

Seed-X: Building Strong Multilingual Translation LLM with 7B Parameters

ByteDance Seed

Full author list in Contributions

简介

本文提出Seed-X，专为多语言翻译(Multilingual Translation)任务设计的开源7B LLM。Seed-X 7B覆盖28种语言，经历了完整的pre-training from scratch、sft和RL三阶段训练，每个阶段的训练都有比较详细的技术细节：1) 预训练阶段，构建了高质量的单语与双语语料，为了平衡单语和双语训练数据，采用三阶段训练策略；2) 在sft阶段，通过多样化的prompt template以及引入CoT template，在将base model转为翻译模型时促使其更好地理解上下文掌握语义细节；3) 在rl阶段，设计了两种reward 分别是基于偏好数据训练的reward model和DuPO reward function(后续会有相关论文)，使用ppo算法训练。从实验效果来看，Seed-X以7B的大小展现了很强的翻译能力。

背景

针对多语言翻译(Multilingual Translation)任务，基于LLM的翻译方案已经比传统翻译模型有了显著提升，当然还可以继续优化，另外一个更大的问题是开源翻译模型比OpenAI/Anthropic/Google等闭源模型的效果差距还是挺大的，尤其在低资源小语种以及考虑文化背景方面比较突出。同时，作者认为目前开源社区在如何构建开源SOTA翻译模型仍缺乏系统性的方法论，为此作者团队提出了Seed-X系列模型，也希望为多语翻译模型的发展道路给出较为明确的一种方案。

预训练

数据包含单语和双语平行语料两种类型。单语言数据约6T，并且刻意排除代码和数学等STEM领域的数据，将LLM能力完全用于多语言理解与翻译建模。平行语料数据的构建则是迭代的过程，简单来说，先找一个高质量的种子数据集，然后训练初始翻译模型然后对单语数据生成伪平行数据，再不断迭代过滤双语数据、重新训练模型，最终得到质量比较高的平行语料。

如何融合单语和平行语料训练也是一个问题，作者设计了三阶段训练：

- 以中英等主要语种的单语数据训练LLM的基础语言能力；
- 逐步引入更多低资源语言的单语和双语数据，提升多语言迁移和泛化；
- 只用平行语料训练全面提升翻译能力，或者理解为LLM转变为垂直领域的翻译模型。

sft和RL

sft阶段是为了强化base model的翻译能力，数据集包含236k条翻译任务指令数据，同时作者认为翻译任务不仅是语言映射，还需要理解语义、文化背景和表达推理，因此额外设计了CoT prompt template。

RL阶段主要目的是让翻译模型考虑到人类偏好，针对并行语料和单语言两种场景设计了两种reward方式：

- 针对资源丰富的语种，收集人类偏好数据，基于base model tuning得到reward model
- 针对资源稀缺的小语种，收集不到人类偏好数据，设计了DuPO reward function，简单来说，执行这样的翻译流程 $A \rightarrow B \rightarrow A'$ ，通过比较A和A'的相似度作为 $A \rightarrow B$ 翻译质量的reward

思考

机器翻译曾是LLM时代之前NLP领域第一大研究方向，诞生了encoder-decoder、attention、Transformer等模型，遗憾的是个人当时没有机会接触到相关的工作，后来LLM颠覆了机器翻译任务，也就没有再继续关注了。本文的Seed-X将机器翻译作为核心目标，经历了预训练、sft和rl的完整训练路径，从实验结果看，7B小模型直追闭源大模型，不得不感慨，技术发展真快啊，估计在其他有重要落地价值的领域，也会出现这类专精特性小模型吧。