

# **Building Generative AI Applications**

## **AN LLM CASE STUDY**

Michelle Yi  
ODSC 2023



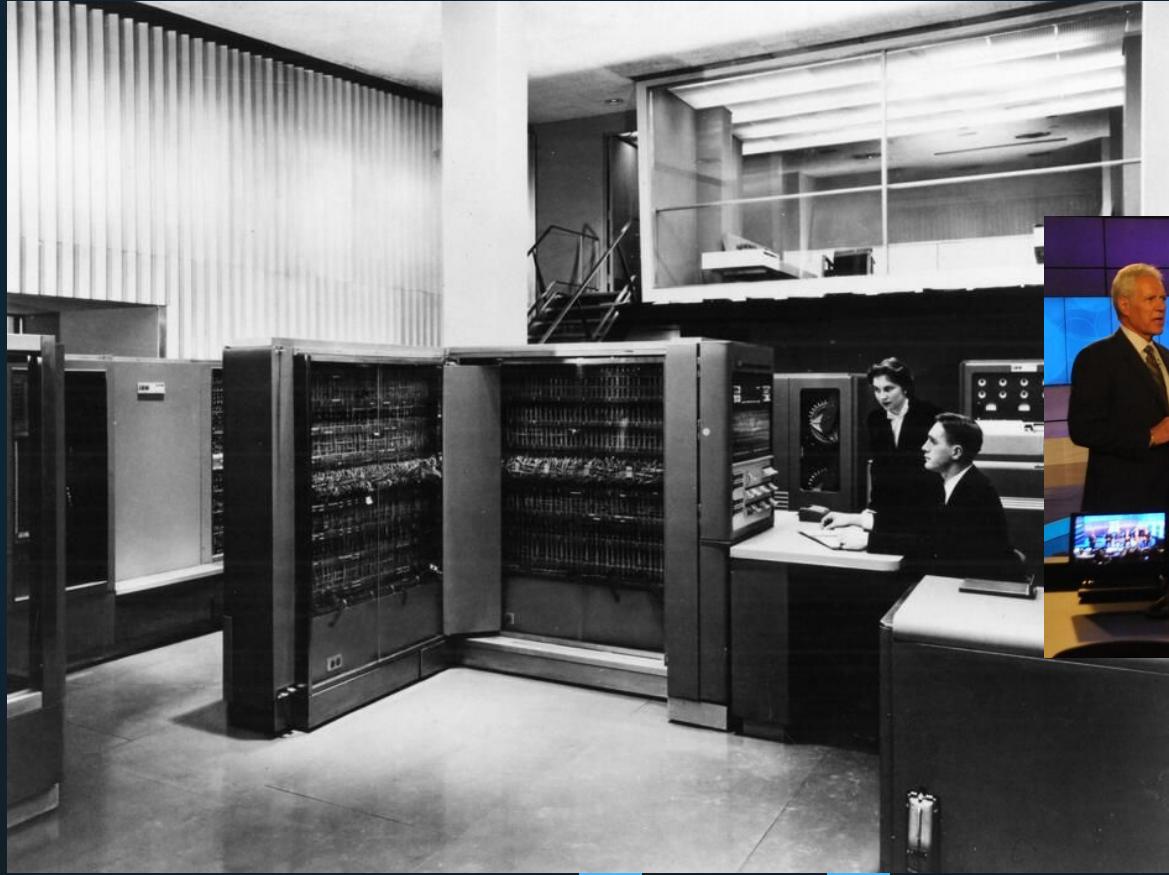
# INTRODUCTION

Michelle Yi | [Michelle@generationship.ai](mailto:Michelle@generationship.ai)

Applied AI, Artist, Speaker

Board Member, Women in Data





<https://www.techrepublic.com/article/ibm-watson-the-inside-story-of-how-the-jeopardy-winning-supercomputer-was-born-and-what-it-wants-to-do-next/>

## **DESIRED OUTCOME**

Take any LLM-related idea that you have and feel empowered to bootstrap it with open source tooling.

# **DESIRED OUTCOME**

1. A framework for building LLM-driven applications
2. An architecture of how to get started

A dark, moody photograph of a person's hands holding a black video game controller, likely a PlayStation DualShock 4. The background is blurred, showing what appears to be a television screen displaying a game. Overlaid on the image are several blue squares of varying sizes, some aligned vertically on the left and others forming a grid pattern on the right.

HOW IT ALL  
STARTED

# HOW IT ALL STARTED



PALADIN



DRUID

# DUNGEONS AND DRAGONS



# DUNGEONS AND DRAGONS





# BASIC RULES

## Player's Basic Rules Version 0.2

### CREDITS

**D&D Lead Designers:** Mike Mearls, Jeremy Crawford

**Design Team:** Christopher Perkins, James Wyatt, Rodney Thompson, Robert J. Schwalb, Peter Lee, Steve Townshend, Bruce R. Cordell

**Editing Team:** Chris Sims, Michele Carter, Scott Fitzgerald Gray  
**Producer:** Greg Bilsland

**Art Directors:** Kate Irwin, Dan Gelon, Jon Schindehette, Mari Kolkowsky, Melissa Rapier, Shauna Narciso

**Graphic Designers:** Bree Heiss, Emi Tanji  
**Interior Illustrator:** Jaime Jones

**Additional Contributors:** Kim Mohan, Matt Sennett, Chris Dupuis, Tom LaPille, Richard Baker, Chris Tulach, Miranda Horner, Jennifer Clarke Wilkes, Steve Winter, Nina Hess

**Project Management:** Neil Shinkle, Kim Graham, John Hay  
**Production Services:** Cynda Callaway, Brian Dumas, Jefferson Dunlap, Anita Williams

Based on the original D&D game created by

E. Gary Gygax and Dave Arneson, with Brian Blume, Rob Kuntz, James Ward, and Don Kaye

Drawing from further development by

J. Eric Holmes, Tom Moldvay, Frank Mentzer, Aaron Allston, Harold Johnson, David "Zeb" Cook, Ed Greenwood, Keith Baker, Tracy Hickman, Margaret Weis, Douglas Niles, Jeff Grubb, Jonathan Tweet, Monte Cook, Skip Williams, Richard Baker, Peter Adkison, Bill Slavicsek, Andy Collins, and Rob Heinsoo

Playtesting provided by  
over 175,000 fans of D&D. Thank you!

Additional consultation provided by

Jeff Grubb, Kenneth Hite, Kevin Kulp, Robin Laws, S. John Ross, the RPGPundit, Vincent Venturella, and Zak S.

<https://dnd.wizards.com/>

### ABILITY SCORE SUMMARY

#### Strength

**Measures:** Natural athleticism, bodily power

**Important for:** Fighter

**Racial Increases:**

Mountain dwarf (+2)

Human (+1)

#### Dexterity

**Measures:** Physical agility, reflexes, balance, poise

**Important for:** Rogue

**Racial Increases:**

Elf (+2)

Stout halfling (+1)

Human (+1)

#### Constitution

**Measures:** Health, stamina, vital force

**Important for:** Everyone

**Racial Increases:**

Dwarf (+2)

Stout halfling (+1)

Human (+1)

This method of determining ability scores enables you to create a set of three high numbers and three low ones (15, 15, 15, 8, 8), a set of numbers that are above average and nearly equal (13, 13, 13, 12, 12, 12), or any set of numbers between those extremes.

