# (An Incomplete List of) Popular LLM Evals

| | |
|---|---|
| Human Preferences for chat | **Chatbot Arena** |

| | |
|---|---|
| LLM as a judge for chat | Alpaca Eval<br>MT Bench<br>**Arena Hard V1 / V2** |

It's easy to improve any one of the benchmarks.

| | |
|---|---|
| Static Benchmarks for Instruct LLM | **LivecodeBench**<br>**AIME 2024 / 2025**<br>GPQA<br>MMLU Pro<br>IFEval |

It's much harder to improve **without degrading other domains**.

| | |
|---|---|
| Function Calling & Agent | BFCL V2 / V3<br>NexusBench V1 / V2<br>**TauBench**<br>ToolSandbox |

but it can be much harder to improve some benchmark
但要提升某些基準測試表現

# Do you really need post-training?

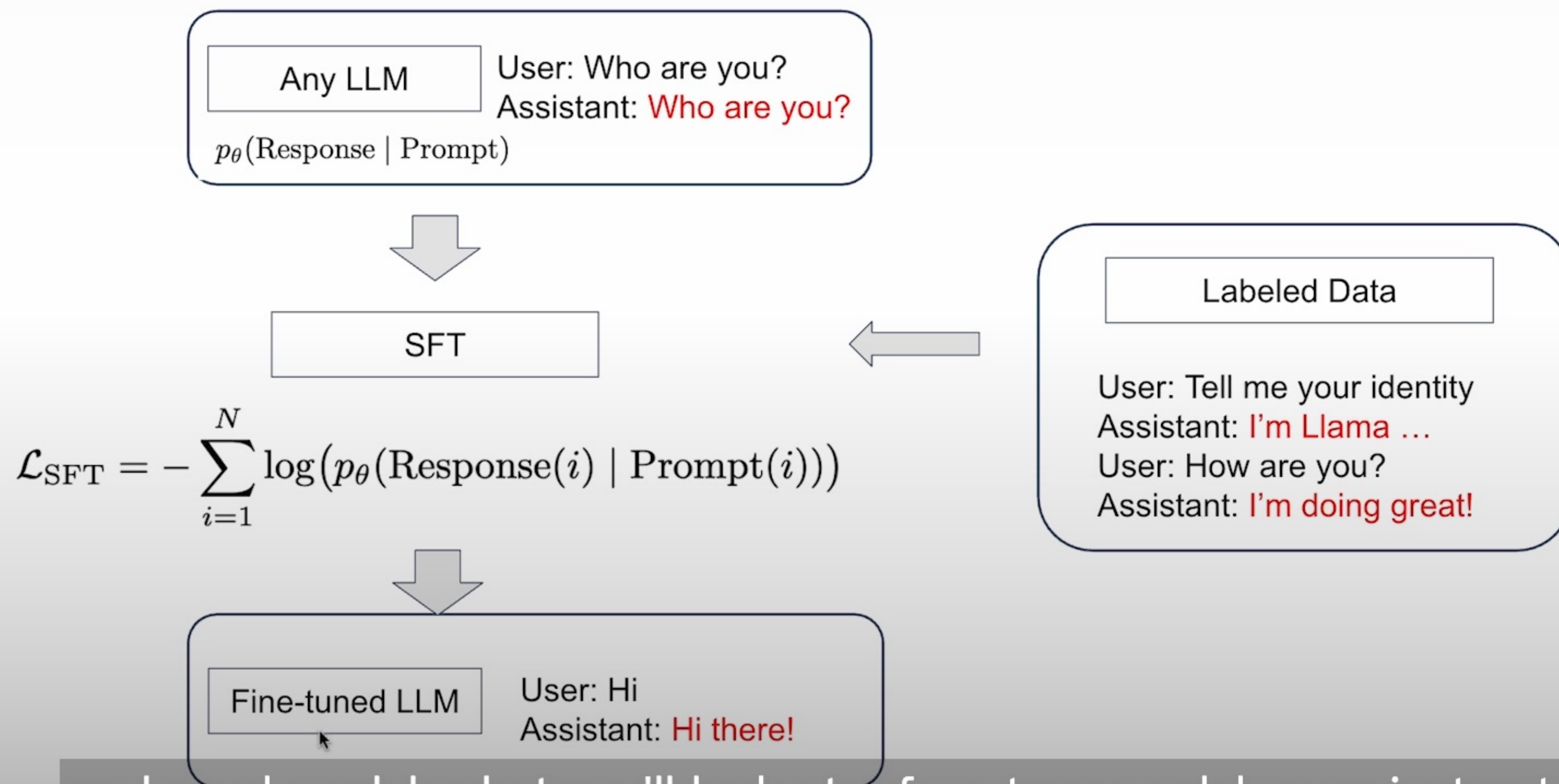| Use Cases | Methods | Characteristics |
|---|---|---|
| Follow a few instructions (do not discuss XXX) | Prompting | Simple yet brittle: models may not always follow all instructions |
| Query real-time database or knowledgebase | Retrieval- Augmented Generation (RAG) or Search | Adapt to rapidly-changing knowledgebase |
| Create a medical LLM / Cybersecurity LLM | Continual Pre-training + Post-training | Inject large-scale domain knowledge (>1B tokens) not seen during pre-training |
| Follow 20+ instructions tightly; Improve targeted capabilities ("Create a strong SQL / function calling / reasoning model") | Post-training | Reliably change model behavior & improve targeted capabilities; May degrade other capabilities if not done right |

# SFT: Imitating Example Responses



Any LLM

User: Who are you?
Assistant: Who are you?

