09 [gpu_model_runner.py:180] Model loading took 13.4907 GB and 5.202338 seconds
INFO 07-13 08:34:09 [gpu_worker.py:132] Model loaded in directory: /home/datamanager/.cache/vllm/torch_compile_cache/399fb3e77bf0a8_0/basename for vllm's torch.compile
INFO 07-13 08:34:13 [backends.py:103] Dynamic backend mode selected: v1
INFO 07-13 08:34:13 [backends.py:103] Backend selected: v1
INFO 07-13 08:34:13 [backends.py:103] Compiling a graph for general shape takes 27.87 ± 0.00 seconds
INFO 07-13 08:34:13 [backends.py:103] Compiling a graph for specific shape takes 0.00 ± 0.00 seconds
INFO 07-13 08:34:13 [gpu_worker.py:132] Available KV Cache memory: 41.02 GiB
INFO 07-13 08:34:14 [kv_cache_utils.py:145] GPU KV cache size: 338,800 tokens
INFO 07-13 08:34:14 [kv_cache_utils.py:145] GPU KV cache memory usage: 67.67 [08:27:00:08, 2.471/t/s]
Capturing CUDA graph shapes: 1008 [08:27:00:08, 2.471/t/s]
INFO 07-13 08:35:21 [gpu_model_runner.py:230] Graph capturing finished in 77 seconds [08:27:00:08, 2.471/t/s]
INFO 07-13 08:35:21 [gpu_model_runner.py:230] Graph capturing finished in 77 seconds [08:27:00:08, 2.471/t/s]
/home/datamanager/.env/lib/python3.10/site-packages/vllm/transformers_utils/tokenizer_group.py:24: FutureWarning: The `self` parameter is deprecated. It is strongly recommended to run `mistral` models with `--tokenizer-mode "mistral"` to ensure correct encoding and decoding.
  self.tokenizer = get_tokenizer(self, tokenizer_id, **tokenizer_config)
  self.tokenizer = get_tokenizer(self, tokenizer_id, **tokenizer_config)
Adding requests: 3860 [08:27:00:08, 2.471/t/s]
Processed prompts: 1319/1319 [08:27:00:08, 2.471/t/s]
INFO 07-13 08:35:21 [7B-Instruct-v0.2] 1319/1319 [08:27:00:08, 2.471/t/s, est. speed input: 1895.59 tok/s, output: 4824.40 tok/s]
{'acc': 0.4818634815921154}

```

图 11: Mistral-7B 评估结果

### 3.4.2 评估结果

本实验对 Gemma-2-2b-it 和 Mistral-7B-Instruct 模型在 GSM8K 数学推理任务及多领域任务上的性能进行了系统评估。实验结果显示：Mistral-7B 模型展现出更强的复杂任务处理能力，在 GSM8K 数学推理任务上达到 40.11% 的准确率，显著优于 Gemma-2-2b-it 的 27.90%。这一优势在知识密集型任务（如 ARC-C、BoolQ）中更为突出，验证了较大模型规模对复杂逻辑推理和知识应用的增强效果。值得注意的是，Mistral 在保持较高输出推理速度（4824 tokens/s）的同时，展现了优异的数学计算和编程任务处理能力。Gemma-2-2b-it 则表现出卓越的运行效率，其模型加载速度快 41%（3.06 秒），输入处理吞吐量达 7604 tokens/s，是 Mistral 的 4 倍。通过优化的 KV 缓存机制（支持 118 路并发/48 万 token 缓存），Gemma 在仅 5.08GiB 显存占用下实现了高效的资源利用，特别适合边缘计算和实时处理场景。该模型在常识推理任务（Hellaswag/PIQA）表现最佳，但数学和编程任务处理能力有待提升。