### 4. DESCRIBE YOUR CHARACTER

#### ABILITY SCORE POINT COST

Score	Cost	Score	Cost
8	0	12	4
9	1	13	5
10	2	14	7
11	3	15	9

Once you know the basic game aspects of your character, it's time to flesh him or her out as a person. Your character needs a name. Spend a few minutes thinking about what he or she looks like and how he or she behaves in general terms.

Using the information in chapter 4, you can flesh out your character's physical appearance and personality traits. Choose your character's **alignment** (the moral

## MONSTERS

Guidelines for understanding the information found in a monster's statistics are presented below.

### STATISTICS

A monster's statistics, sometimes referred to as its **stat block**, provide the essential information that you need to run the monster.

#### SIZE

A monster can be Tiny, Small, Medium, Large, Huge, or Gargantuan. The Size Categories table shows how much space a creature of a particular size controls in combat. See the player's D&D basic rules or the *Player's Handbook* for more information on creature size and space.

#### SIZE CATEGORIES

Size	Space	Examples
Tiny	2½ by 2½ ft.	Imp, sprite
Small	5 by 5 ft.	Giant rat, goblin
Medium	5 by 5 ft.	Orc, werewolf
Large	10 by 10 ft.	Hippogriff, ogre
Huge	15 by 15 ft.	Fire giant, treant
Gargantuan	20 by 20 ft. or larger	Kraken, purple worm

iconic constructs. Many creatures native to the outer planes of Mechanus, such as modrons, are constructs shaped from the raw material of the plane by the will of more powerful creatures.

**Dragons** are large reptilian creatures of ancient origin and tremendous power. True dragons, including the good metallic dragons and the evil chromatic dragons, are highly intelligent and have innate magic. Also in this category are creatures distantly related to true dragons, but less powerful, less intelligent, and less magical, such as wyverns and pseudodragons.

**Elementals** are creatures native to the elemental planes. Some creatures of this type are little more than animate masses of their respective elements, including the creatures simply called elementals. Others have biological forms infused with elemental energy. The races of genies, including djinn and efreet, form the most important civilizations on the elemental planes. Other elemental creatures include azers, invisible stalkers, and water weirds.

**Fey** are magical creatures closely tied to the forces of nature. They dwell in twilight groves and misty forests. In some worlds, they are closely tied to the Feywild, also called the Plane of Faerie. Some are also found in the Outer Planes, particularly the planes of Arborea and the Beastlands. Fey include dryads, pixies, and satyrs.

# Campaign

Name: \_\_\_\_\_  
 City: \_\_\_\_\_  
 Location: \_\_\_\_\_  
 Date: \_\_\_\_\_  
 Notes: \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_

Name: \_\_\_\_\_  
 City: \_\_\_\_\_  
 Location: \_\_\_\_\_  
 Date: \_\_\_\_\_  
 Notes: \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_

Name: \_\_\_\_\_  
 City: \_\_\_\_\_  
 Location: \_\_\_\_\_  
 Date: \_\_\_\_\_  
 Notes: \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_

Name: \_\_\_\_\_  
 City: \_\_\_\_\_  
 Location: \_\_\_\_\_  
 Date: \_\_\_\_\_  
 Notes: \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_  
 \_\_\_\_\_

\_\_\_\_\_

	Level 1	Level 2	Level 3	Level 4	Level 5	Level 6	Level 7	Level 8	Level 9	Level 10
<b>A - Required</b>										
Resources from Head Business	\$ 170,000	\$ 401,000	\$ 711,000	\$ 1,031,000	\$ 1,351,000	\$ 1,671,000	\$ 2,000,000	\$ 2,320,000	\$ 2,640,000	\$ 2,960,000
Business Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Other Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
<b>B - Recommended</b>										
Business Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Software Expenses	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000	\$ 110,000
Software Support	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Office Assets	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
<b>C - Operating Expenses</b>										
Business Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Software Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Marketing and Advertising Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Equipment and Infrastructure Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Office Assets	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Travel Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Entertainment Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Food and Beverage Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Professional Services Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Personal Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
<b>D - Marketing and Advertising Expenses</b>										
Business Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Software Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Marketing and Advertising Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Equipment and Infrastructure Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Office Assets	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Travel Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Entertainment Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Food and Beverage Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Professional Services Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Personal Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
<b>E - Marketing and Advertising Expenses</b>										
Business Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Software Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Marketing and Advertising Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Equipment and Infrastructure Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Office Assets	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Travel Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Entertainment Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Food and Beverage Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Professional Services Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000
Personal Expenses	\$ 10,000	\$ 20,000	\$ 30,000	\$ 40,000	\$ 50,000	\$ 60,000	\$ 70,000	\$ 80,000	\$ 90,000	\$ 100,000

Part 5

#### Q5 General Questions

what would the action sequence be like for sheathing one weapon, drawing another, and attacking? If so, action? yeah - the intent is to avoid punishing players for that stuff by changing an action - M

How does the term "per turn" work? Could a PC who can make an attack on everyone else's turn? yes, though you'd still need to use your action or reaction. M

Is being "attacked" ranged? So you give away your location, does this cancel "hiding"? Doesn't say you are "seen" or "heard". Yes, others are aware of you. M

Can a 1st-level character ready an action, then take that action if triggered by a high-level character in the next round? yes. M

Can an action be a bonus action? Can a character use a bonus action to draw? "Drawing" doesn't mean "draw" in the sense of the word? Actions to draw aren't necessarily drawn, so that kind could use bonus action to draw, or free hand, if not both are a turn. M

#### Hiding Thing

Does a character that stealths with a meal roll 50% or 60%? There was a question because the Stealth ability says one die, just one die - great news! Y'all - M

Do characters can only use 2nd to attack one creature since bonus action reflects, or can attack 2 diff enemies with 2nd? Can attack two targets. -M



# 01

## CREATING AN ASSISTANT DUNGEON MASTER

# EXISTING SOLUTIONS

CHATGPT

BARD

LLAMA (2)

- Dexterity: 15 (+2 racial bonus for Wood Elf = 17)
- Wisdom: 12
- Constitution: 12
- Strength: 10

Features:

**Hallucination**

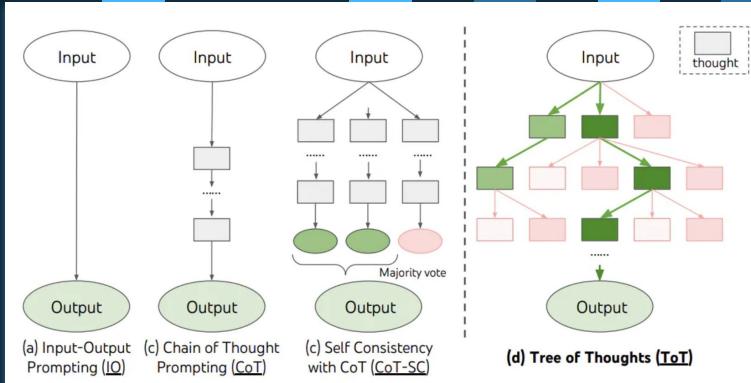
- Martial Arts
- Flurry of Blows

“I can see the world  
as it should be.”

