

序章

背景介绍

【缘起】：OpenLLM Talk 这个事情起源于 20230603 OpenLLM 交流群中大家的一场讨论，本着心动不如行动的想法，我们花了一点时间来将其落地，希望可以为大家提供一个 LLM/NLP 领域的交流平台。——**我们或许不够 AI，但尽量足够 Open；我们也不知道能走多远，但尽量比自己想的更远。**

【结构】：整体上分成本周新闻、本周推荐、本周经典（可选）、本周实践（可选）、free talk 等版块，建议后续最好采用每个版块每期由 1-2 人认领+多人参与贡献+自由讨论的形式。

本期记录

【编号】：OpenLLM Talk 006 (三位数是希望 LLM 的热度+我们的热情+读者的热情可以支撑我们做到三位数)

【时间】：20230729 晚上九点（每周六晚上九点，节假日顺延）

【本期提要】：ICML2023 杰出论文；attention 机制的 bug；llama2 的汉化；多轮对话；DPO 论文阅读；LLM 评估；text2vec；量化；NOPE；长度外推；OpenLLMAI 与实践计划；

【本期贡献者】 - 排名不分先后：

【主持人】：羡鱼（后续每期由大家自行认领）

【编辑】：羡鱼（最好由主持人兼任）

【版块负责人】：多人（后续每期由大家自行认领）

【具体内容贡献者】：请查看具体内容后面的署名，比如问题、回答和观点的来源

【talk 视频】：

注意事项

【talk 模板】：<https://zhuanlan.zhihu.com/p/640522290>；可参考模板进行贡献

【小要求】：主持人及版块负责人认领之后尽量准时参加，其余同学可自行选择是否参与；

本周新闻

【本周新闻】：LLM/AI news，包括但不限于学术、项目、工业界新闻和进展；多人认领或者直接在此添加，由 **1-2 人认领并汇总**；建议大家都参与进来，相互补充，尽量**减少信息冗余和缺漏**；共~10 分钟；

【贡献者】：

【建议区】：可以考虑 GitHub 的讨论区，看个人习惯；论文可以写个摘要；

学术

注：论文+重点

ICML2023 杰出论文出炉

https://mp.weixin.qq.com/s/cDzUidwGZjIG_KyMi3M9BQ

Attention 机制竟有 bug，Softmax 是罪魁祸首，影响所有 Transformer

<https://mp.weixin.qq.com/s/cSwWapqFhXu9zafzPUeVEw>

项目

chinese 版 llama2

<https://github.com/LinkSoul-AI/Chinese-Llama-2-7b>

ziya 公开直播训练

<https://huggingface.co/spaces/IDEA-CCNL/Ziya-LLaMA2-13B-Pretrain>

工业界

Baby llama2

<https://github.com/karpathy/llama2.c>

SMP 2023 ChatGLM 金融大模型挑战赛

<https://tianchi.aliyun.com/competition/entrance/532126/introduction>

本周推荐

【本周推荐】：本周重点内容推荐和介绍，模型、开源项目、好的资料或课程，建议 1-3 项；共 15 分钟；

【贡献者】：

【提名区】：

【建议区】：

【本期主题】：

资料

一文看懂：如何充分高效训练多轮对话大模型

<https://mp.weixin.qq.com/s/KsbgRNTwXE86kCGTJhu0WQ>

【LLM 系列】对行业大模型的思考 - 黄文灏的文章 - 知乎

<https://zhuanlan.zhihu.com/p/643805698>

DPO——RLHF 的替代之《Direct Preference Optimization: Your Language Model is Secretly a Reward Model》论文阅读

<https://zhuanlan.zhihu.com/p/634705904>

中文 LLaMA&Alpaca 大语言模型词表扩充+预训练+指令精调

<https://zhuanlan.zhihu.com/p/631360711>

模型

项目

Efficient and Effective Text Encoding for Chinese LLaMA and Alpaca

<https://arxiv.org/abs/2304.08177>

杂项

refs:

本周经典-optional

【本周经典】：NLP/LLM 领域的经典话题探讨；~15 分钟；

【贡献者】：

【提名区】：量化

【本周主题】：

本周实践-optional

【本周实践】：NLP/LLM 领域实践经验分享，可以分享自己的实践经验或者他人的实践经验，后面群里也会组织一些实践内容；~15 分钟；

【贡献者】：

【提名区】：

【建议区】：coding 搞起来；后续拉个 read_code/paper 分支，LLM 精读、注释；专门建一个数据专题；

LLaMA2 框架

成员

初七, qwang, Sine, donny, 风吹草地见牛 ...