$p_\theta(\text{Response} \mid \text{Prompt})$

SFT

Labeled Data

User: Tell me your identity
Assistant: I'm Llama …
User: How are you?
Assistant: I'm doing great!

$$\mathcal{L}_{\text{SFT}} = -\sum_{i=1}^{N} \log\big(p_\theta(\text{Response}(i) \mid \text{Prompt}(i))\big)$$

Fine-tuned LLM

User: Hi
Assistant: Hi there!

# Best Use Cases for SFT

- **Jumpstarting new model behavior**
  - Pre-trained models -> Instruct models
  - Non-reasoning models -> reasoning models
  - Let the model uses certain tools without providing tool descriptions in the prompt

- **Improving model capabilities**
  - Distilling capabilities for small models by training on high-quality synthetic data generated from larger models
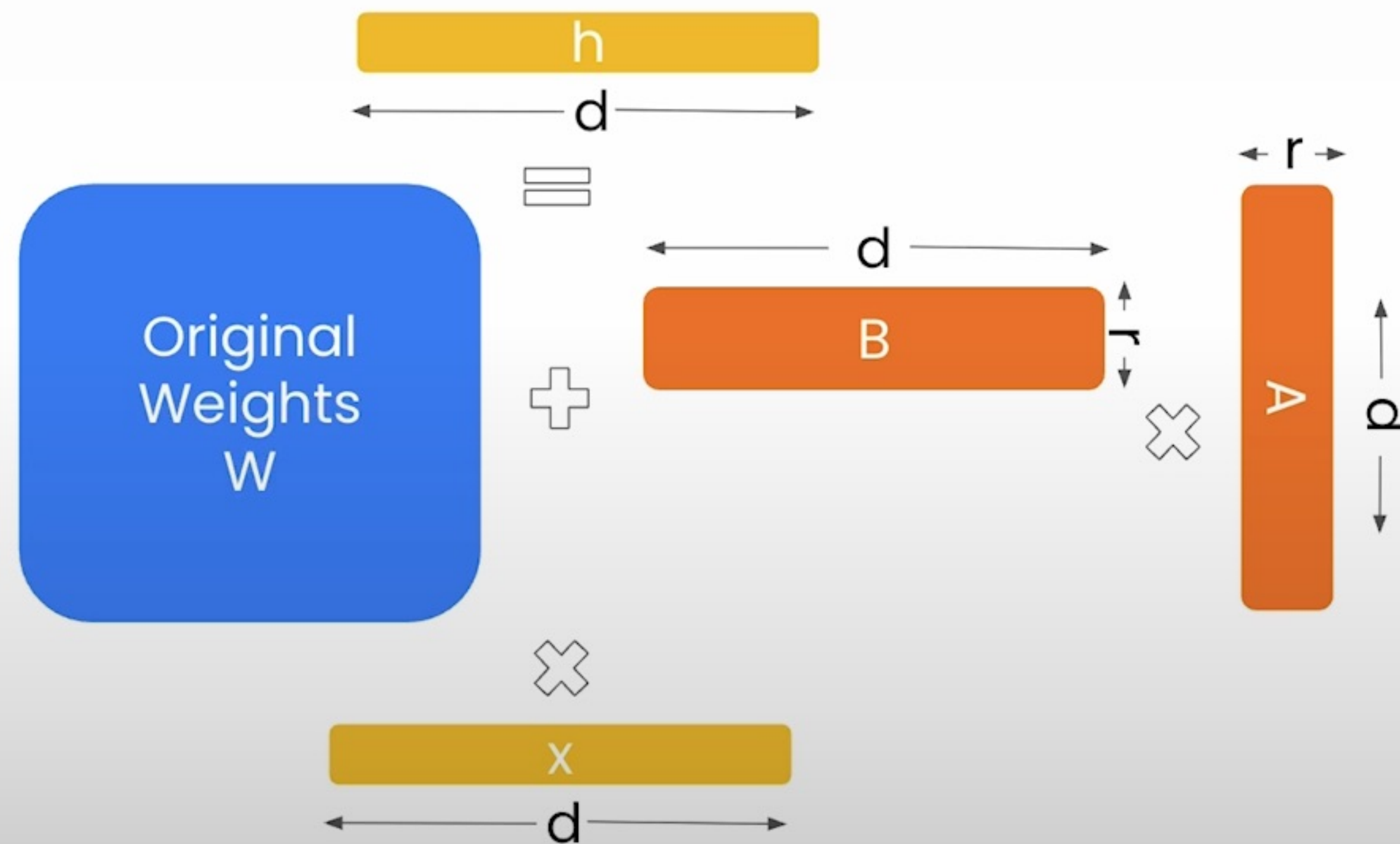
# Full Fine-tuning vs Parameter Efficient Fine-tuning (PEFT)



$$h = (W + \Delta W)x$$
$$W, \Delta W \in \mathbb{R}^{d \times d}, h, x \in \mathbb{R}^{d \times 1}$$

$$h = (W + BA)x$$
$$B \in \mathbb{R}^{d \times r}, A \in \mathbb{R}^{r \times d}$$

Both full-finetuning and PEFT can be used in any of the post-training methods.
PEFT like Lora saves memory, learns less while forgets less [1]

[1] Biderman, Dan, Jacob Portes, Jose Javier Gonzalez Ortiz, Mansheej Paul, Philip Greengard, Connor Jennings, Daniel King et al. "Lora learns less and forgets less." *arXiv preprint arXiv:2405.09673* (2024).