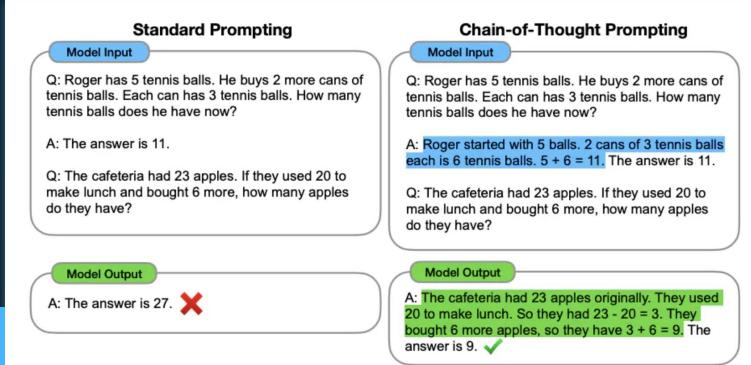
# THEN WE DO PROMPT ENGINEERING

We try to hack commercial LLMS by using:

- Chain-of-thought prompting
- Few-shot prompting
- Tree of thoughts
- Retrieval augmented generation (RAG)
- Other trial and error



<https://arxiv.org/abs/2305.10601>



<https://arxiv.org/abs/2201.11903>

# THE PROBLEM



Prompt engineering and trying to solve for limited context windows may be the second step and not the first!

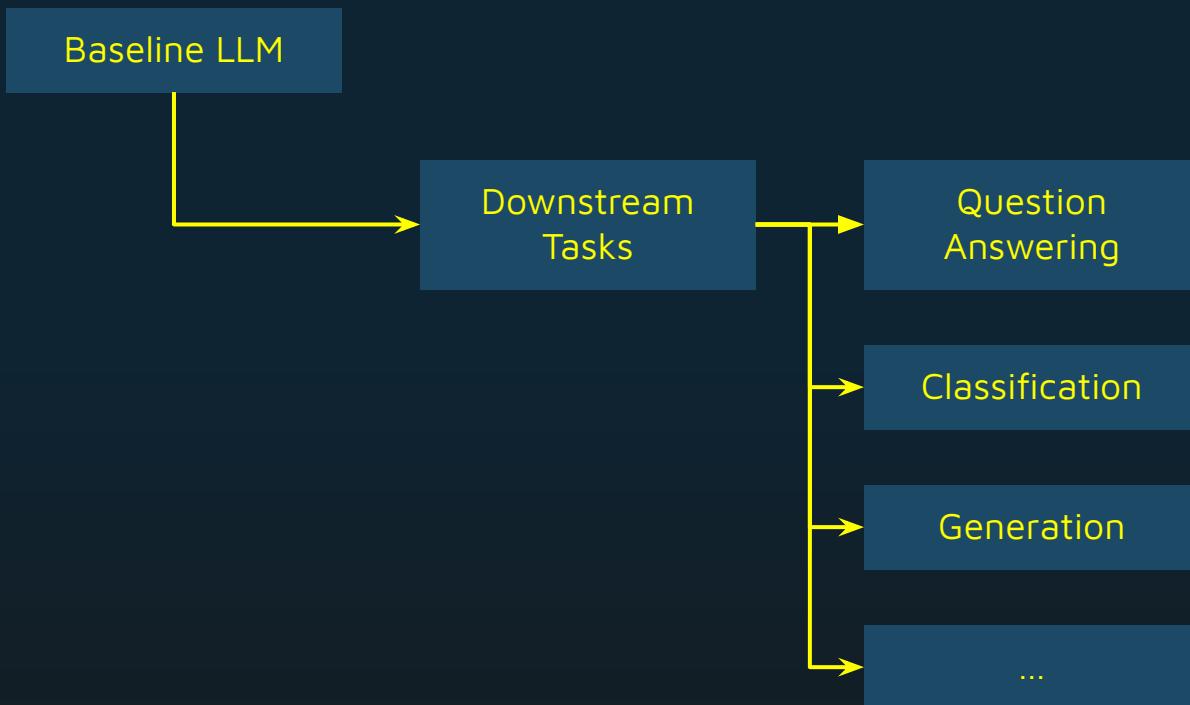
**We can instead ask: How can we give LLMs better domain knowledge?**