技术讨论

相关工具: HF/DeepSpeed/Megatron/Ray/RLHF/LLM

参考框架: DeepSpeed Chat ;

Ray: qwang

RLHF: 初七、Yiran (周末)

SFT: Hope(可以写一些代码)、羡鱼 (周末)

pretrain: 羡鱼 (周末)

test datasets:

Debug machine: 用 300m 模型单卡测试, 后期用集群 perf 测试

预期产出:

llama2 架子

垂直领域的 llama2、

第一次会议主要讨论技术方案, 项目开发组织方式, 分工

然后起一个好听的名字

Free Talk

【Free Talk】自由提问, 自由讨论; 在文档里提问或者在群里提问, 建议尽量在此汇总; 如果群里已经有比较好的讨论结果, 也可以将讨论结果搬运过来; 时间不限;

【贡献者】: 羡鱼 (编辑) +OpenLLM 群友

1.Evaluation of LLM\MLLM

对大模型的评测是比较困难的事情，一般人工或者依靠 GPT4 来评价，都是比较费钱的。一些新的评测基准（MMBench）是否可以用 llama2 chat 这种 rlhf 之后的模型来评测，是否可以起到完全相同的效果。

相关领域：模型评测，LLM

答：GPT4 也不一定准确，特别，不建议在被评测的回答中有 GPT4 存在的情况下用 GPT4 评价（既当运动员又当裁判员）；

张拳石老师：可解释评测；

<https://arxiv.org/abs/2304.01083>

质量好不好，哪里好，哪里差？

做一个通用的奖励模型？

感觉奖励模型往往是针对特定模型来训的；

LLM 竞技场，几个 LLM，人类打分；

一个想法：分层的 RM，general-》domain-》task-- 羡鱼

迭代式的 RM，目前 OpenAI、anthropic 等都采用的是多轮迭代式的 RLHF 流程

llama2 70b rlhf 基本上是开源的最强的模型，在一些方面有接近 GPT3.5 的能力

中文方面是不是不一定？llama 基本没做中文；

主观问题往往回答越长越好，直接 len() 都比较靠谱？ :)

2. 关于垂直领域的 text2vec，各位有没有相关的数据集构建或者其他的模型选择的经验？

答：OpenSE，有空我会放个 repo 出来，做一个 SE/text2vec，大体的流程：基础模型--》无监督--》自监督——》细粒度监督训练；

3. FP8 的软硬件支持现状？

答：

为什么 FP16 的值域这么窄？

BF16？

smoothquant

<https://github.com/mit-han-lab/smoothquant>

量化的效果：

FP16、FP8，还有 8 位、4 位、甚至于 3 位、两位？

Qlora：int4 量化；3090 微调 13B；

Yiran

quant后面有很多工作，GPTQ这类的

Yiran

主要是做推理的，节省显存

Yiran

我之前做过Flexgen，16G 卡的能推理170b

宋省身

int8->qlora

宋省身

gptq是支持4bit甚至3bit

宋省身

qlora的推理速度一般，二次量化的计算成本太高

宋省身

llm.int8()和qlora的作者是同一个人

4. 想问下各位大佬，之前有讨论过 NoPE 这篇文章吗？

arxiv.org/pdf/2305.19466.pdf

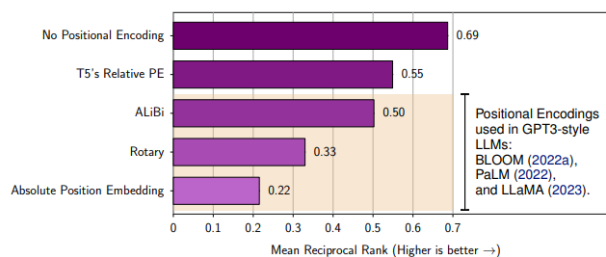


Figure 1: No positional encoding (NoPE) outperforms all other positional encodings at length generalization of decoder-only Transformers (GPT-style) trained from scratch and evaluated on a battery of reasoning-like downstream tasks. This figure shows aggregate ranking of positional encoding methods across 10 tasks.

5. Claude100k 是怎么搞的？

答：

llama 训练时 2k，微调 1000 步到 32k；

微软有个十亿 token 的；

外推最近进展到什么地步？貌似 NTK dynamic 效果超过 16 倍不太行；

[浅谈 LLM 的长度外推 - 知乎 \(zhihu.com\)](https://zhuanlan.zhihu.com/p/64444444)

6. 指令数据里面如果**消解矛盾**？尤其是 GPT4 这种传言上百万的指令数据集。另外，如果保持预训练、SFT、RLHF 部分的一致性，至少像数据层面的一致性？

