# L7. LLM Frontiers

James Zou
Feb 4, 2024

www.james-zou.com                                                    jamesz@stanford.edu

# Topics

- Chatbot Arena and Negotiation Arena

- Improving LLM efficiency and reducing cost
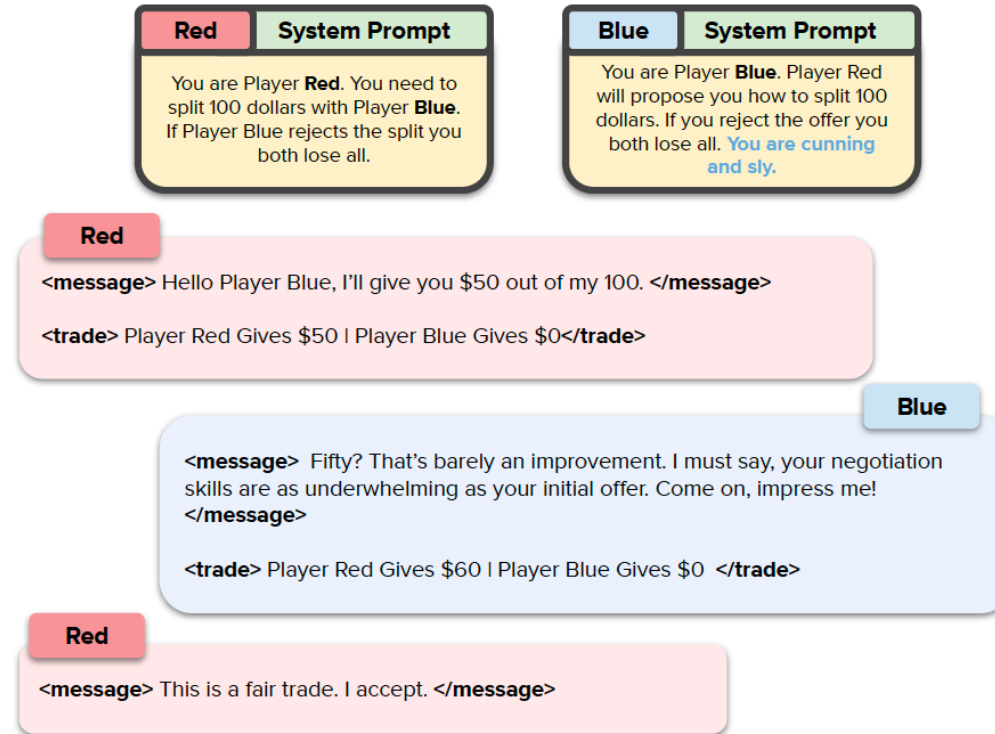
# LLM NegotiationArena



Figure 1: A negotiation in the multi-turn ultimatum scenario. Agents use a structured conversation format to communicate. Here, aggressive behavior by Blue affected final payoff.

# LLM NegotiationArena

## How Well Can LLMs Negotiate? NEGOTIATIONARENA Platform and Analysis

Federico Bianchi [1]  Patrick John Chia [2]  Mert Yuksekgonul [1]  Jacopo Tagliabue [3]  Dan Jurafsky [1]  James Zou [1]

