# LLM Best Practises

ODSC Conference: Module 2

Oct 31, 2023

Sanyam Bhutani, Sr Data Scientist H2O.ai

H2O.ai

# Agenda

- Top Tricks from LLM Research

- [Case Study] Implementing these ideas

- Understanding Agents and AI Landscape

- [Hands-on] Creating LLM Apps and AI Agent Apps

# Quiz!

## https://tinyurl.com/ODSCLLMQ2

# Building blocks of LLMs

H2O.ai

**Why Large?**

○ Large Training Dataset: Trained on massive datasets
○ Large Architectures : Billions of parameters
○ Large Computing Power: Requires massive GPUs

## 01
### Foundation

Enormous amount of text data trained in an autoregressive manner

## 02
### Fine-tuning

Supervised fine-tuning on appropriate and well curated datasets to teach desired output behaviour.

## 03
### RLHF

Next token loss function replaced or combined with a reward model trained on Human Feedback.
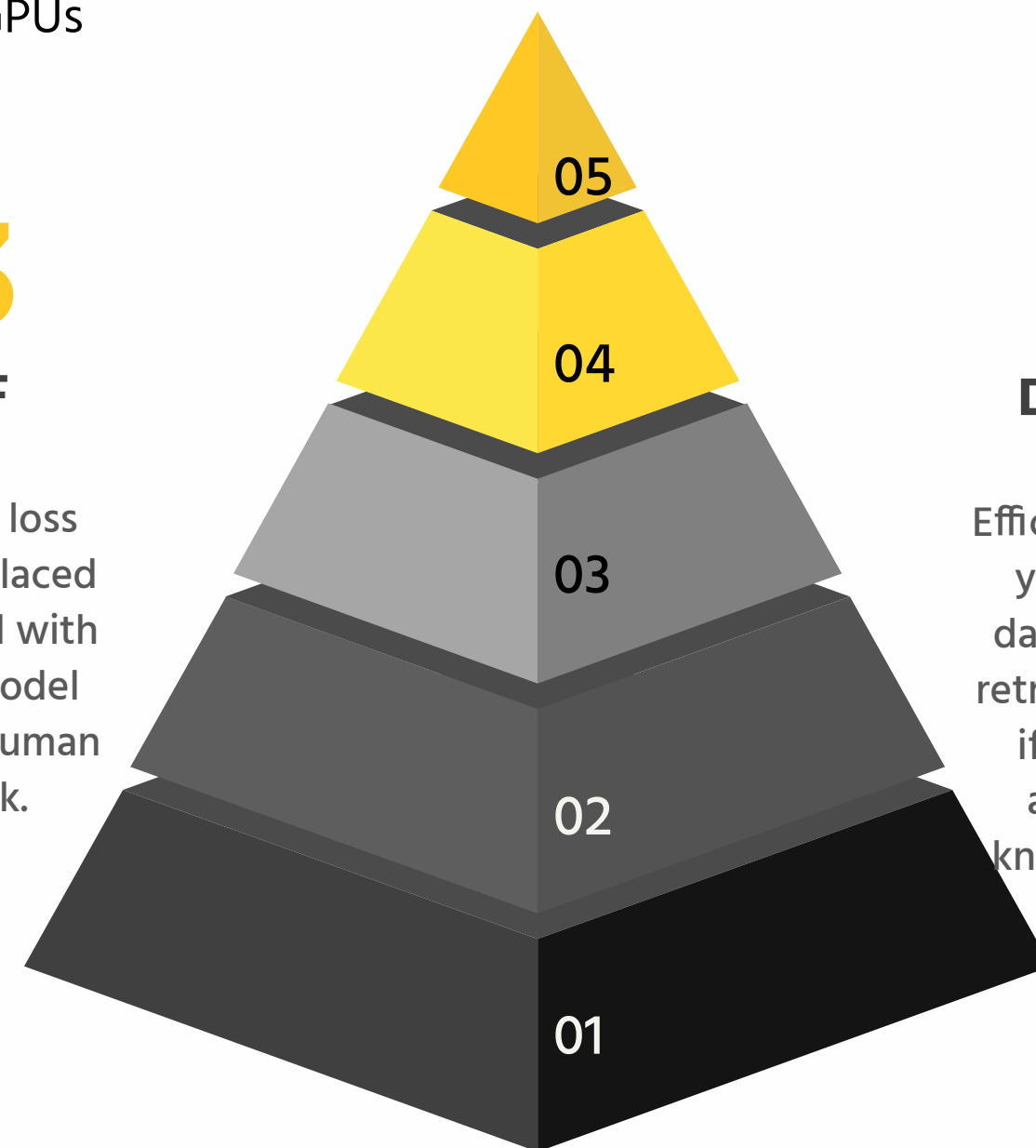
## 04
### Database

Efficiently leverage your company data. No need to retrain your model if a new pdf is added to the knowledge base.

