

lab-lora-tuning-peft

Goal:

Fine-tune the GPT-2 model using LoRA with different configurations to improve text generation and save training time.

Model & Data:

- Base model: GPT-2
- Dataset: 50 prompts from fka/awesome-chatgpt-prompts
- Task: Generate answers to user prompts

LoRA Configs Tested:

- Config 1 (r=4): Output was repetitive (“to to to...”).
- Config 2 (r=8): Slightly better, but still repetitive.
- Config 3 (r=16): Training still running, expected to improve.

Findings:

- GPT-2 doesn't support query_key_value, so that module should be changed.
- Better results with higher r and lora_alpha.
- LoRA works better on newer models like bloom-560m.

Recommendation:

- Use a compatible model (e.g., Bloom).
- Increase training data and epochs.
- Tune generation settings like temperature or top_p.