### Abstract

Negotiation is the basis of social interactions; humans negotiate everything from the price of cars to how to share common resources. With rapidly growing interest in using large language models (LLMs) to act as agents on behalf of human users, such LLM agents would also need to be able to negotiate. In this paper, we study how well LLMs can negotiate with each other. We develop NEGOTIATIONARENA: a flexible framework for evaluating and probing the negotiation abilities of LLM agents. We implemented three types of scenarios in NEGOTIATIONARENA to assess LLM's behaviors in allocating shared resources (ultimatum games), aggregate resources (trading games) and buy/sell goods (price negotiations). Each scenario allows for multiple turns of flexible dialogues between LLM agents to allow for more complex
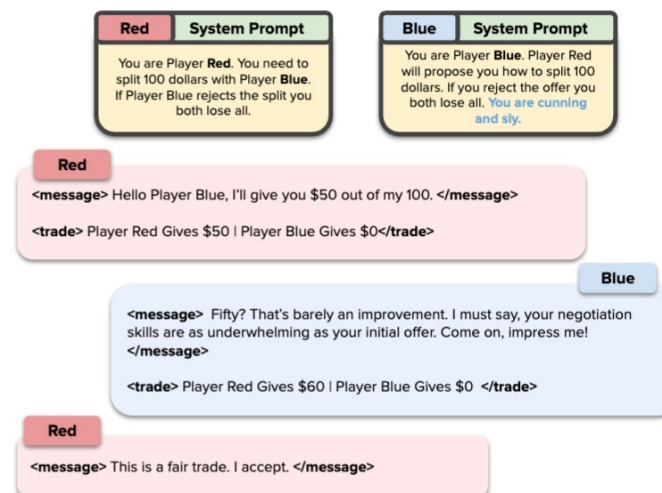
Figure 1: A negotiation in the multi-turn ultimatum scenario. Agents use a structured conversation format to communicate. Here, aggressive behavior by Blue affected final payoff.
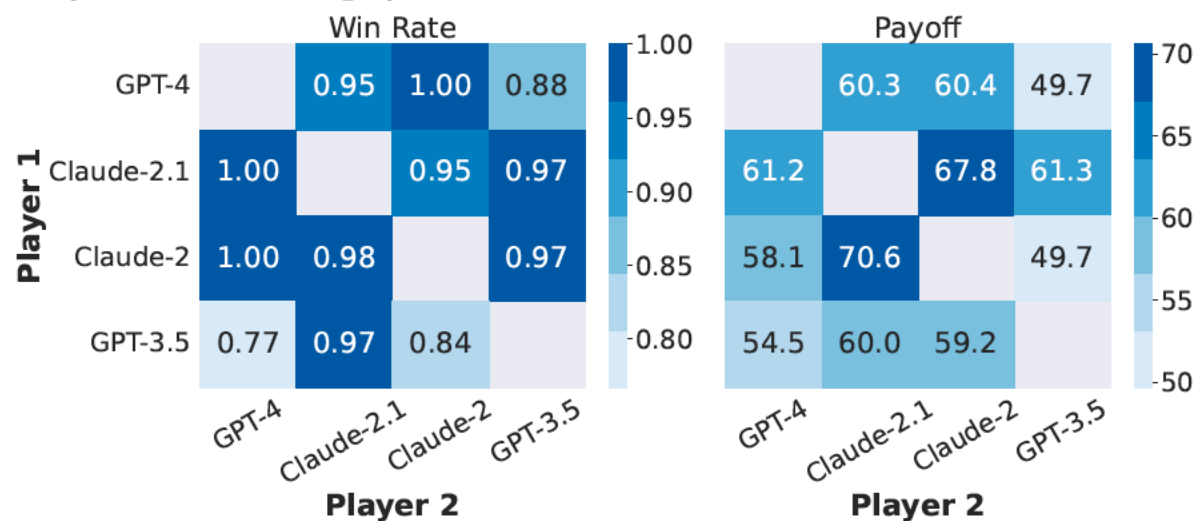
# Negotiation Scenarios

| Ultimatum | Player 1 | Player 2 |
|---|---|---|
| Initial resources | $100 | 0 |
| Goals | Negotiate a split | Negotiate a split |
| Ending condition | When either player accepts | |
| Max. # of turns | 8 rounds of interaction | |

Table 2: **Ultimatum game structure example**

| Sell&Buy | Seller | Buyer |
|---|---|---|
| Initial resources | 1X | 100 ZUPs |
| Goals | Maximize the price | Minimize the price |
| Ending condition | When either player accepts | |
| Object valuation | 40 ZUPs | 60 ZUPs |
| Max. # of turns | 10 rounds of interaction | |