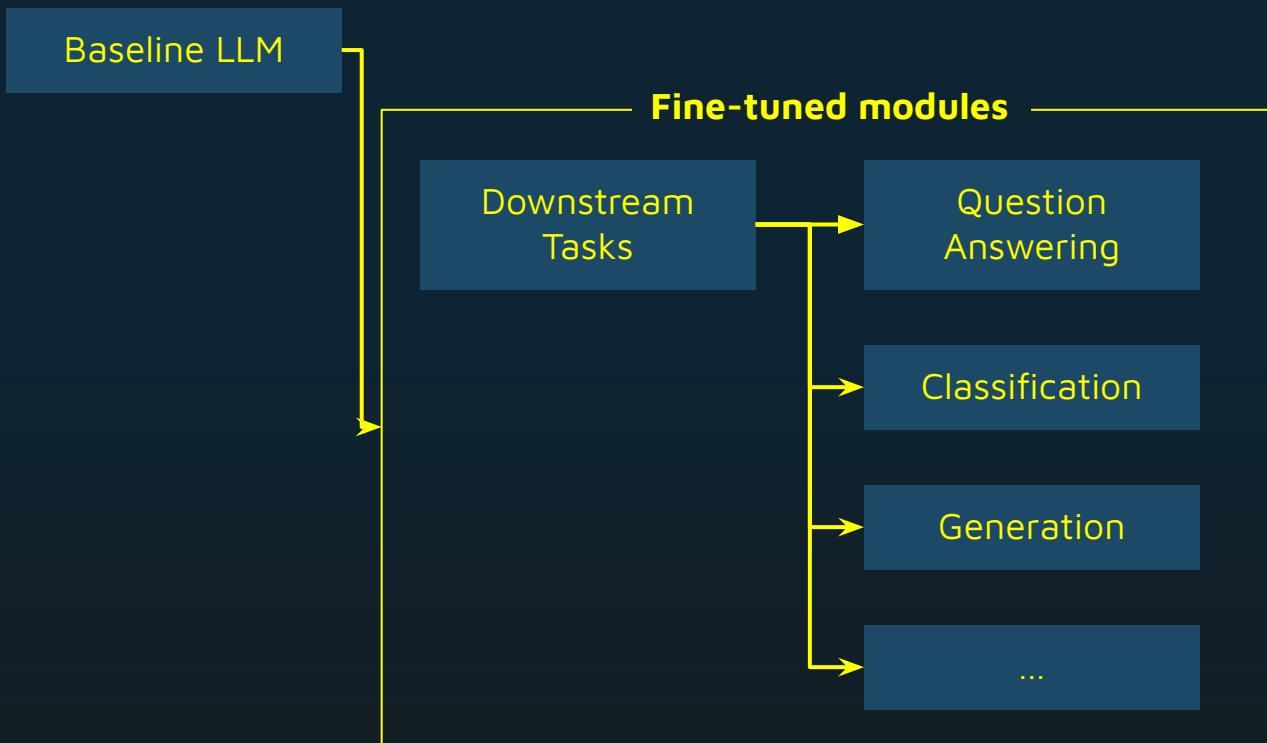
02

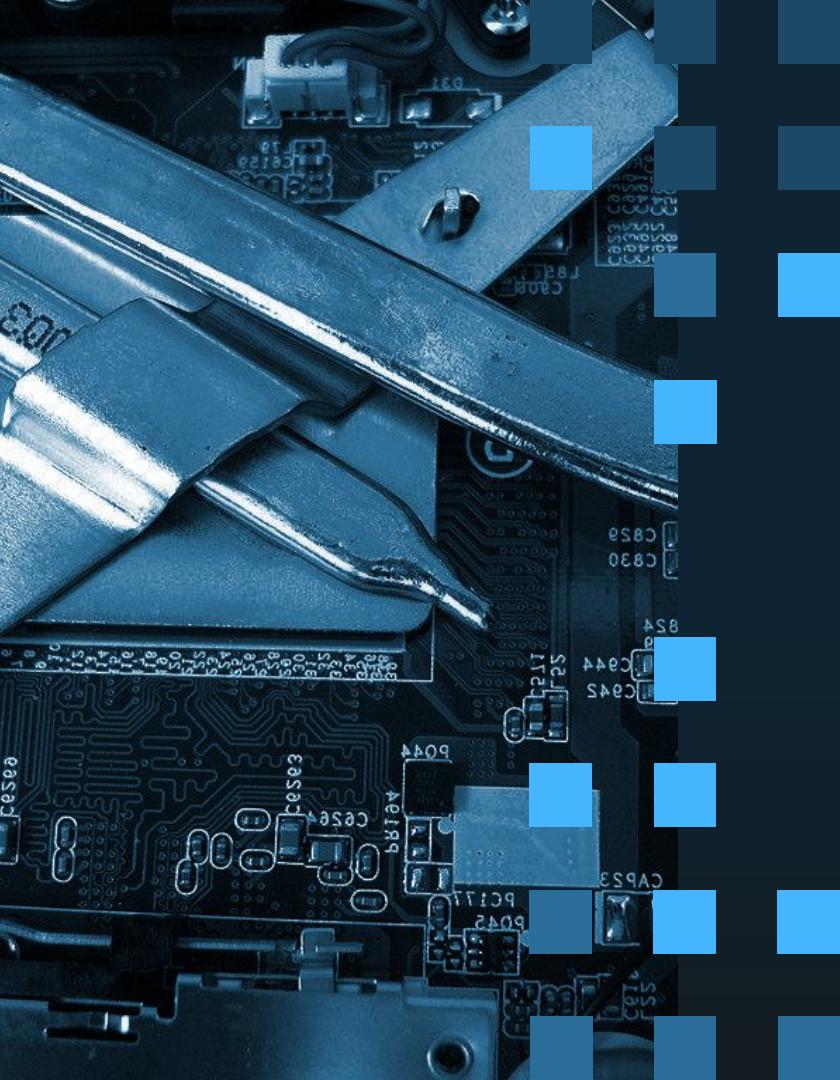
# FINE-TUNING A MODEL

# WHERE TO START



# WHERE TO START





**BUT WAIT, WHAT  
ABOUT GPUs?**

Don't I need a thousand of these?

# LOW-RANK ADAPTATION OF LLMS

## LORA: LOW-RANK ADAPTATION OF LARGE LANGUAGE MODELS

Edward Hu\* Yelong Shen\* Phillip Wallis Zeyuan Allen-Zhu  
Yuanzhi Li Shean Wang Lu Wang Weizhu Chen

Microsoft Corporation

{edwardhu, yeshe, phwallis, zeyuana,  
yuanzhil, swang, luw, wzchen}@microsoft.com  
yuanzhil@andrew.cmu.edu  
(Version 2)

### ABSTRACT

An important paradigm of natural language processing consists of large-scale pre-training on general domain data and adaptation to particular tasks or domains. As we pre-train larger models, full fine-tuning, which retrains all model parameters, becomes less feasible. Using GPT-3 175B as an example – deploying independent instances of fine-tuned models, each with 175B parameters, is prohibitively expensive. We propose Low-Rank Adaptation, or LoRA, which freezes the pre-trained model weights and injects trainable rank decomposition matrices into each layer of the Transformer architecture, greatly reducing the number of trainable parameters for downstream tasks. Compared to GPT-3 175B fine-tuned with Adam, LoRA can reduce the number of trainable parameters by 10,000 times and the GPU memory requirement by 3 times. LoRA performs on-par or better than fine-tuning in model quality on RoBERTa, DeBERTa, GPT-2, and GPT-3, despite having fewer trainable parameters, a higher training throughput, and, unlike adapters, no additional inference latency. We also provide an empirical investigation into rank-deficiency in language model adaptation, which sheds light on the efficacy of LoRA. We release a package that facilitates the integration of LoRA with PyTorch models and provide our implementations and model checkpoints for RoBERTa, DeBERTa, and GPT-2 at <https://github.com/microsoft/Lora>.

