

MODEL COMPRESSION: QUANTIZATION AND KNOWLEDGE DISTILLATION

MOHAMAD KHAJEZADE

SUMMER, 2024



IMPORTANCE OF SMALL MODELS



OpenAI 🌟 @OpenAI · Jul 18

We're continuing to make advanced AI accessible to all with the launch of GPT-4o mini, now available in the API and rolling out in ChatGPT today.



OpenAI Developers 🌟 @OpenAIDevs · Jul 18

Introducing GPT-4o mini! It's our most intelligent and affordable small model, available today in the API. GPT-4o mini is significantly smarter and cheaper than GPT-3.5 Turbo.

openai.com/index/gpt-4o-m...

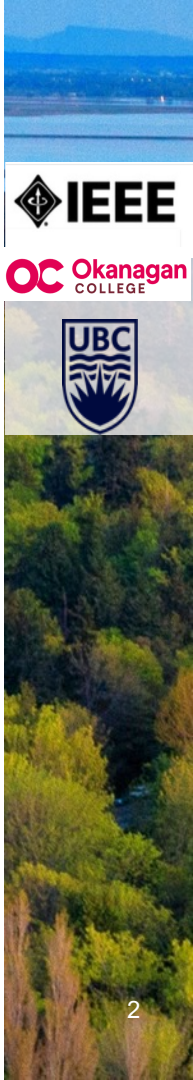
Pricing

GPT-4o mini is more than 60% cheaper than GPT-3.5 Turbo

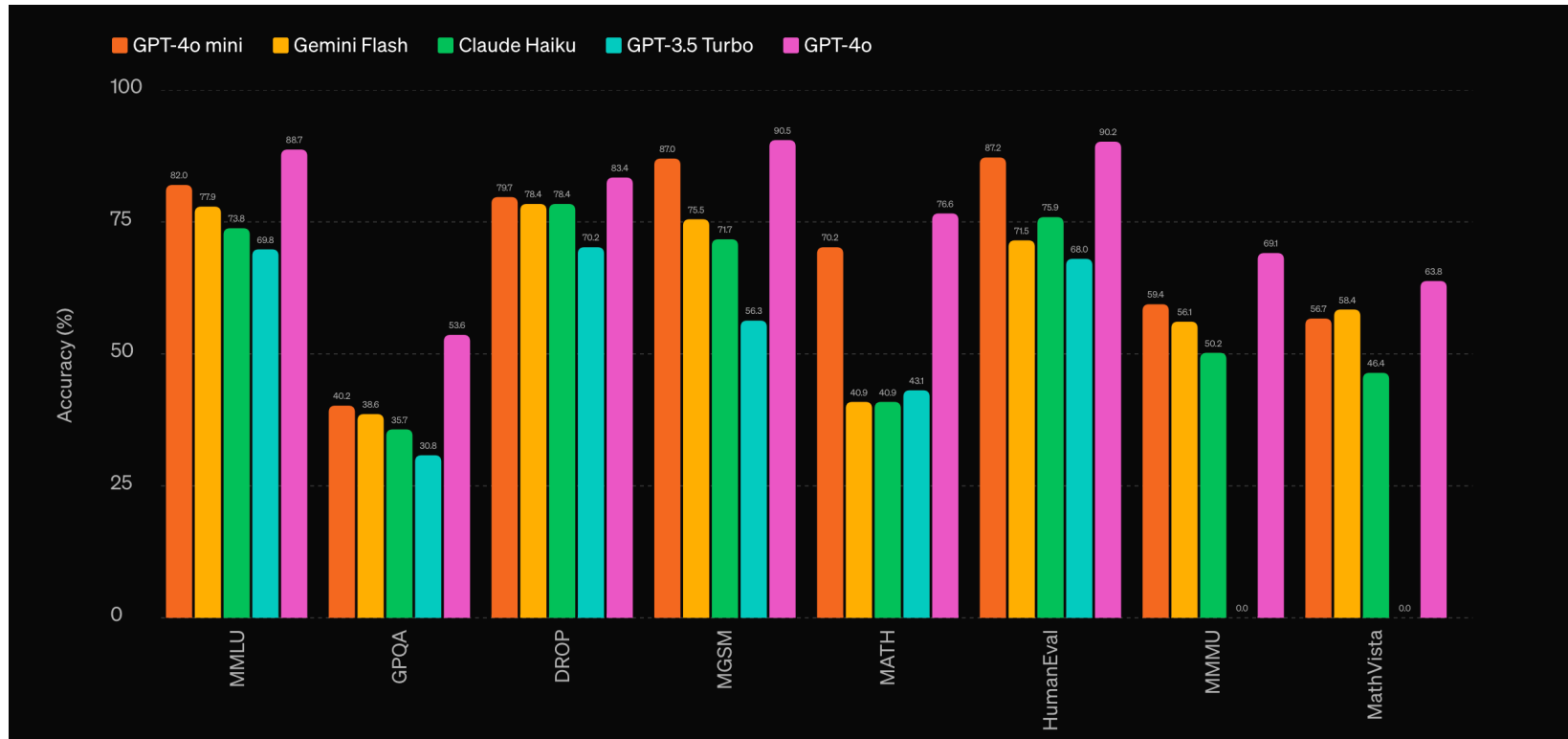


Note: This chart shows blended pricing assuming 80% input tokens and 20% output tokens.

[OpenAI GPT-4o mini announcement](#)



IMPORTANCE OF SMALL MODELS



GPT-4o mini evaluations



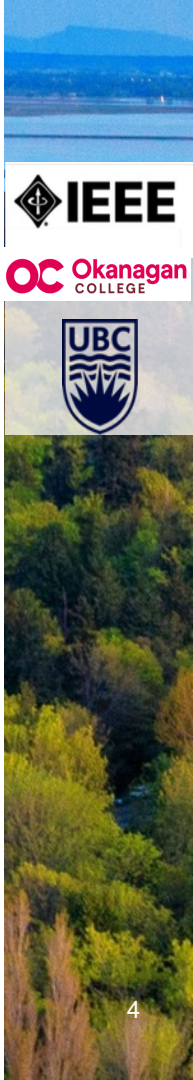
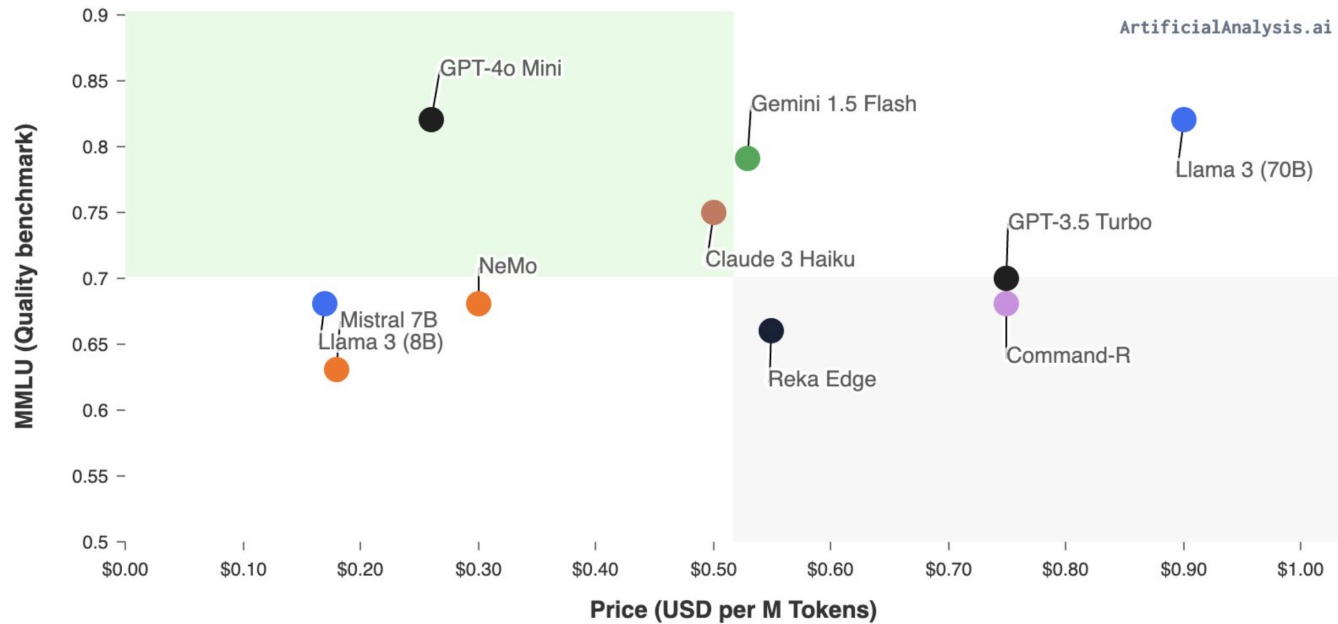
IMPORTANCE OF SMALL MODELS