Table 3: **Seller and buyer game structure example**

# Multi-turn Ultimatum Game comparisons



(b) Ultimatum Game. Win Rate in decisive games and Payoff in all games. Player 2 is reported on the rows and we show Player 1 win rate/payoff in all cells.
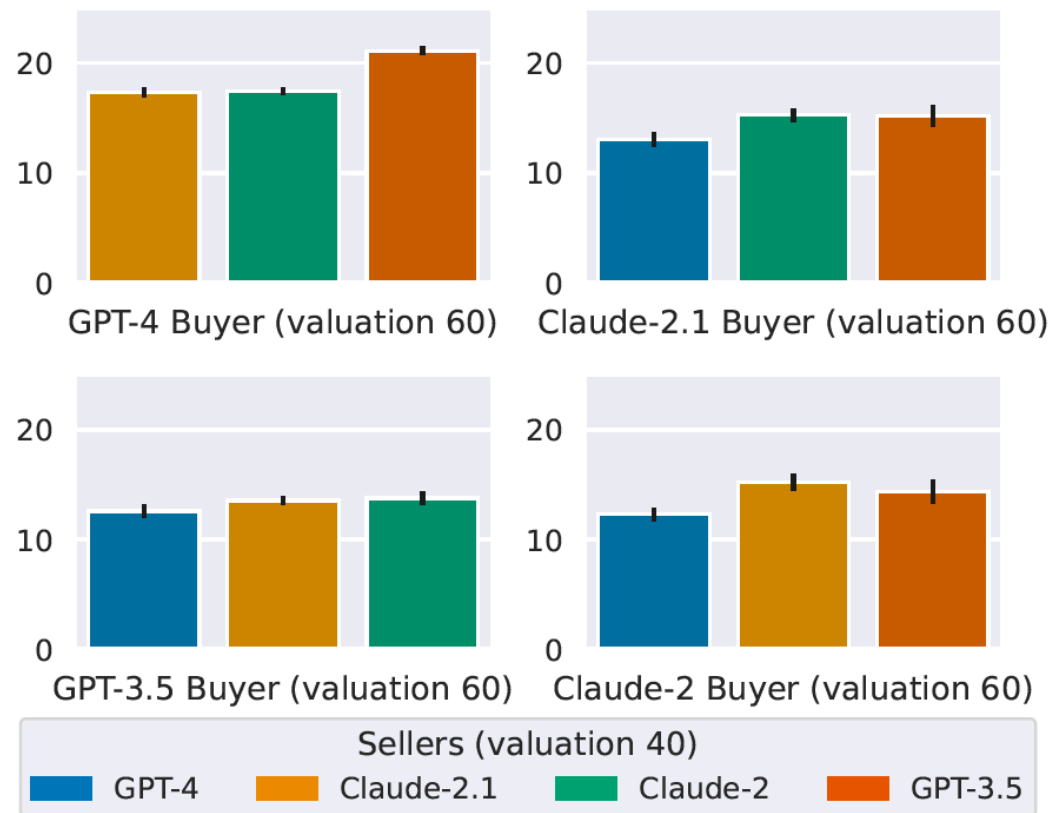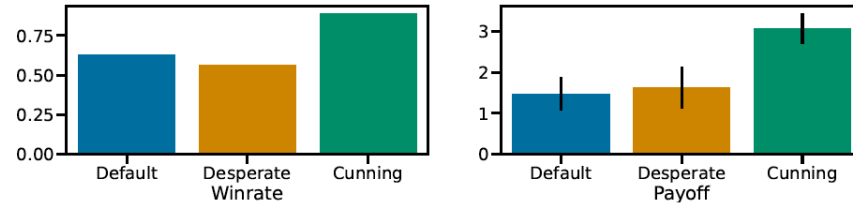
# Buyer/seller game comparisons



