

- ✓ 测试集准确率94%。
- ✓ 支持日均3.2千次查询，API P99延迟<1.2s。
- ✓ 技术咨询重复问题减少67%，人工抽检满意度93%。

工作经验

广州小飞侠教育科技有限公司 - 大模型算法工程师	2024-07 ~ 2025-06
<div>1. 文档分块与向量化，设计分块策略，选择嵌入模型，构建向量数据库（pgsql）。</div> <div>2. 检索算法优化，实现混合检索（关键词BM25+向量检索），调整相似度阈值，加入元数据过滤（如文档来源、时效性）。</div> <div>3. 性能调优，量化嵌入模型（FP16→INT8），测试HNSW/PQ等索引算法。</div> <div>4. 实现监督微调（SFT）、RLHF（PPO/DPO）或对比学习（CPT）</div> <div>5. 调试LoRA/QLoRA/P-Tuning等参数高效微调方法</div> <div>6. 实现量化（GPTQ/AWQ）、模型剪枝、知识蒸馏（如DistilBERT→TinyLLM）</div> <div>7. 部署vLLM/TensorRT-LLM等推理框架，优化KV Cache</div>	

腾讯微保 - 大模型算法实习生	2023-10 ~ 2024-01
<div>1. 参与大规模语料清洗、数据预处理（去噪、去重、标准化）</div> <div>2. 构建高质量的指令数据集（如Self-Instruct数据生成）</div> <div>3. 调试LoRA/QLoRA/P-Tuning等参数高效微调方法</div>	

上海梦孚教育科技有限公司 - NLP算法实习生	2023-06 ~ 2023-08
<div>1. 参与大规模语料清洗、数据预处理（去噪、去重、标准化）</div> <div>2. 对Bret、CILP等模型进行微调</div> <div>3. 设计课程学习（Curriculum Learning）策略或数据采样策略</div> <div>4. 实验跟踪：Loss曲线分析、训练效率优化（如梯度裁剪、学习率调度）</div>	

技能特长

- 熟练使用Python编程，熟悉Linux开发环境及Shell脚本，了解C/C++基础开发。
- 熟悉PyTorch框架，熟悉Hugging Face生态，掌握LangChain等大模型应用开发工具链。
- 熟悉大语言模型（LLM）核心技术，包括检索增强生成（RAG）、提示词工程（Prompt Engineering）、模型微调与对齐方法、Agent。
- 熟悉DeepSpeed分布式训练框架，具备多机多卡GPU集群训练经验。
- 掌握高效微调技术（LoRA、QLoRA、P-Tuning），熟悉LLaMA-Factory等轻量化训练框架。
- 熟悉大模型高性能推理框架（vLLM），具备模型服务API开发经验。
- 掌握Docker容器化部署，能独立完成模型封装、服务部署及运维。

教育背景

湖南农业大学 - 智能科学与技术 - 本科	2020-09 ~ 2024-06
<div>• 大学生创新创业项目负责人</div> <div>• 人工智能实验室负责人</div>	