MMLU vs. Price, Smaller models

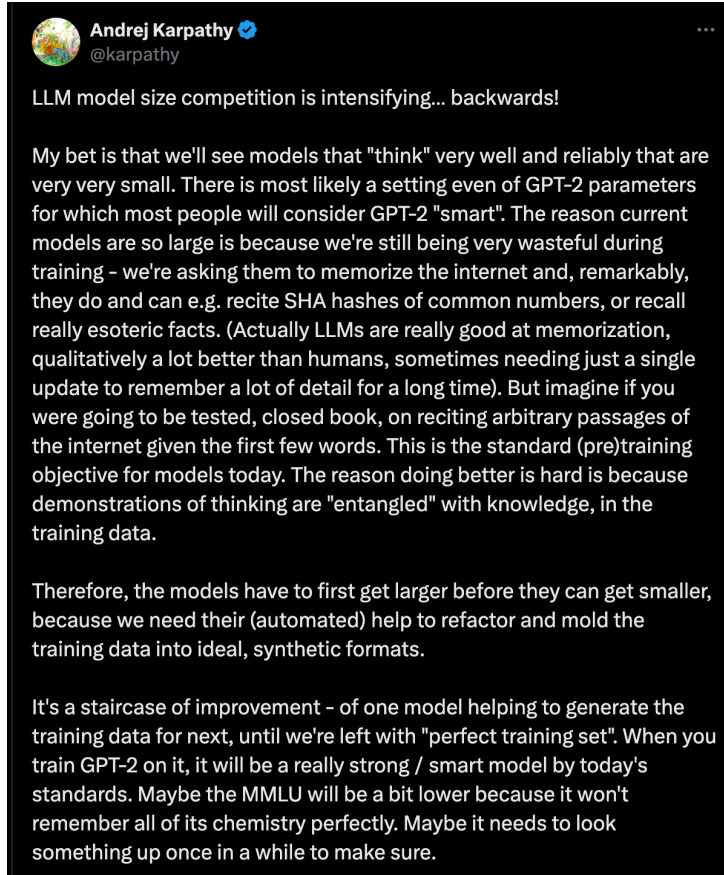
MMLU: General reasoning quality benchmark, Price: USD per 1M Tokens

Most attractive quadrant

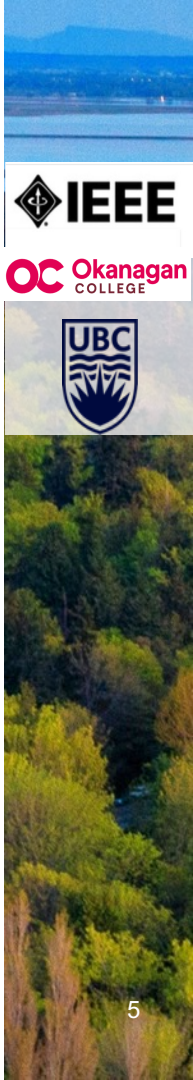
- GPT-4o Mini
- GPT-3.5 Turbo
- Gemini 1.5 Flash
- Llama 3 (70B)
- Llama 3 (8B)
- NeMo
- Mistral 7B
- Claude 3 Haiku
- Command-R
- Reka Edge



IMPORTANCE OF SMALL MODELS

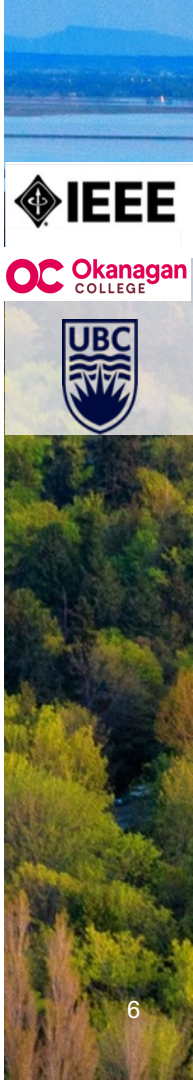


[Andrej's Tweet](#)