## 05
### Memory

LLMs can have a huge context length and keep previous questions/tasks in memory for superior context understanding.

05

04

03

02

01

# LIMA: Less is More Alignment

- 1,000 carefully curated prompts and examples

- LLaMA-1 was fine-tuned on these to outperform all other models
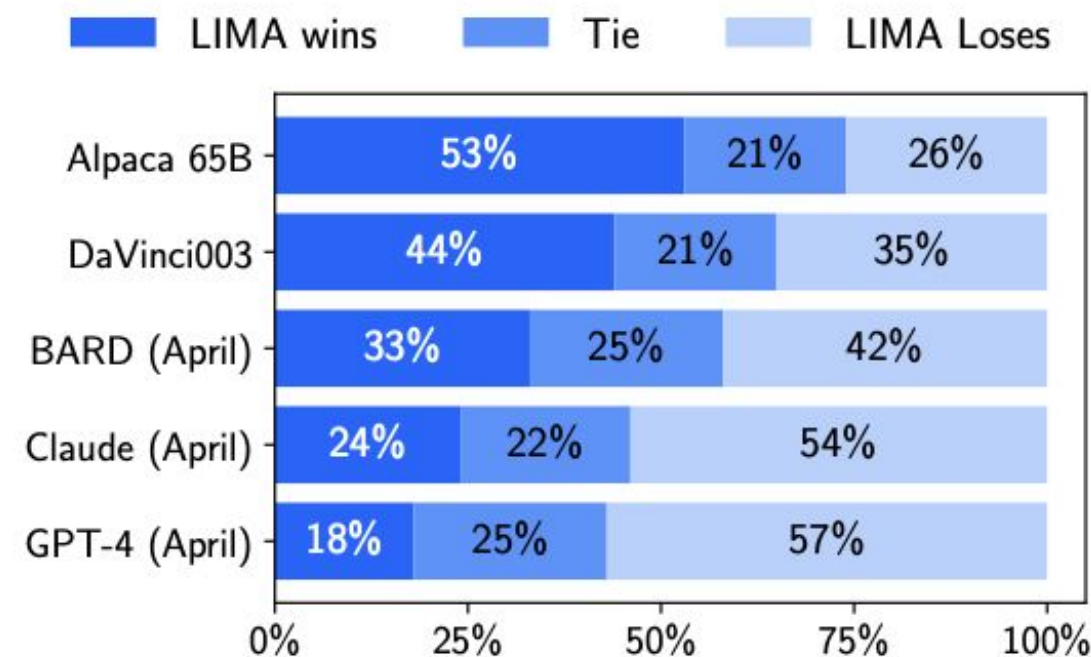
- Note: 65B model was used