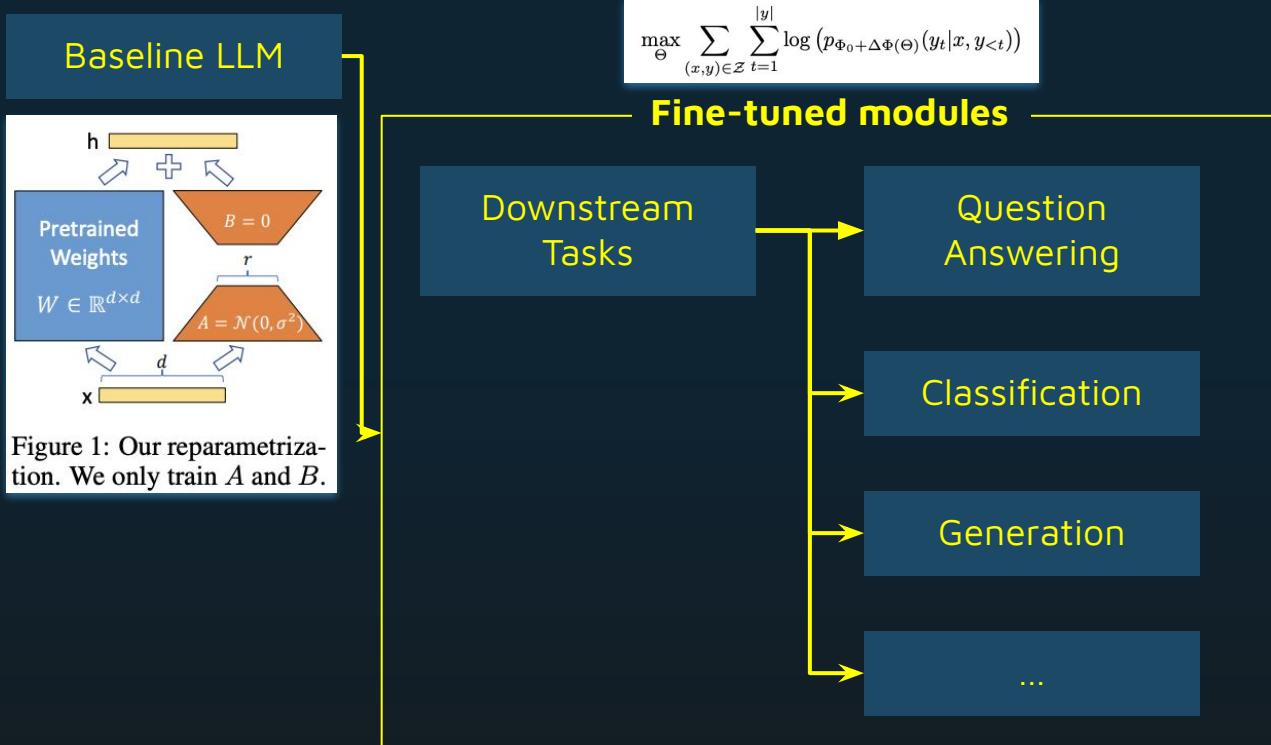
### 1 INTRODUCTION

Model & Method	# Trainable Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
RoB <sub>base</sub> (FT)*	125.0M	<b>87.6</b>	94.8	90.2	<b>63.6</b>	92.8	<b>91.9</b>	78.7	91.2	86.4
RoB <sub>base</sub> (BitFit)*	0.1M	84.7	93.7	<b>92.7</b>	62.0	91.8	84.0	81.5	90.8	85.2
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.3M	87.1 <sub>±.0</sub>	94.2 <sub>±.1</sub>	88.5 <sub>±1.1</sub>	60.8 <sub>±.4</sub>	93.1 <sub>±.1</sub>	90.2 <sub>±.0</sub>	71.5 <sub>±2.7</sub>	89.7 <sub>±.3</sub>	84.4
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.9M	87.3 <sub>±.1</sub>	94.7 <sub>±.3</sub>	88.4 <sub>±.1</sub>	62.6 <sub>±.9</sub>	93.0 <sub>±.2</sub>	90.6 <sub>±.0</sub>	75.9 <sub>±2.2</sub>	90.3 <sub>±.1</sub>	85.4
RoB <sub>base</sub> (LoRA)	0.3M	87.5 <sub>±.3</sub>	<b>95.1<sub>±.2</sub></b>	89.7 <sub>±.7</sub>	63.4 <sub>±1.2</sub>	<b>93.3<sub>±.3</sub></b>	90.8 <sub>±.1</sub>	<b>86.6<sub>±.7</sub></b>	<b>91.5<sub>±.2</sub></b>	<b>87.2</b>
RoB <sub>large</sub> (FT)*	355.0M	90.2	<b>96.4</b>	<b>90.9</b>	68.0	94.7	<b>92.2</b>	86.6	92.4	88.9
RoB <sub>large</sub> (LoRA)†	0.8M	90.3 <sub>±.3</sub>	96.5 <sub>±.5</sub>	87.7 <sub>±1.7</sub>	66.3 <sub>±2.0</sub>	94.7 <sub>±.2</sub>	91.5 <sub>±.1</sub>	72.9 <sub>±2.9</sub>	91.5 <sub>±.5</sub>	86.4
RoB <sub>large</sub> (Adpt <sup>D</sup> )†	0.8M	<b>90.6<sub>±.2</sub></b>	96.2 <sub>±.5</sub>	<b>90.2<sub>±1.0</sub></b>	68.2 <sub>±1.9</sub>	<b>94.8<sub>±.3</sub></b>	91.6 <sub>±.2</sub>	<b>85.2<sub>±1.1</sub></b>	<b>92.3<sub>±.5</sub></b>	<b>88.6</b>
DeB <sub>XXL</sub> (FT)*	1500.0M	91.8	<b>97.2</b>	92.0	72.0	<b>96.0</b>	92.7	93.9	92.9	91.1
DeB <sub>XXL</sub> (LoRA)	4.7M	<b>91.9<sub>±.2</sub></b>	96.9 <sub>±.2</sub>	<b>92.6<sub>±.6</sub></b>	<b>72.4<sub>±1.1</sub></b>	<b>96.0<sub>±.1</sub></b>	<b>92.9<sub>±.1</sub></b>	<b>94.9<sub>±.4</sub></b>	<b>93.0<sub>±.2</sub></b>	<b>91.3</b>

10,000x reduction in trainable parameters for GPT-3, without any loss in performance.

Table 2: RoBERTa<sub>base</sub>, RoBERTa<sub>large</sub>, and DeBERTa<sub>XXL</sub> with different adaptation methods on the GLUE benchmark. We report the overall (matched and mismatched) accuracy for MNLI, Matthew’s correlation for CoLA, Pearson correlation for STS-B, and accuracy for other tasks. Higher is better for all metrics. \* indicates numbers published in prior works. † indicates runs configured in a setup similar to Houlsby et al. (2019) for a fair comparison.

# WHAT LORA DOES



## WHY IT MATTERS

1. A pre-trained model can be shared and used to build many small LoRA modules for different tasks.
2. LoRA can boost training efficiency and lower hardware requirements up to three times by using adaptive optimizers that eliminate the need to calculate gradients or maintain optimizer states for most parameters.
3. No inference latency.
4. Interoperable with other techniques. After conducting a thorough analysis of previous research in this field, the authors discovered that this approach could be integrated with other adapter or quantization methods.

# EASY TO USE OPEN-SOURCE

[← Back to blog](#)

## 😊 PEFT: Parameter-Efficient Fine-Tuning of Billion-Scale Models on Low-Resource Hardware

Published February 10, 2023

[Update on GitHub](#)



[smangrul](#)  
Sourab Mangrulkar



[sayakpaul](#)  
Sayak Paul

### Motivation

Large Language Models (LLMs) based on the transformer architecture, like GPT, T5, and BERT

<https://huggingface.co/blog/peft>

The screenshot shows a GitHub repository named "LoRA" with the "Public" badge. It displays the main branch with 29 branches and 4 tags. A commit by "edwardjhu" titled "Update setup.py" is shown, along with other commits related to examples, loralib, .gitignore, LICENSE.md, README.md, SECURITY.md, and setup.py. The repository has 3f5c193 on Aug 2.