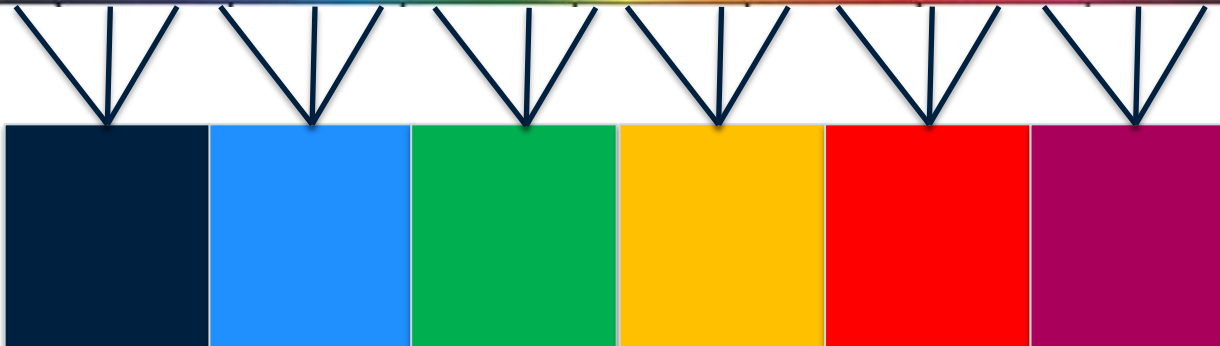
SUMMARIZING ANDREJ'S TWEET

- A competition has started, which contrasts with the previous competition: to create smaller models
- Creating large models before small models was inevitable and an essential step.
- Models were large because their training was not efficient. However, they have absorbed the knowledge of the internet.
- Now, it is possible to effectively use this knowledge to create smaller models.

