

< Improving Accuracy of LLM Application >

partnership with Llama & Meta

LLM : trained on "public dataset" and task

↳ they not perform well on application
requiring "proprietary data"



But : Llama : open-source ! → try me on fine-tune the model.

When fine-tune doesn't work?

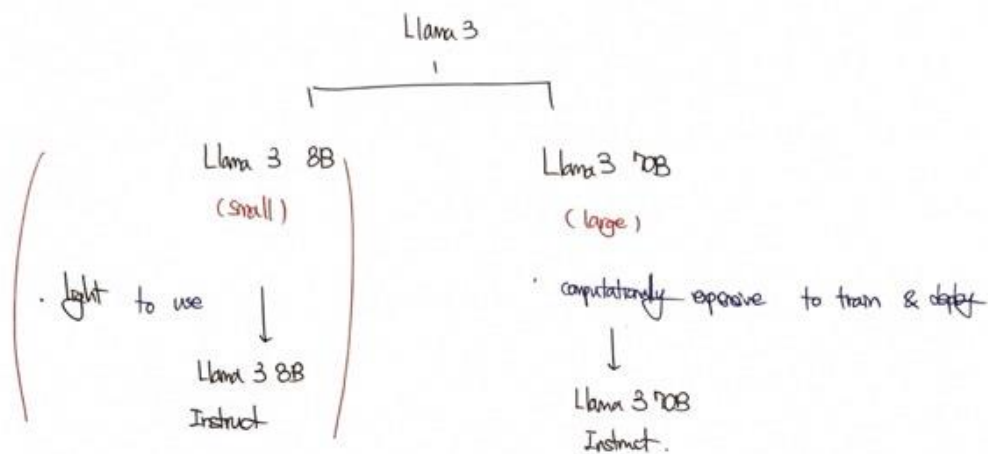
↳ think as : not enough data → no! , the way you fine-tune should be revised.

Fine-tune : slow & costly → But (LoRA , parameter efficient fine-tune .
or
Llama memory tune)

In this course , learn the way

↓
cost ↓ , fast !

Improve Accuracy of LLM by SGL



→ usually use 8B in laptop!

*** hallucinations

- Incorrect or fabricated information generated by the model
- LLM thinks: something slightly right = right.
- Detrimental with facts & when you need precision

*** Common approaches help but not enough:

- Prompt engineering → 26%
- Self-reflection → 26-40%
- Retrieval Augmented Generation → 50%
(RAG)
- Instruction Fine-tune → 40 ~ 60%

How Can fine-tuning reduce hallucination?

- Embed facts into the model to Remove
- of) Hw, instruction fine-tuning isn't the tool to remove the hallucination.

- Memory Tuning (by Llama) allows the models to recall a lot of facts precisely

Important! Without Compromising generalization & instruction-following

- Memory Tuning \rightarrow 95%

\rightarrow Reason: Sampling from a distribution that looks more like this (A)

EX) A SQL Agent.

LLM: can be SQL Agent

\rightarrow but makes hallucinations

\rightarrow fine-tuning \rightarrow Lift the Accuracy!

Is SQL Agent useful.?

1. Optimal UX

faster answer for business users.

2. Operational Efficiency

Less data team time spent answering simple questions.

3. Reliability

More reliable in breadth than business users writing their own queries.

Hallucination In SQL Agent.

1. Invalid SQL

Mixed column names / IDs, formats, or functions.

2. Malformed SQL

Valid, but semantically wrong.

< How to quickly improve accuracy >

↳ New way: Llama Merge Tuning: embeds facts about the database into Llama 3 very efficient.



SQL Agent's workflow with Llama 3

User Question → SQL LLM → Run query → Response user sees

< Diagnose hallucination >

SQL query [mixed
malformed

→ diagnose hallucination: Interactive.

- As you improve the model, you'll continue to dig for hallucination.
- this is how you point the LLM to what it can improve on.

Goal of Evaluation

- Is it improving?
- where is it hallucinating? → where can I fix more?
(improve)

Initial evaluation dataset.

- Start small: 20 ~ 100 data
- Quality > Quantity. Iterative expansion
- Focus on areas of improvement

Strong SQL LLM Specifically

- generated SQL can be a correct semantic match but not exact match
- Evaluate generated SQL against DB and compare SQL exact match

Good Evaluation

1. Quantitative
2. Point out what you can improve
3. Scalable & automated

Practical tips for your evaluation dataset.

- Easiest examples that still fails.
- Use an "adversarial playground"
↓
let your LLM get break down.
And check the boundary.
- Set a next accuracy target for your LLM.
small target → big target.

~~xxxx~~ Fine-Tuning

~~xxxx~~ why Fine-Tuning?

- Fit more data than what fits into the prompt → storing info in weights
- Learn from the data, rather than just get access to it
- Deeper Control of the LLM to achieve what you want
- No accuracy ceiling

Types of Fine-Tuning

Instruction Finetuning

↓
Get a pre-trained LLM to follow Instruction

Memory Tuning

↓
Get an LLM to Not hallucinate

~~xxxx~~ why Pre-training is not enough?

Pre-training

- LLM reads the internet, one token at a time
- Reduces average error over example (generalization error)
- It creates powerful foundation models, learning a lot.

↓

But, there's still a gap in usefulness

- LLM cannot follow instructions
→ It doesn't have a chat in "chatGPT"
- LLM hallucinates on facts
→ Pretty good at everything, but perfect at nothing.
- nearly right answers are still wrong

Memory-Tuning : Teach perfect recall on facts

- Reduce error to zero on facts
- Near-perfect on facts, pretty good at everything else.
- Before : a probability distribution over the answer
↓
- After : only commit to the right answer

Myth1 : Prompting & RAG will solve all my problems?

Following Instructions : Use prompting to get the model to do a little bit better

Fact Recall : Use RAG to get the model to do a little bit better → shift probabilities conditioned on retrieved prompt, but still sampling from a distribution of similar but wrong fact.

Myth2 : Fine-Tuning is still too expensive?

Reality : cheaper than running huge prompts in RAG

PEFT (Parameter-efficient Fine-tuning) : Reduces cost dramatically!

MoE (Mixture of Memory Experts) : turns any LLM into a million-way mixture-of-expert adapters, reducing time by 240x

< Large Multimodal Model Prompting With Gemini >

Gemini : multimodal model from google Deepmind .

< Model Parameter > - In Gemini

Topk : The Top k possible next words

Top P : The Top P possible next words with a cumulative probability $\geq P$

Temperature : choose next word based on it

└─ high Temperature → Randomness : ↑
low " → " : ↓

< Image Requirements , Image Cap >

Image to tokens : Use up to 3000 images

Image type : PNG or JPEG

Image Pixel : unlimited , but large images are scaled down and padded

Image to token : Each Image counts for 256 tokens