**LoRA: Low-Rank Adaptation of Large Language Models**

<https://github.com/microsoft/LoRA>

# STEP 1: DATA

1. Select ethically and responsibly (no proprietary D&D books were used for training)
2. Engineer into the right format (**HINT: Use an LLM!**)
3. Human expert review (Dungeon Master)

## Sample JSON outputs:

```
[{"instruction": "What is the name of the game mentioned in the text? ",  
 "input": "", "output": "Dungeons and Dragons."},  
 {"instruction": "What is the meaning of the numbers and signs listed after the six ability scores? ",  
 "input": "", "output": "They show the score and the modifier for each ability."},  
 {"instruction": "What are the skills listed under the behir's stat block? ",  
 "input": "",  
 "output": "Perception and Stealth."},....]
```

# STEP 2: MODEL SELECTION

Models 10,614

- CausalLM/14B**  
Text Generation • Updated about 9 hours ago • ↓ 384 • ❤ 169
- meta-llama/Llama-2-7b-chat-hf**  
Text Generation • Updated 6 days ago • ↓ 947k • ❤ 1.57k
- CausalLM/7B**  
Text Generation • Updated about 9 hours ago • ↓ 326 • ❤ 83
- meta-llama/Llama-2-7b**  
Text Generation • Updated Jul 19 • ❤ 2.88k
- TheBloke/CausalLM-14B-GGUF**  
Text Generation • Updated 5 days ago • ↓ 36 • ❤ 39
- THUDM/agentlm-70b**  
Text Generation • Updated 9 days ago • ↓ 171 • ❤ 50
- oobabooga/CodeBooga-34B-v0.1**  
Text Generation • Updated 6 days ago • ↓ 110 • ❤ 49
- meta-llama/Llama-2-7b-hf**  
Text Generation • Updated Aug 9 • ↓ 507k • ❤ 739

**Start small** - for this example we start with Llama 7B open source, which we can use LoRA to fit into a single free-tier Colab GPU

# STEP 3: FINE-TUNING

1. Set-up (Colab in this example)
2. Size the training job
3. Read Llama model (in this example: 7B-hf)
4. Configure LoRA parameters

```
1 import os
import sys

import torch
import torch.nn as nn
import bitsandbytes as bnb
from datasets import load_dataset
import transformers
from transformers import LlamaForCausalLM, LlamaTokenizer
from peft import (
    prepare_model_for_int8_training,
    LoraConfig,
    get_peft_model,
    get_peft_model_state_dict,
)
```

```
2 # Set random seed for reproducibility
RANDOM_SEED = 1234
transformers.set_seed(RANDOM_SEED)

# Fix into Kaggle T4+2
MICRO_BATCH_SIZE = 8
BATCH_SIZE = 128
GRADIENT_ACCUMULATION_STEPS = BATCH_SIZE // MICRO_BATCH_SIZE
EPOCHS = 3 # One epoch takes ~6 hours, and 2 epochs may exceed Kaggle 12-hour limit
LEARNING_RATE = 3e-4 # Following stanford_alpaca
CUTOFF_LEN = 512 # 256 accounts for about 96% of the data. Shorter input, faster training/less VRAM
LORA_R = 16 # Some LoRA parameters
LORA_ALPHA = 10
LORA_DROPOUT = 0.05
VAL_SET_SIZE = 2000
TARGET_MODULES = [
    'q_proj',
    'v_proj',
    'k_proj',
    'o_proj'
]
```

```
3 DATA_PATH = '# LoRA config.
OUTPUT_DIR =
config = LoraConfig(
    r=LORA_R,
    lora_alpha=LORA_ALPHA,
    target_modules=TARGET_MODULES,
    lora_dropout=LORA_DROPOUT,
    bias='none',
    task_type='CAUSAL_LM',
```

1

2

3

4

# STEP 3: FINE-TUNING, Cont'd.

1. Generate and tokenize prompts
2. Training/evaluation data split
3. Train the model & peft setup
4. Save the fine-tuned model

```
def generate_and_tokenize_prompt(data_point):
    """This function masks out the labels for the input, so that our loss is computed on the response."""
    if data_point['input']:
        user_prompt = 'Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.\n\n'
        user_prompt += f'### Instruction:{data_point["instruction"]}\n\n'
        user_prompt += f'### Input:{data_point["input"]}\n\n'
        user_prompt += f'### Response:\n'
    else:
        user_prompt = 'Below is an instruction that describes a task. Write a response that appropriately completes the request.'
        user_prompt += f'### Instruction:{data_point["instruction"]}\n\n'
        user_prompt += f'### Response:\n'

# Count the length of prompt tokens
len_user_prompt_tokens = len(tokenizer(user_prompt,
                                        truncation=True,
                                        max_length=CUTOFF_LEN + 1,
                                        padding='max_length')['input_ids'])

len_user_prompt_tokens -= 1 # Minus 1 (one) for eos token

# Tokenize the input, both prompt and output
full_tokens = tokenizer(
    user_prompt,
    data_point['output'],
    truncation=True,
    max_length=CUTOFF_LEN + 1,
    padding='max_length',
)
if 'input_ids' in full_tokens:
    return {
        'input_ids': full_tokens['input_ids'],
        'labels': [-100] * len_user_prompt_tokens + full_tokens['input_ids'][len_user_prompt_tokens:],
        'attention_mask': [1] * (len(full_tokens)),
    }
else:
    train_data = data['train'].map(generate_and_tokenize_prompt)
    val_data = None

# Train/val split
if VAL_SET_SIZE > 0:
    train_val = data['train'].train_test_split(
        test_size=VAL_SET_SIZE, shuffle=False, seed=RANDOM_SEED
    )
    train_data = train_val['train'].map(generate_and_tokenize_prompt)
    val_data = train_val['test'].map(generate_and_tokenize_prompt)
else:
    train_data = data['train'].map(generate_and_tokenize_prompt)
    val_data = None

# HuggingFace Trainer
trainer = transformers.Trainer(
    model=model,
    train_dataset=train_data,
    eval_dataset=val_data,
    args=transformers.TrainingArguments(
        seed=RANDOM_SEED, # Reproducibility
        data_seed=RANDOM_SEED,
        per_device_train_batch_size=MICRO_BATCH_SIZE,
        gradient_accumulation_steps=GRADIENT_ACCUMULATION_STEPS,
        warmup_steps=100,
        num_train_epochs=EPOCHS,
        learning_rate=LEARNING_RATE,
        fp16=True,
        logging_steps=20,
        eval_steps=steps if VAL_SET_SIZE > 0 else 'no',
        save_steps=steps,
        save_steps=50,
        eval_steps=50 if VAL_SET_SIZE > 0 else None,
        output_dir=OUTPUT_DIR,
        save_total_limit=3,
        load_best_model_at_end=True if VAL_SET_SIZE > 0 else False,
        ddp_find_unused_parameters=False if ddp else None,
    ),
    data_collator=transformers.DataCollatorForLanguageModeling(tokenizer, mlm=False),
)
model.config.use_cache = False

# PEFT setup
old_state_dict = model.state_dict
model.state_dict = (
    lambda self, *, **kwargs: get_peft_model_state_dict(self, old_state_dict)
).__get__(model, type(model))
```