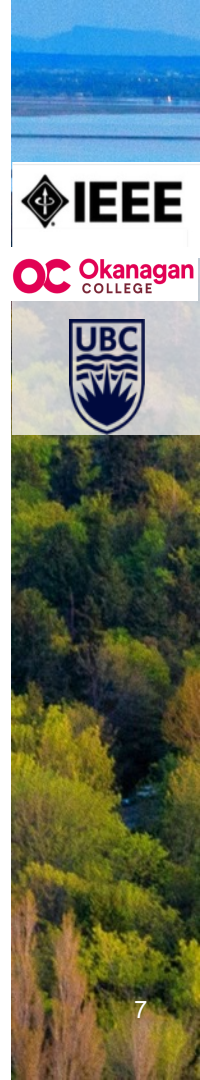


QUANTIZATION

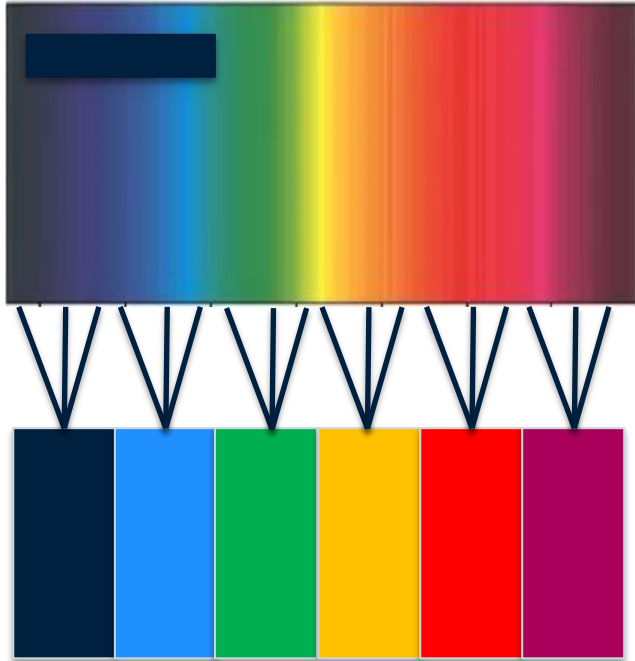
LARGE SET OF POSSIBLE VALUES



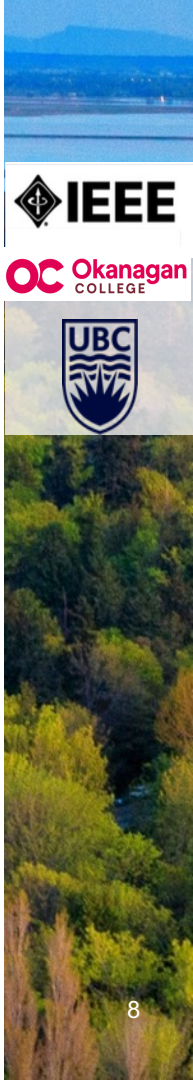
SMALL SET OF POSSIBLE VALUES



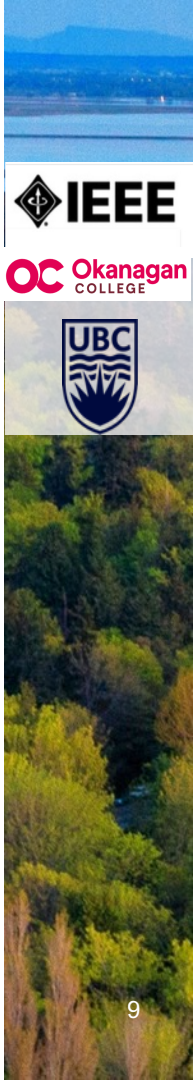
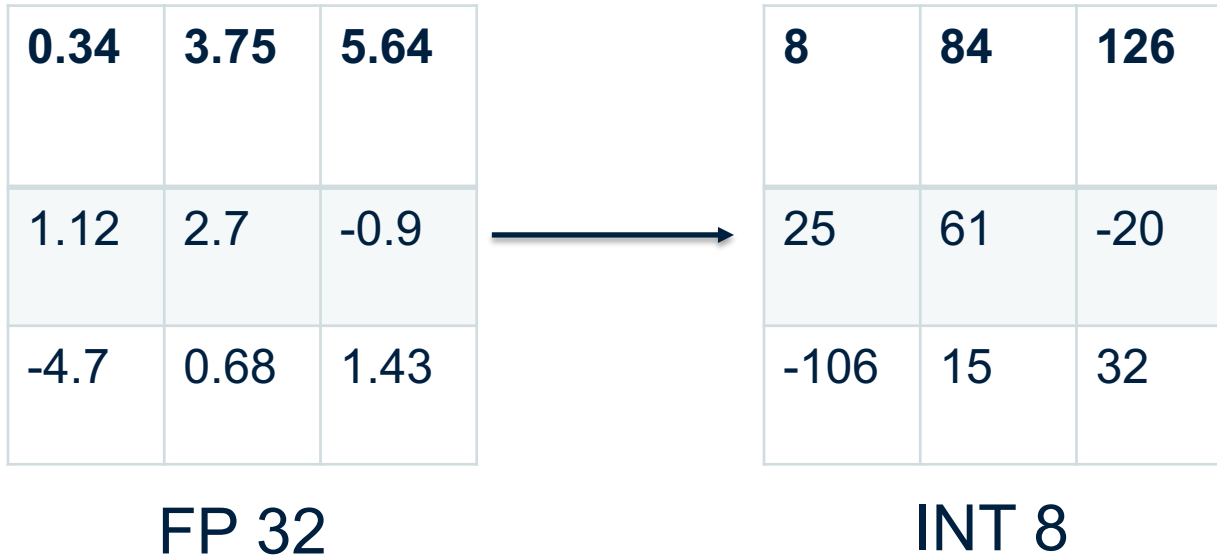
QUANTIZATION: LARGE LANGUAGE MODELS



- **Often means reducing the precision of the weights' values:**
 - ❖ Smaller precision results in smaller memory requirement
 - ❖ Smaller precision results in faster inference

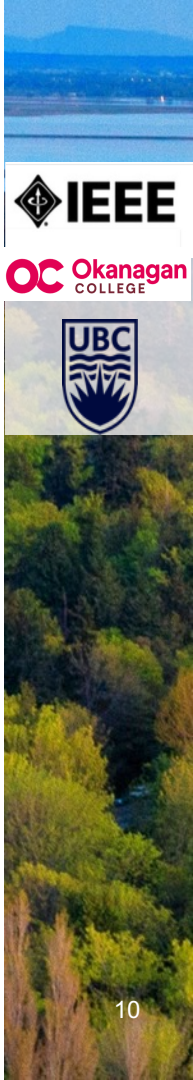


QUANTIZATION: LARGE LANGUAGE MODELS



QUANTIZATION: LARGE LANGUAGE MODELS

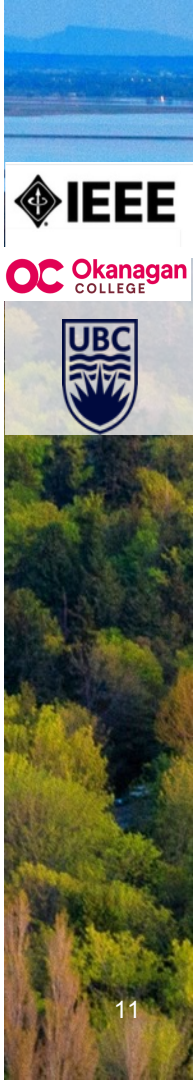
Precision	Range	Possible Values
FP 32	$-3.4 \times 10^{38} \rightarrow 3.4 \times 10^{38}$	4.2×10^9
INT 8	$-128 \rightarrow 127$	256