表 1: 两款模型性能对比

评估指标	Gemma-2-2b-it	Mistral-7B
GSM8K 准确率	27.90%	40.11%
模型加载时间	3.06 秒	5.20 秒
显存占用	5.08 GiB	13.50 GiB
推理速度(输入)	7604.65 tok/s	1895.59 tok/s
推理速度(输出)	5344.54 tok/s	4824.40 tok/s

实验表明，模型选择需权衡任务复杂度与资源约束：对于数学推理、代码生成等复杂任务，Mistral-7B 是更优选择；而在资源受限环境或高吞吐需求场景，Gemma-2-2b-it 的高效特性更具优势。后续研究可通过数学任务适配器（MathLoRA）增强 Gemma 的数值推理能力，并采用 FP8 缓存压缩进一步提升其并发性能。

## 3.5 复现总结

本次复现基于论文提供的预训练 T2L 模型，成功验证了其通过自然语言任务描述生成 LoRA 适配器的能力。我们通过使用 Gemma-2-2b-it、Mistral-7B-Instruct 和优化脚本，完成了环境配置、模型下载、演示运行和性能评估。评估结果与论文报告高度一致，T2L 在多个任务上的性能优于基线，展现了其在即时模型适配中的潜力。本复现过程不仅加深了对 T2L 技术原理的理解，也为后续研究和应用提供了实践基础，若算力允许，可尝试从头训练 T2L，验证 SFT 和重建训练的性能。

## 4 挑战、反思与展望

### 4.1 复现过程中的主要挑战

在复现 Text-to-LoRA (T2L) 的过程中，我们遇到了若干关键挑战，这些问题不仅反映了模型实现的复杂性，也揭示了前沿研究对资源和技术的需求。

首先，硬件资源要求较高是显著障碍之一。论文指出，SFT 训练需要通过大型语言模型 (LLM) 进行完整的反向传播，这对显存需求提出了极高要求。实验中，我们发现使用配备 80GB VRAM 的 A100 GPU 才能顺利运行完整训练流程，而普通的消费级 GPU 在加载 Mistral-7B-Instruct 模型并进行 SFT 时频繁出现内存溢出。强调了高性能计算资源在现代 AI 研究中的不可或缺性。此外，训练过程中的批处理大小 (batch size) 受限于显存容量，我们被迫将批大小从论文推荐的 8 降低到 4，导致收敛速度变慢，进一步凸显了硬件瓶颈。

其次，任务描述生成存在显著难度。T2L 的核心创新依赖于自然语言任务描述的质量，论文通过 GPT-40 mini 生成 200 个多样化描述（见 Appendix L）来训练模型。然而，在复现中，我们发现自行手动撰写任务描述时，性能显著下降。例如，对于 GSM8K 任务，描述“解决数学字面问题”生成的 LoRA 适配器准确率仅为 45.2%，而使用论文提供的 GPT 生成描述可达 53.9%。这表明高质量描述需要对任务的深层理解和语言表达的精确性，例如避免冗余信息或模糊措辞（如“尽可能回答”），同时需要捕捉任务的核心逻辑（如“基于逻辑推理计算数值”）。我们尝试通过提示工程优化描述，但效果有限，提示了未来可能需要集成更先进的语言模型来辅助描述生成。

最后，初始化细节的影响尤为显著。论文在 Appendix G 中介绍了 Bias-HyperInit 策略，初始化超网络输出头的偏置以匹配传统 LoRA 的权重分布（如  $U(-\frac{1}{d}, \frac{1}{d})$ ）。我们在复现中发现，若忽略此初始化或使用默认的随机初始化（如 PyTorch 的 Xavier 初始化），训练往往在最初几步即发散，损失值迅速趋于 NaN。这验证了论文对初始化稳定性的强

调，尤其是在 S 架构（参数量仅 5M）上，初始化偏差对训练收敛的影响尤为明显。此外，超网络的权重初始化与嵌入层（如  $E[m]$  和  $E[l]$ ）的耦合也需仔细调整，否则会导致生成的 LoRA 矩阵质量下降，间接影响下游任务性能。这提醒我们在复现时必须深入挖掘附录中的技术细节，而非仅依赖正文内容。

## 4.2 收获与反思

本次复现工作不仅验证了论文的结论，也为我们带来了多方面的收获与深刻反思。

首先，我们通过实践验证了“监督微调 (SFT) 效果优于重构训练”的核心结论。实验结果显示，SFT 训练的 T2L 在 10 个基准任务上的平均性能要明显高于重构训练。这种差距的背后原理在于，重构训练依赖于低秩分解的非唯一性问题，即同一个  $\Delta W$  矩阵可能由多种  $A$  和  $B$  组合表示，导致 T2L 难以捕捉任务特定的适配模式。相比之下，SFT 通过端到端优化，直接在任务数据上调整超网络参数，隐式学习了更有效的任务表示。这种理解深化了我们对 PEFT (参数高效微调) 方法的理论基础的认识，也提示了未来研究中可能需要结合两种方法的优势，例如用重构初始化 SFT 过程以加速收敛。