# STEP 4: INFERENCE

1. Load model
2. Create prompt
3. Generate response
4. Formatting output (tokenizer)
5. Conversation storage/retrieval (vector DB, llama-index)

```
"""## Create Prompt"""
prompt_base = f"""
Answer the generic questions without any help from the context.
But your goals are to give answer based on the relevant information given to you.

## READ FROM HERE FOR THE CONTEXT:
<<CONTEXT>>

## PREVIOUS CONVERSATION:
<<CONVERSATION>>

## Instruction: <<QUESTION>>
Now answer here as:
### Response:
"""

PROMPT_TEMPLATE = f"""
Below is an instruction that describes a task. Write a response that appropriately completes the r
"""

# ## Instruction: <<QUESTION>>

def generate_response(prompt: str, model: PeftModel) -> GreedySearchDecoderOnlyOutput:
    encoding = tokenizer(prompt, return_tensors="pt")
    input_ids = encoding["input_ids"].to(DEVICE)

    generation_config = GenerationConfig(
        temperature=0.1,
        top_p=0.75,
        repetition_penalty=1.1,
    )
    with torch.inference_mode():
        return model.generate(
            input_ids=input_ids,
            generation_config=generation_config,
            return_dict_in_generate=True,
            output_scores=True,
            max_new_tokens=256,
        )

def format_response(response: GreedySearchDecoderOnlyOutput) -> str:
    decoded_output = tokenizer.decode(response.sequences[0])
    # print(decoded_output)
    response = decoded_output.split("### Response:")[1].strip()
    return "\n".join(textwrap.wrap(response))

def ask_alpaca(prompt: str, model: PeftModel) -> str:
    prompt = create_prompt(prompt)
    response = generate_response(prompt, model)
    return format_response(response)
```

# STEP 5: TESTING

1. Asking out of distribution questions
2. Unit testing to make sure outputs do not change (depending on stability of model)
3. Iterate on data and fine-tuning based on results

# STEP 6: DEPLOYMENT OPTIONS

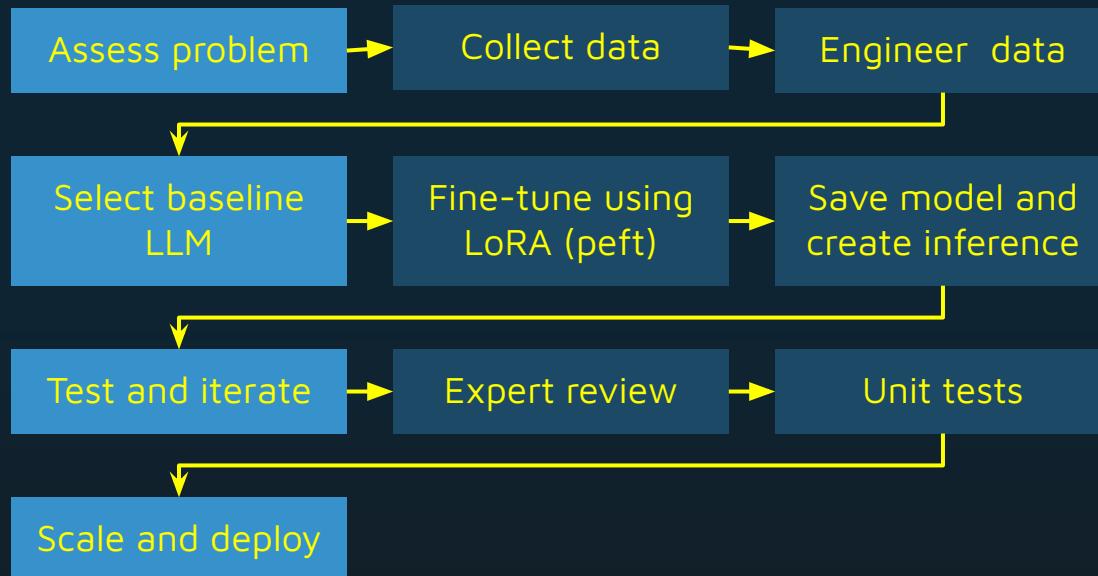
1. For fast iteration and testing, a notebook such as Colab is a good start.
2. When you have a front-end and more robust applications around it, consider kubernetes and containerizing your model to deploy: HTTPS, pub-sub, batch depending on what downstream task you need the module to do and latency requirements.
3. Consider adding a vector DB for quick caching, auditing, and longer-term information retrieval based on how the assistant DM helps shape the world ("Remember this...").