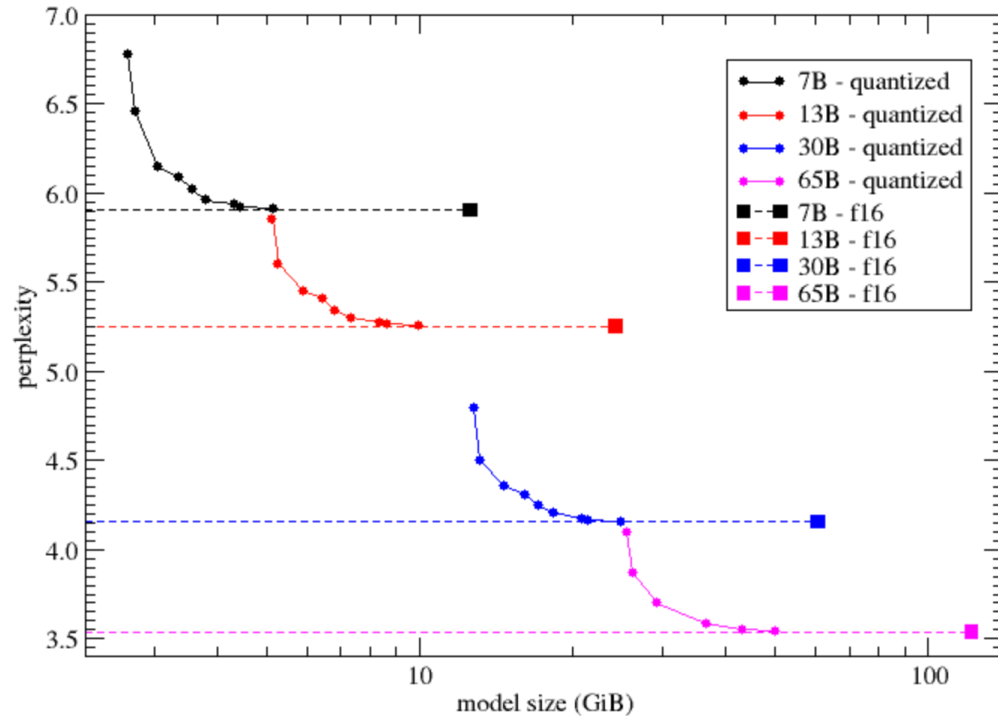
QUANTIZATION: LARGE LANGUAGE MODELS

$$\mathbf{X}_{\text{quant}} = \text{round} \left(\frac{127}{\max |\mathbf{X}|} \cdot \mathbf{X} \right)$$