Figure 1: Human preference evaluation, comparing LIMA to 5 different baselines across 300 test prompts.
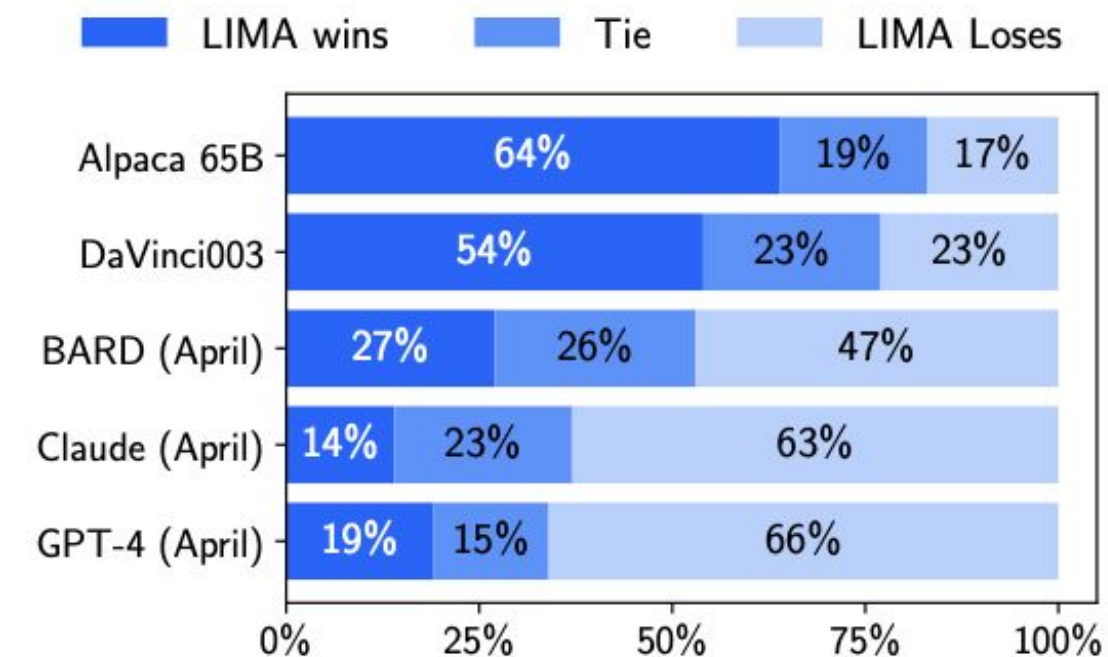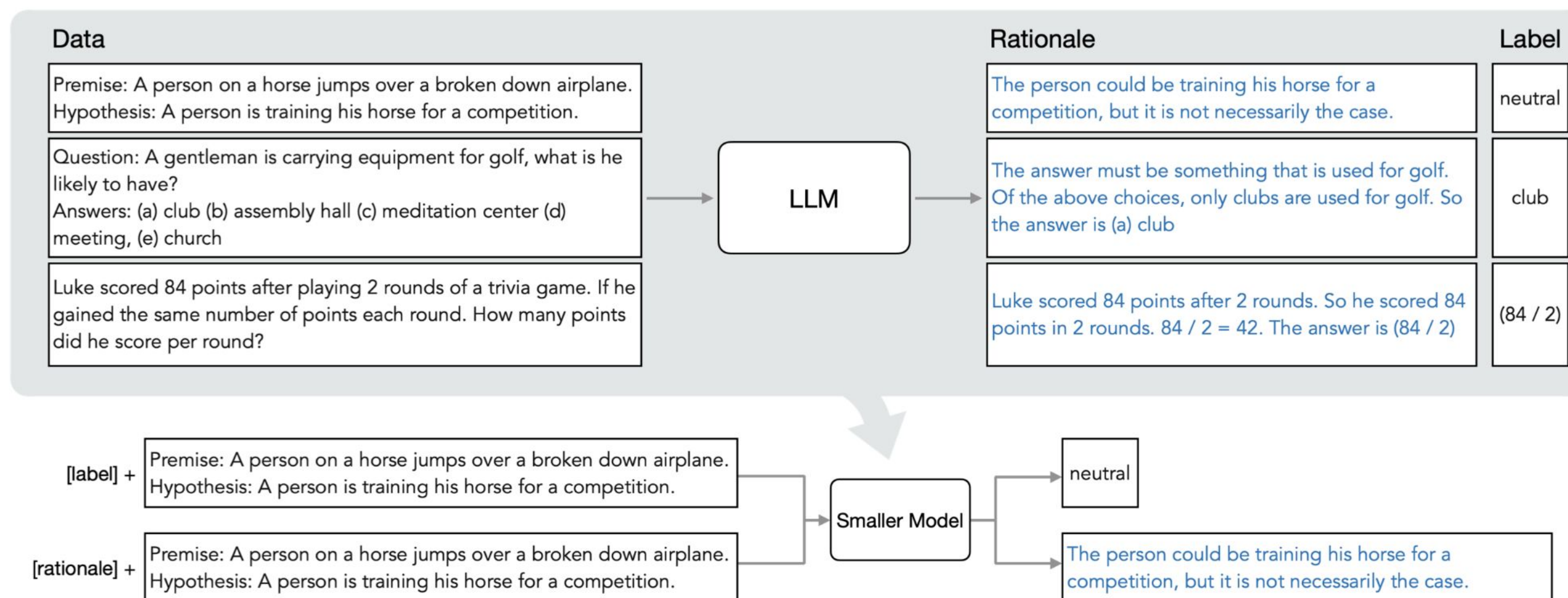
Figure 2: Preference evaluation using GPT-4 as the annotator, given the same instructions provided to humans.

