

**Project Name:** Improving performance of large language models in an affordable way

**Project description:**

As large language models show their ability to handle a variety of tasks, people start thinking about incorporating large language models into their own business. However, there are still many concerns for using language models. First, training your own large language model from scratch takes a lot of GPU power which is expensive to startup companies. Second, large language models often provide false or toxic answers which are unforgivable from user's perspective. The main goal of this project is to improve the performance of current language models in an affordable way.

**Assignments:**

1. Weekly reflections on paper reading
2. Final report
3. Potential demo to show improvement

**Related reading:**

1. Transformer: 《Attention Is All You Need》
2. GPT: 《Improving Language Understanding by Generative Pre-Training》
3. GPT-2: 《Language Models are Unsupervised Multitask Learners》
4. GPT-3: 《Language Models are Few-Shot Learners》
5. InstructGPT: 《Training language models to follow instructions with human》
6. Sparrow: Improving alignment of dialogue agents via targeted human judgements
7. LLaMA: 《LLaMA: Open and Efficient Foundation Language Models》
8. Llama 2: 《Llama 2: Open Foundation and Fine-Tuned Chat Models》
9. RLHF:《Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback》
10. RLHF 《A General Language Assistant as a Laboratory for Alignment》
- 11.《Emergent Abilities of Large Language Models》
12. 《Scaling Laws for Neural Language Models》
13. 《Training Compute-Optimal Large Language Models》
14. 《Language Models (Mostly) Know What They Know》
15. 《Scaling Instruction-Finetuned Language Models》
16. 《Proximal Policy Optimization Algorithms》

**Schedule:**

1. I will read through two papers above per week and write reflection on what I learned from these paper and the potential insight for my final report.
2. Final report will be submitted before August 20th.