



北京邮电大学  
BEIJING UNIVERSITY OF POSTS AND TELECOMMUNICATIONS

韩竞骁

13624298637 | hanjingxiao@bupt.edu.cn | 北京  
微信: bupt\_ee\_hjx | 岗位: 机器学习



## 教育经历

北京邮电大学	2023年09月 - 2026年06月
通信工程 硕士 信息与通信工程学院	北京
一等学业奖学金 GPA: 90.5/100	
哈尔滨工业大学	2018年08月 - 2022年06月
通信工程 本科 电子与信息工程学院	哈尔滨
一等人民奖学金 GPA: 89.5/100	

## 实习经历

科大讯飞股份有限公司	2025年01月 - 2025年03月
助理算法工程师 核心研发平台	北京
工作描述: 微调多模态大模型, 提升其理解文档图片的能力。	
工作内容:	
1. 评测主流开源模型在文档QA的效果, 研究视觉token数/ViT参数量/切图方式与文档理解能力的关系。	
2. 两阶段训练InternVL2.5-8b:	
• 阶段1-视觉文本和文本对齐: 冻结LLM, 增量训练ViT和MLP在表格解析、文本定位等任务;	
• 阶段2-下游任务指令微调: 解冻LLM, 冻结其他。视觉QA、要点总结、文本阅读等多任务学习。	
3. 设计特殊token解决跨行跨列表格不能转markdown代码的问题。	
4. 设计多页文档的文本定位和解析特定页文字的任务, 提升模型区分多页文档的能力。	
北京罗克维尔斯科技有限公司 (理想汽车)	2024年08月 - 2024年12月
大模型算法实习生 空间AI (理想同学)	北京
工作描述:	
1. 后训练阶段微调数据冗余的背景下, 设计数据去重和高质量数据选择pipeline。	
2. 设计训练策略, 提高车机模型的汽车专业知识水平和模糊指令回答能力。	
工作内容:	
1. 两阶段的多样性约束的数据集去重研究:	
• 阶段1-基于相似度的粗筛: 先后使用TF-IDF和bge-base-zh-v1.5生成的相似度向量之后聚类过滤。	
• 阶段2-基于半监督的分布外数据筛选: 已标数据记为1, 未标数据标为0。构建二分类模型做k折交叉验证, 选出每一折测试集为0的数据。	
2. 微调数据选择方法研究: 训练指令完整度奖励模型; 高PPL高reward的指令数据被召回。	
3. 构造对抗样本, 提高模型在用户问题不清晰的情况下的回答能力。	
4. 学习率退火阶段的领域知识注入实验, 得出专业知识学习的最佳实践。	

## 竞赛经历

WWW2025-阿里天池 多模态对话意图识别挑战赛	2024年12月 - 2025年01月
题目背景: 给定淘宝客服中的图文对话数据, 使用10B以内的Qwen系列模型, 在有限集合内输出用户意图。	
题目难点: 1.直接生成的准确率低。2.类别数目过多。3.大量未标注数据。	
解决方案:	
1. Qwen2-VL-7b隐藏层接分类头, 不生成文字, 直接输出类别。	
2. 采用两阶段的训练策略, 提升长尾类别的识别性能。	
• 阶段1: 均匀采样数据, 训练LLM和分类器, 目的是学习到长尾数据集上的最佳特征表示。	
• 阶段2: 逆采样数据, 得到平衡子集。冻结LLM, 训练分类器。	
3. 伪标签法: 阶段1模型在未标注数据中筛选高置信度的尾部类别数据, 将其添加到阶段2的数据集中。	

## 科研经历

北京电信-北邮校企合作项目: 运营商网络流量预测算法研究	2023年09月 - 2024年07月
项目描述: 设计针对运营商城域网流量的时空预测算法。高精度预测算法赋能资源分配、异常检测等下游任务。	
项目成果:	
1. 发表两篇论文, 一篇TII(IEEE Trans on Industrial Informatics)在投, 一篇发表在ICCIP24, 获best paper。	
2. 预测算法部署到电信城域网运维系统中, 预测误差5%范围内节点占比98%。	
3. 该项目评选为北邮研究生科研创新 A 级项目。	

## 个人总结

- 本科辅修计算机, 学习了数据结构、CSAPP、计算机网络课程。
- 两段大模型实习, 一篇时序方向Trans在投和一段多模态比赛经历。在LoRA/分布式训练/多模态方面有实践经验。