答：比如说，SFT 部分如果有一些预训练的知识盲区，容易加剧胡说八道的情况。

7. 将 linear transformer scale 到 175B？

<https://arxiv.org/abs/2307.14995>

那个不是没有测eval数据集？

Xue

<https://arxiv.org/abs/2307.14995>

葛春江

transNormer

葛春江

去掉了softmax

Yiran

好像没有做效果相关的实验

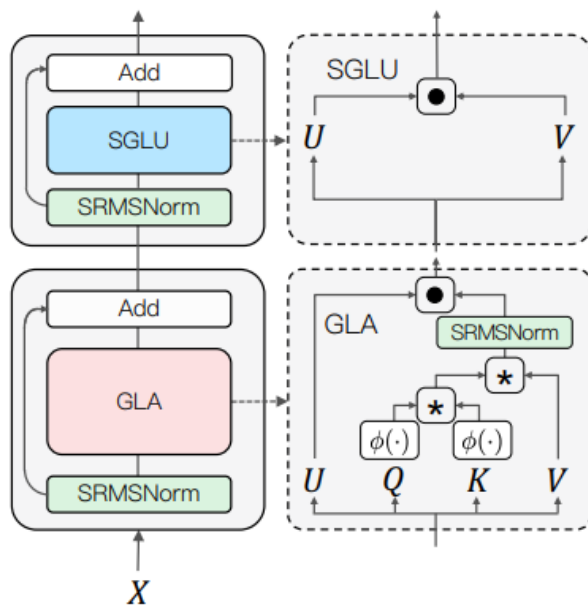


Figure 1: Architecture overview of the proposed model. Each transformer block is composed of a Simple Gated Linear Unit (SGLU) for channel mixing and a Gated Linear Attention for token mixing. We apply pre-norm for both modules.

8.

线上讨论:

1.

群里讨论：

有空会同步，取决于人力，希望大家积极认领~

参考资料

后续计划

我们正式升级为一个不太正式的组织了！叫做 OpenLLMAI。

<https://github.com/OpenLLMAI>

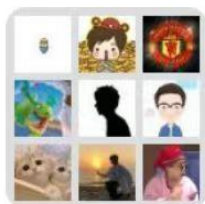
- 正式开启 OpenLLM talk 的运营，P1；
- ChatPiXiu 项目：陆续有一些实践计划，P0；
- <https://github.com/OpenLLMAI/OpenLLaMA2>，P0，doing
- <https://github.com/OpenLLMAI/chinese-llama2>，P0，doing
- OpenSE：检索项目，字符检索+语义检索，P0；
- OpenLLM：LLM 学习和实践项目，P0；
- OpenAIWiki：AI wiki for everyone；
- ChatLover：模拟恋人+爱情助手，P1；

加入/赞助我们

我们非常缺人，也非常缺时间和算力，希望能有越来越多的朋友参与进来，认领 talk 的组织者、主持人、板块的负责人；参与项目后续的开发和讨论等等。

微信群：（请优先加入微信群，如果失效则加入 QQ 群再私聊我进微信群）

（二维码过期了！）



群聊：羡鱼智能-OpenLLM 技术
交流群



该二维码7天内(7月7日前)有效，重新进入将更新

QQ 群：



羡鱼智能-OpenLL...

群号: 740679327



扫一扫二维码，入群聊。



往期精彩

【OpenLLM Talk 005】本期提要：llama2；FreeWilly；LLM 推理与评估；LLM 八股；RetNet；DPO；数据配比 - 羡鱼智能的文章 - 知乎

<https://zhuanlan.zhihu.com/p/645679737>

【OpenLLM Talk 004】本期提要：外挂知识；抱抱脸每日论文；MOSS-RLHF；GPT4 细节；OpenAI 代码解释器；百川 13B；LLM 面经；多轮对话；数学能力；反思；LLM 中的知识 - 羡鱼智能的文章 - 知乎

<https://zhuanlan.zhihu.com/p/643960837>

【OpenLLM Talk 003】本期提要：SuperCLUE-Open；文心盘古；chatlaw；LLM 综述；NTK-Aware Scaled RoPE；10 亿上下文；InternLM；GLM 讲座 - 羡鱼智能的文章 - 知乎

<https://zhuanlan.zhihu.com/p/642376781>

【OpenLLM Talk 003】SuperCLUE-Open；文心盘古；chatlaw；LLM 综述；NTK-Aware Scaled RoPE；GLM 讲座】 【精准空降到 10:10】

https://www.bilibili.com/video/BV1Kh4y1E7nX/?share_source=copy_web&vd_source=9e7882f0ef2735e23d66a6f128612943&t=610

【OpenLLM Talk 002】本期提要：chatgpt 增速放缓；gorilla-cli；RoPE 外推；vllm vs llama.cpp；lora 融合；模型参数和数据之比；OpenSE 计划 - 羡鱼智能的文章 - 知乎

<https://zhuanlan.zhihu.com/p/641285737>

【OpenLLM Talk 001】本期提要：长程记忆；OpenAI 上新；百川智能 7B 模型；State of GPT；位置编码；deepspeed-rlhf；RLHF 数据 - 羡鱼智能的文章 - 知乎

<https://zhuanlan.zhihu.com/p/640275116>

【OpenLLM Talk 000】我们做了一个 LLM 领域的交流平台 - 羡鱼智能的文章 - 知

乎

<https://zhuanlan.zhihu.com/p/636350755>

【OpenLLM Talk 模版】兴趣和热爱胜过一切，OpenLLM 就从这里开始吧！欢迎加入！ - 羡鱼智能的文章 - 知乎

<https://zhuanlan.zhihu.com/p/640522290>