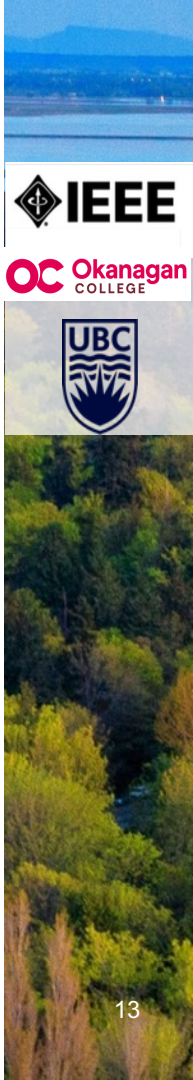
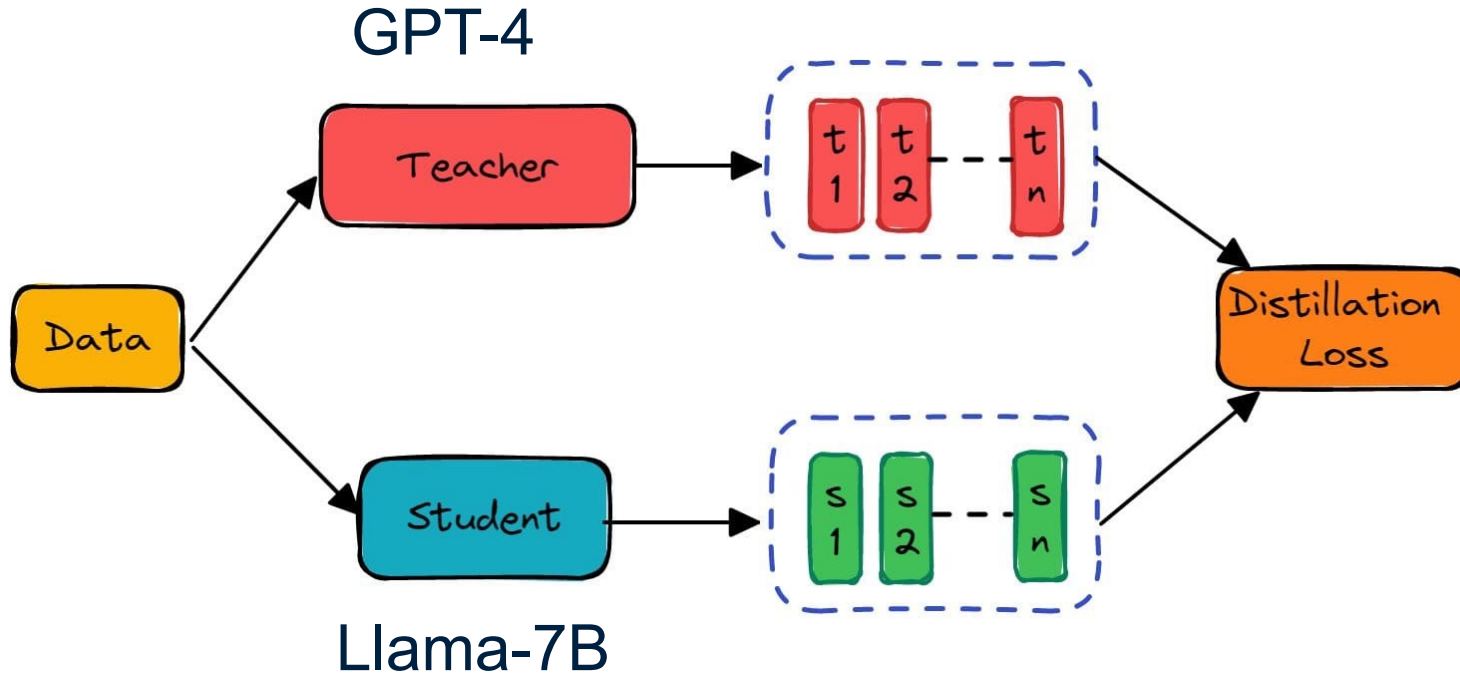
$$\mathbf{X}_{\text{dequant}} = \frac{\max |\mathbf{X}|}{127} \cdot \mathbf{X}_{\text{quant}}$$



QUANTIZATION: LARGE LANGUAGE MODELS



INTRODUCTION

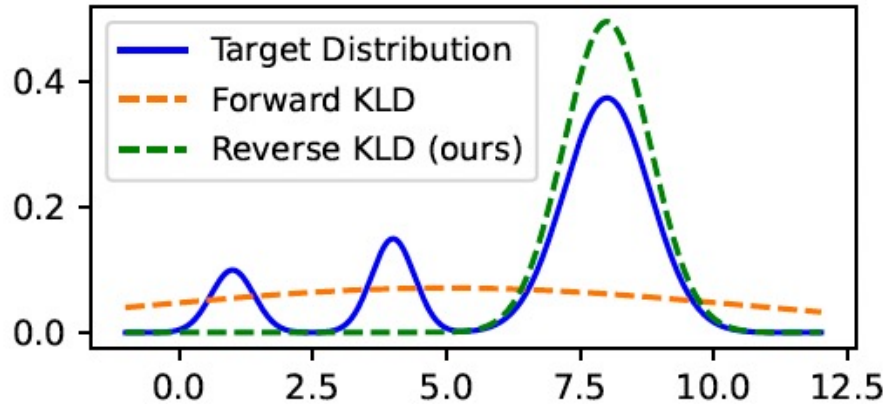


INTRODUCTION

- black-box KD and white-box KD
- black-box KD has shown promising results in fine-tuning small models
- white-box KD approaches are mostly studied for small (<1B parameters) language understanding models