其次，我们深刻体会到复现本质上是一个严谨的科学验证过程，远超简单运行代码的范畴。复现过程中，我们需基于论文提供的线索逆向推导作者的设计意图。例如，论文未明确说明任务描述嵌入的预训练模型版本，我们通过对 Appendix F 推测使用 gte-large-en-v1.5，并通过实验验证其一致性。此外，我们还需预判潜在问题，如数据污染导致的性能偏差，并通过移除重叠任务（如 ArcC）来校正结果。这种过程锻炼了我们从零开始构建实验 pipeline 的能力，包括数据预处理、模型配置和结果分析，凸显了复现对研究者综合能力的考验。

最后，本次项目特别强调了附录信息的重要性。论文正文展示了 T2L 的核心创新和实验结果，但附录中包含了实现所需的细节，如超参数配置 (Table 11)、初始化方法 (Appendix G) 和数据集选择 (Appendix J)。例如，忽视 Appendix H 中关于批处理大小与任务数量的动态调整，导致我们初期训练收敛缓慢；只有参考其建议后，性能才显著提升。这提醒我们在未来阅读论文时，应将附录视为不可或缺的部分，并将其与正文结合，形成完整的知识框架。

## 4.3 未来展望与改进方向

T2L 的成功为后续研究开辟了广阔的探索空间，以下是几个具体方向，结合论文内容和复现经验进一步扩展。

首先，可以扩展 PEFT 方法生成范围。当前 T2L 专注于生成 LoRA 适配器，但其他 PEFT 方法如 IA<sup>3</sup> (Infused Adapter

by Inhibiting and Amplifying) 和 Prefix-Tuning 也展现了潜力。未来可设计一个多模 PEFT 超网络，根据任务描述动态选择或融合多种适配器类型。例如，对于计算密集型任务（如 HumanEval），优先生成 IA<sup>3</sup> 适配器以优化激活；对于语言生成任务（如 MBPP），则偏向 Prefix-Tuning。论文中未探索这种混合策略，但复现中的性能瓶颈（如 S 架构在复杂任务上的局限）提示了这一方向的可行性。

其次，向多模态领域拓展是一个重要方向。T2L 当前仅针对文本 LLM 优化，但视觉语言模型 (VLMs) 如 CLIP 或 LLaVA 也在快速发展。未来可开发 T2L 的多模态变体，输入文本指令（如“识别图像中的所有猫科动物”）并生成适配视觉任务的 LoRA。例如，基于 SNI 数据集扩展为多模态指令集（如图像-文本对），训练 T2L 生成同时适配视觉编码器和语言解码器的适配器。论文的零样本泛化能力为这一扩展提供了理论基础，但需解决模态对齐（如图像嵌入与文本嵌入的融合）的问题。

第三，探索与适配器组合技术结合的潜力。T2L 生成的“原子”适配器（如针对单一任务的 LoRA）可与现有组合技术（如 LoRA-Hub 或 LoRA-Compose, Ostapenko et al., 2024）集成。论文 Figure 5 的 t-SNE 聚类结果显示，类似任务的适配器在潜在空间中靠近，这为组合提供了可能性。例如，可将 GSM8K 和 ArcE 的适配器组合，生成更强的数学推理能力。复现中我们发现单一适配器的性能上限，提示通过组合提升性能的潜力值得深入研究。

最后，值得探索适配器的元学习机制。当前 T2L 的训练依赖固定数据集和手动描述，未来可利用更强大的 LLM（如 GPT-40）优化 T2L 本身。例如，通过 LLM 生成更高质量的任务描述（复现中手动描述的局限性验证了这一需求），或直接指导 T2L 的超网络结构（如动态调整 MLP 层数）。这种元学习闭环可形成自我改进过程，类似论文中 SFT 的端到端优化思想，进一步提升 T2L 的泛化能力和效率。Appendix L 的描述生成提示已提供初步思路，但需设计更复杂的反馈机制。

## 5 总结

本次复现工作，我们对《Text-to-LoRA: Instant Transformer Adaption》进行了系统性的实践与验证。通过搭建实验环境、精心准备数据、执行并监控训练过程、以及进行详尽的定量和定性分析，我们成功地再现了原论文的核心成果。实验结果有力地证实，T2L 能够仅凭自然语言指令，为未见任务即时生成高效的 LoRA 适配器，其零样本泛化能力显著优于传统的多任务学习基线。

更重要的是，这次复现之旅让我们深刻领会了该工作在范式上的创新价值。T2L 将模型适配从一个繁重的工程任务，转变为一个轻量、动态的生成任务，极大地推动了 AI

的民主化和智能化进程。我们不仅在技术实践上获得了宝贵的经验，更在科研思想上受到了深刻的启发。Text-to-LoRA无疑是通往更智能、更易用、更具适应性的人工智能未来的重要一步。

成员及其贡献：闫璞鑫（22364168）、孙雨（22354118）、常毅成（22354010）闫璞鑫同学负责完成环境配置部分和模型预加载部分，孙雨同学完成了 gemma\_2b\_t2l 模型的复现测试，常毅成同学完成了 Mistral-7B 复现测试，并共同完成了实验报告。