

大模型微调范式：监督、指令与参数高效策略 深度解析

2025 年 10 月 25 日

1 监督微调（Supervised Fine-Tuning, SFT）

1.1 目标定义与损失设计

监督微调旨在将预训练语言模型适配到具体任务，通过少量标注样本更新模型参数。常见设置包括：

- **条件生成任务**：例如摘要、翻译、代码生成，采用自回归交叉熵损失，保持与预训练期望一致。
- **分类与检索**：加入任务特定的线性头或提示模板，将输出序列映射为标签概率或稠密向量。
- **多任务联合**：采用加权损失或多头架构，在同一模型中兼容问答、理由生成与工具调用。

针对小数据集，常使用标签平滑、Mixout、R-Drop 等正则化方法缓解过拟合。实践中还需规划适当的学习率和 warmup，以防止对预训练知识的灾难性遗忘。

1.2 数据工程与样本构建

高质量的标注数据直接决定 SFT 效果：

- **任务分析**：明确输入模式（问题、上下文、历史对话）与输出格式（简短答案、逐步推理），制定模板。
- **多阶段标注**：先由模型自动生成候选，再由人类审校；结合双人标注和仲裁流程衡量一致性。
- **样本增强**：通过 paraphrase、反事实编辑、知识扩写，提高数据多样性；保持难例与简单样本的平衡。

数据治理环节需要记录标注质量、审核意见与版本历史，以便对微调后的行为进行责任追溯。

1.3 训练流程与评估指标

SFT 通常采用较小 batch size 与较短训练轮次。关键实践包括：

- **分层学习率**：对底层 Transformer 施加较小学习率，对新加入的分类头使用更大更新比例。
- **梯度裁剪**：控制梯度范数，防止少量噪声样本导致发散。
- **持续评估**：监测任务准确率、BLEU/ROUGE、困惑度，同时进行人工抽检，评估可解释性、事实性与安全性。

2 指令微调 (Instruction Tuning, Chat Tuning)

2.1 指令语料构建与对齐维度

指令微调通过引入多样化任务指令与示例对话，使模型学会遵循人类意图。数据来源包括 FLAN、Self-Instruct、OpenOrca 等。关键覆盖维度：

- **任务多样性**：包含分类、抽取、生成、推理、工具调用、数学等，以提升模型泛化能力。
- **角色与语气**：指令中添加“你是一名导师/律师/医生”等角色设定，帮助模型适应不同上下文。
- **对齐层级**：覆盖直接回答、追问澄清、拒答敏感需求、解释推理过程等行为模式。

2.2 自监督扩展与半自动标注

为了降低人力成本，常采用模型辅助生成 instruction 数据：

- **Self-Instruct**：预训练模型根据 seed 指令自动扩展新任务，再由人类筛选与修订。
- **反事实指令**：生成要求模型拒绝或给出警告的场景，提升安全性。
- **难例挖掘**：通过评估模型失误案例，生成针对性对话补充。

需要设计质量控制指标，如平均响应长度、引用来源比例、拒答准确率，确保数据覆盖真实应用需求。

2.3 聊天调优与多轮上下文管理

Chat Tuning 专注于多轮对话场景：

- **系统指令与记忆：**在上下文前置系统消息定义行为准则，并维护有限历史窗口或外部记忆模块。
- **对话状态压缩：**使用摘要或关键词提取的方式，将长对话压缩为短上下文，降低推理成本。
- **拒答与安全策略：**结合 RLHF 或 DPO，确保在越界请求时回复恰当的拒绝或引导。

评估时需模拟真实用户对话，关注连贯性、任务成功率与安全响应比率。

3 Prompt 模板与系统角色 (ChatML, Alpaca Format)

3.1 模板设计原则

Prompt 模板决定模型的输入结构和上下文提示：

- **结构化包裹：**明确分隔系统、用户、助手角色，保证解析稳定；ChatML 使用 `<|system|>`、`<|user|>` 标签。
- **指令显式化：**采用 “Instruction + Input + Output” 格式 (Alpaca Format) 强调任务说明与上下文。
- **约束与示例：**在模板中嵌入风格说明、禁止事项、示例对话，减少模糊理解。

3.2 模板适配与多任务共存

面对多任务训练，需要设计可扩展模板体系：

- **字段可选：**对于无输入任务，提供空输入占位；对复杂任务，支持附加元数据如语言、长度限制。
- **自动化生成：**通过模板渲染器将结构化数据映射为文本，避免人工拼接造成错误。
- **跨平台一致性：**保证训练与线上推理使用同一模板，避免指令偏移。

3.3 系统角色与安全护栏

系统角色定义模型的全局行为边界：

- **行为准则：**说明模型价值观、优先级，如“遵守安全政策”“优先提供事实信息”。
- **工具调用：**在系统消息中描述可用工具、调用格式与错误处理流程。
- **角色切换：**为客服、写作助手、编程助手等预设不同系统消息，实现场景化响应。

设计系统消息时，要防止用户提示覆盖系统约束，可加入优先级声明与拒答策略。

4 参数高效微调 (PEFT: LoRA, QLoRA, Prefix-Tuning)

4.1 LoRA 与低秩适配

LoRA 通过为权重矩阵添加低秩分解并仅训练低秩参数，实现显存友好的调优：

- **原理：**在权重 W 上叠加 BA ，其中 $A \in \mathbb{R}^{r \times d}$ 、 $B \in \mathbb{R}^{k \times r}$ ，训练期间仅更新 A, B 。
- **优势：**显著降低参数量；支持快速切换任务特定 LoRA 模块；兼容混合精度。
- **部署：**可将 LoRA 参数与主模型合并，或按需加载以支持多任务路由。

4.2 QLoRA 与低比特量化

QLoRA 在 LoRA 基础上结合 4-bit 量化，实现消费级 GPU 上的高效微调：

- **量化策略：**使用 NF4 (Normalized Float 4) 量化权重，保持幅度信息，结合 double quantization 减少量化误差。
- **优化器状态：**通过 paged optimizer 将状态存储在 CPU 内存，配合分块加载缓解显存压力。
- **训练流程：**采用 4-bit 权重 + 16-bit 激活的组合，梯度仍以 FP16/FP32 累积，保证稳定。

4.3 Prefix/Prompt-Tuning 与可插拔控制

前缀微调通过优化附加的连续提示向量实现任务适配：

- **Prefix-Tuning：**为每一层注意力引入可学习的前缀键值对，保持原始权重冻结。

- **P-Tuning v2:** 将可学习提示插入到 embedding 层，与虚拟 token 共享，适配更深层模型。
- **组合策略:** 将 Prefix 与 LoRA 结合，实现细粒度任务控制；或通过 Router 根据输入选择合适前缀。

PEFT 方法适合多任务部署，通过模块化加载实现快速切换。需要在评估阶段比较完整微调与 PEFT 的性能差距，并监测泛化与安全性。

参考文献

- Wei et al. “Finetuned Language Models Are Zero-Shot Learners.” ICLR, 2022.
- Chung et al. “Scaling Instruction-Finetuned Language Models.” arXiv, 2022.
- Dettmers et al. “QLoRA: Efficient Finetuning of Quantized LLMs.” arXiv, 2023.
- Hu et al. “LoRA: Low-Rank Adaptation of Large Language Models.” ICLR, 2022.
- Lester et al. “The Power of Scale for Parameter-Efficient Prompt Tuning.” EMNLP, 2021.