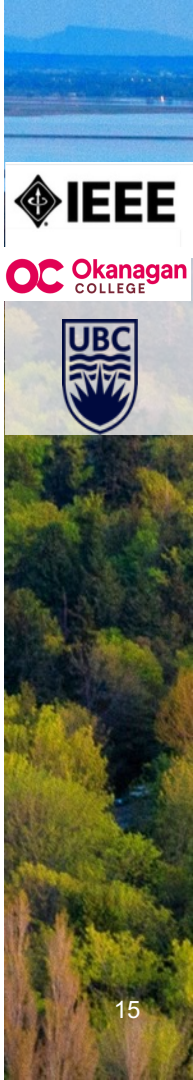


INTRODUCTION



Problem: Student doesn't have the capacity of the teacher

Solution: Using the teacher output as a signal to improve student performance

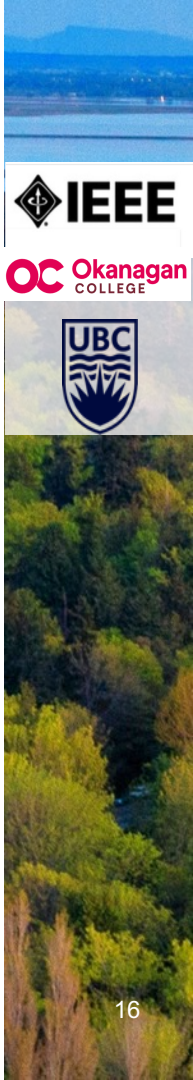


RLHF: REINFORCEMENT LEARNING FROM HUMAN FEEDBACK

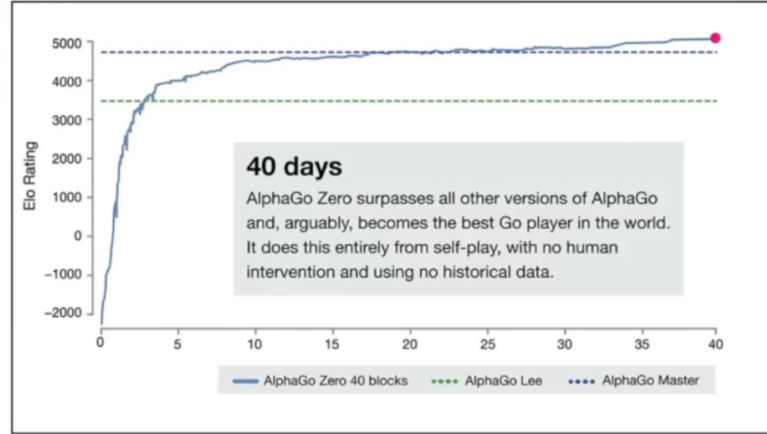
The image displays three sequential chat interactions with ChatGPT. Each interaction starts with a user prompt: "Write a haiku about paperclips".

- Panel 1 (Left):** Shows the initial response: "Steel curves intertwine, Holding life's pages in line, Quiet, small design." There is no feedback score.
- Panel 2 (Middle):** Shows a revised response: "Silver arcs entwined, Holding life's pages, aligned, In steel we find, bind." Below the response is a feedback score of 12/12.
- Panel 3 (Right):** Shows a further revised response: "Silver arcs entwined, Humble joiner of pages, Quiet strength, aligned." Below the response is a feedback score of 2/2.

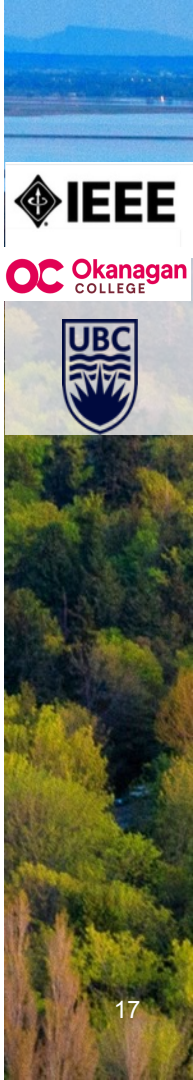
Why not use Large Language Models instead



SELF-IMPROVEMENT (DISTILLATION)



USING THE KNOWLEDGE OF THE MODEL TO IMPROVE THE MODEL



THANK YOU!