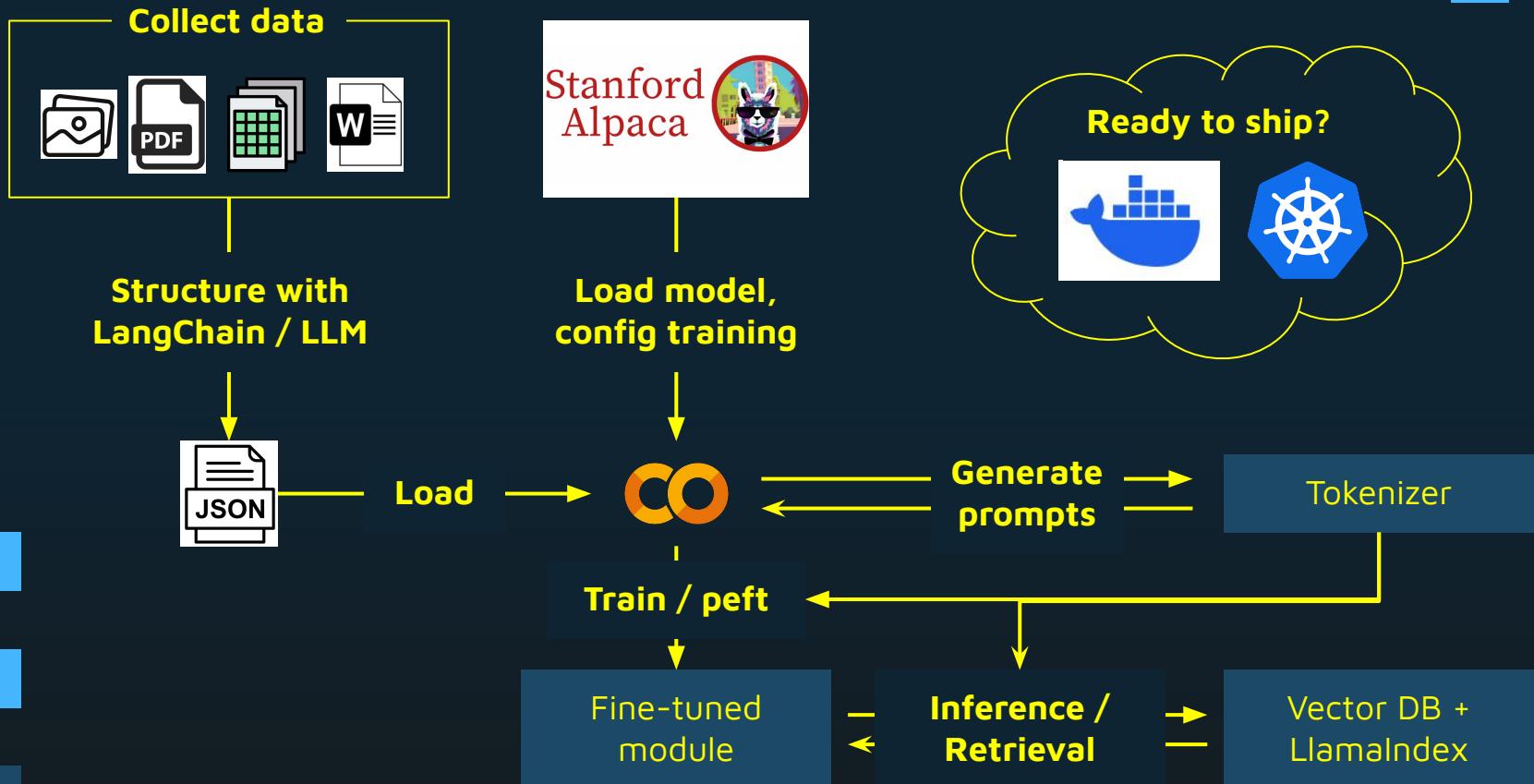
03

PUTTING  
IT ALL  
TOGETHER

# SUMMARY



# SUMMARY, Cont'd.



# SCUFFED DM IS BORN



So this is the story of both how “Scuffed DM” was born, but also how you can bootstrap your own fine-tuning projects.

**Coming soon to an 8-bit front-end near you.**

**But wait, why scuffed?**

# DOWNSTREAM TASK: INFORMATION RETRIEVAL/Q&A

USER: What are the powers of a monk in d&d?

BOT: A monk in D&D has several powers. They include unarmed strikes, monk weapons, ki points, monastic traditions, evasion, stillness of mind, timeless body, perfect self, empty body, and monastic tradition feature.

USER: What are the Damage Immunities of a monster in d&d?

BOT: Monsters in D&D have different kinds of damage immunities. Some monsters are immune to all forms of damage, while others are immune to certain types of damage. For example, a stone golem is immune to acid, fire, and cold damage, while a dragon is immune to acid, cold, and electricity damage.

USER: what is monster (E) in d&d?

BOT: Monster (E) in D&D refers to a creature that is not a player character. It could be a non-player character, a monster, a villain, or a creature that is part of the environment.

<>

≡

USER:

Executing (7m 32s) <cell line: 12> > raw\_in

# DOWNSTREAM TASK: GENERATION

USER: create a monk character based on base attributes

BOT: Name: Master Kung Fu Race/Class: Human/Monk Alignment: Lawful Good

Base Attributes: STR 12, DEX 14, CON 12, INT 12, WIS 12, CHA 12 Base

Skills: Acrobatics (Dex), Athletics (Str), History (Int), Insight  
(Wis), Religion (Wis), Stealth (Dex) Starting Equipment: Dagger  
(Simple Melee Weapon), Backpack (Explorer's Pack), 10 Darts

USER:



# DOWNSTREAM TASK: STEP-BY-STEP COMBAT



## Step 1: Initial Positioning

\* Monk: HP 120/120, AC 18, Attack Bonus +7

\* Paladin: HP 90/90, AC 16, Attack Bonus +5

## Step 2: Paladin Casts Bless

\* Monk: HP 120/120, AC 18, Attack Bonus +7

\* Paladin: HP 90/90, AC 16, Attack Bonus +5, bless active

## Step 3: Monk Dodges Paladin's Attack

\* Monk: HP 120/120, AC 18, Attack Bonus +7

\* Paladin: HP 90/90, AC 16, Attack Bonus +5

## Step 4: Monk Lands a Flurry of Punches and Kicks

\* Monk: HP 120/120, AC 18, Attack Bonus +7

\* Paladin: HP 87/90, AC 16, Attack Bonus (damaged by monk's flurry)

## Step 5: Paladin Casts Healing Word

\* Monk: HP 120/120, AC 18, Attack Bonus +7

\* Paladin: HP 90/90, AC 16, Attack Bonus +5 (healed by healing word)

## Step 6: Monk Lands a Devastating Finishing Blow

\* Monk: HP 120/120, AC 18, Attack Bonus +7

\* Paladin: HP 63/90, AC 16, Attack Bonus +5 (damaged by monk's finishing blow)



# DOWNSTREAM TASK: STEP-BY-STEP COMBAT



## Step 1: Initial Positioning

\* Monk: HP 120/120, AC 18, Attack Bonus +7

\* Paladin: HP 90/90, AC 16, Attack Bonus +5

## Step 2: Paladin Casts Bless

\* Monk: HP 120/120, AC 18, Attack Bonus +7

\* Paladin: HP 90/90, AC 16, Attack Bonus +5, bless active

## Step 3: Monk Dodges Paladin's Attack

\* Monk: HP 120/120, AC 18, Attack Bonus +7

\* Paladin: HP 90/90, AC 16, Attack Bonus +5

## Step 4: Monk Lands a Flurry of Punches and Kicks

\* Monk: HP 120/120, AC 18, Attack Bonus +7

\* Paladin: HP 87/90, AC 16, Attack Bonus (damaged by monk's flurry)

## Step 5: Paladin Casts Healing Word

\* Monk: HP 120/120, AC 18, Attack Bonus +7

\* Paladin: HP 90/90, AC 16, Attack Bonus +5 (healed by healing word)

## Step 6: Monk Lands a Devastating Finishing Blow

\* Monk: HP 120/120, AC 18, Attack Bonus +7

\* Paladin: HP 63/90, AC 16, Attack Bonus +5 (damaged by monk's finishing blow)



# WORKSHOP ON CAUSAL GRAPHS

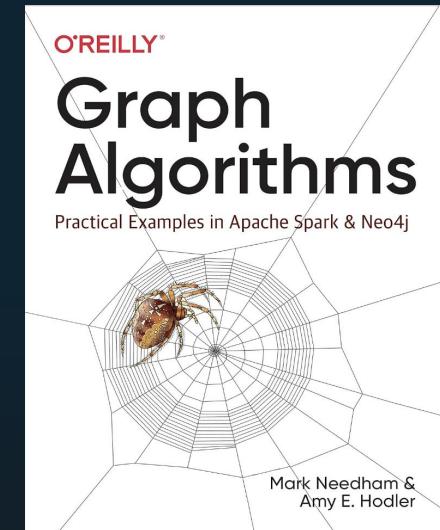
## GRAPHS

The Next Frontier of  
GenAI Explainability

Michelle Yi & Amy Hodler  
Fall 2023



3:45 PM Regency C



# THANKS

Do you have any questions or want to collaborate?

Reach out below:

michelle@generationship.ai

<https://www.linkedin.com/in/michelleyulleyi/>

**CREDITS:** This presentation template was created by **Slidesgo**, including icons by **Flaticon**, and infographics & images by **Freepik** and illustrations by **Stories**