# Distil: Step by Step

- Outperform 2000x Larger Models

- CoT to give logic to outputs and high quality tokens

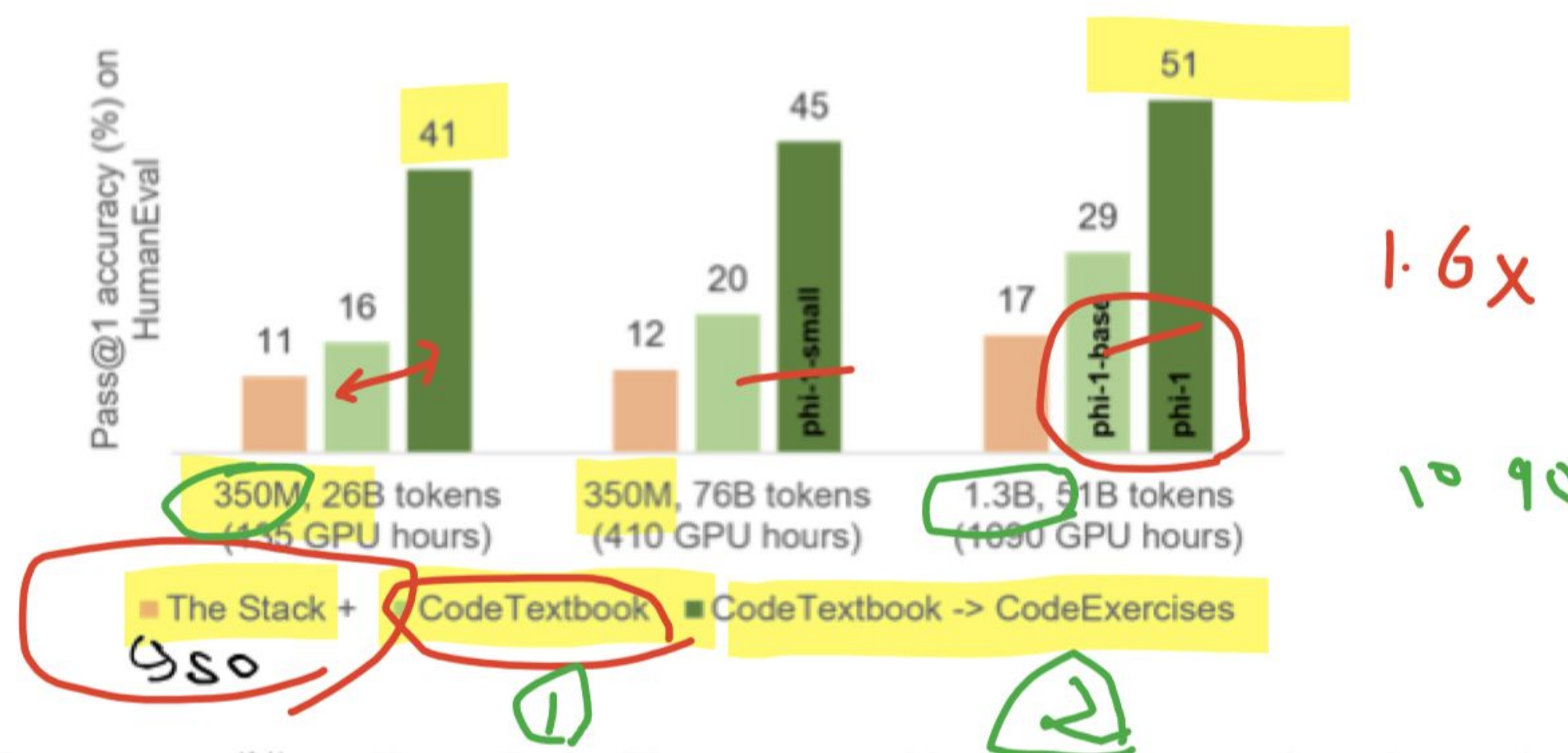- Outperforms both fine-tuned and distilled models

# Instruction BackTranslation

- Pseudo Labelling: Using Model to label data and perform SSL

- LLMs require to be converted to a "chatbot" where they are fine-tuned with chats

- This needs question-answer pairs

- We perform "backtranslation": LLaMA is used to create Qs from answers

- 3200 answers are enough to outperform everything else

# Textbooks are all you Need

- Smallest Model to generate Python Code

- Key: First Train on Task

- Later: Fine-Tune to questions

- The above step causes Emergent Abilities

# Quiz!

https://tinyurl.com/ODSCLLMQ3

# Demo in Practise

# Case Study: Simulacra Paper

# Hands-on: Implementing AI Agents

# Thank You

H2O.ai