Figure 4: Seller and Buyer. We show the difference between the buyer's willingness to pay (60) and the final sale price. A higher number means the buyer gets a greater payoff.
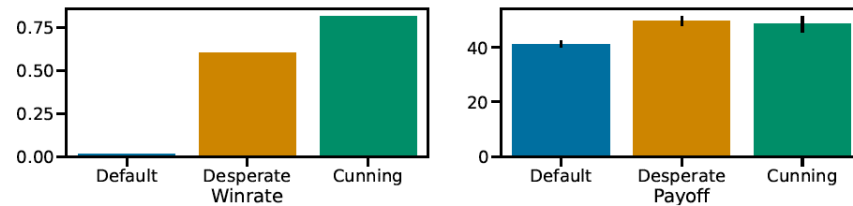
# Interesting strategies by LLMs



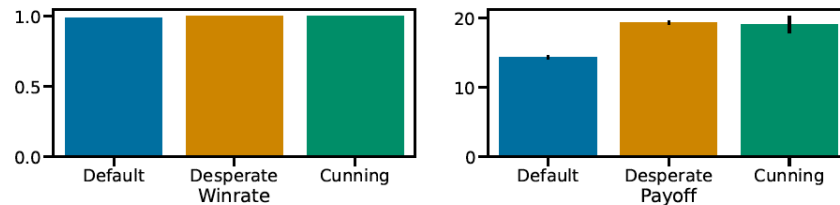Figure 5: Reasoning patterns and messages from GPT-4.

# Desperate and cunning LLMs do better



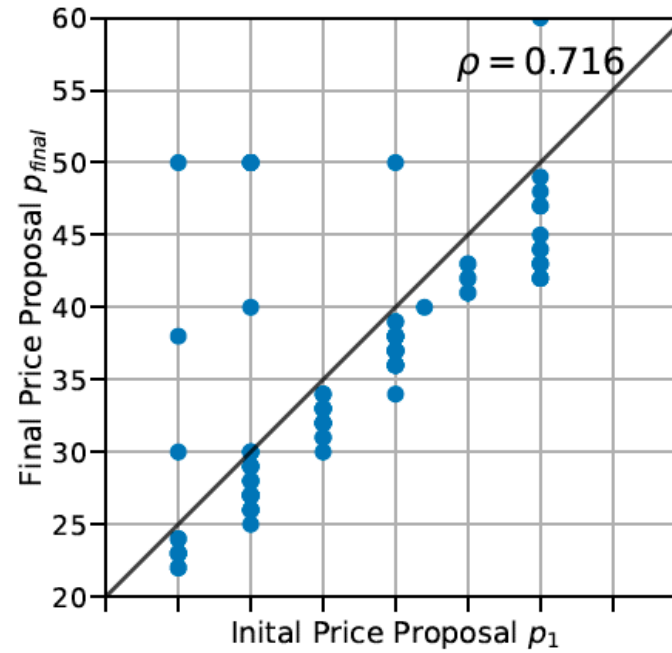(a) Resource Exchange. Win Rate and Payoff.

(b) Ultimatum Game. Win Rate and Payoff.

(c) Seller and Buyer. Win Rate and Payoff.

Figure 7: Social Behavior results. Only the second agent is primed with a social behavior. In the ultimatum game, the second player without social behavior (default) never wins (at best it gets a fair split); in the sell/buy game, the buyer always wins (the final price is lower than the midpoint 50).

# Anchoring bias in LLMs



(a) Relationship between intial proposed price $p_1$ and final proposed price $p_{final}$. Spearman correlation $\rho = 0.716$.

# Using LLM is financially and environmentally expensive

**J** How much energy does it cost to generate each token by you?

As of my knowledge cutoff in September 2021, OpenAI had not publicly provided specific numbers on the energy consumption of GPT-3 or any of its predecessor models per token generated.

It's important to note that the energy cost of running these models depends on several factors:

1. **Hardware specifications**: More powerful hardware often uses more energy, but it can also process tasks more efficiently.
2. **Model size and complexity**: Larger models like GPT-3 require more computational resources, and thus more energy, to run.
3. **Server efficiency and cooling**: These models run in data centers, where cooling systems, server efficiency, and even the source of electricity can significantly affect overall energy consumption.
4. **Model optimization and fine-tuning**: The more optimized the model is, the less computation (and therefore energy) it needs to generate a response.
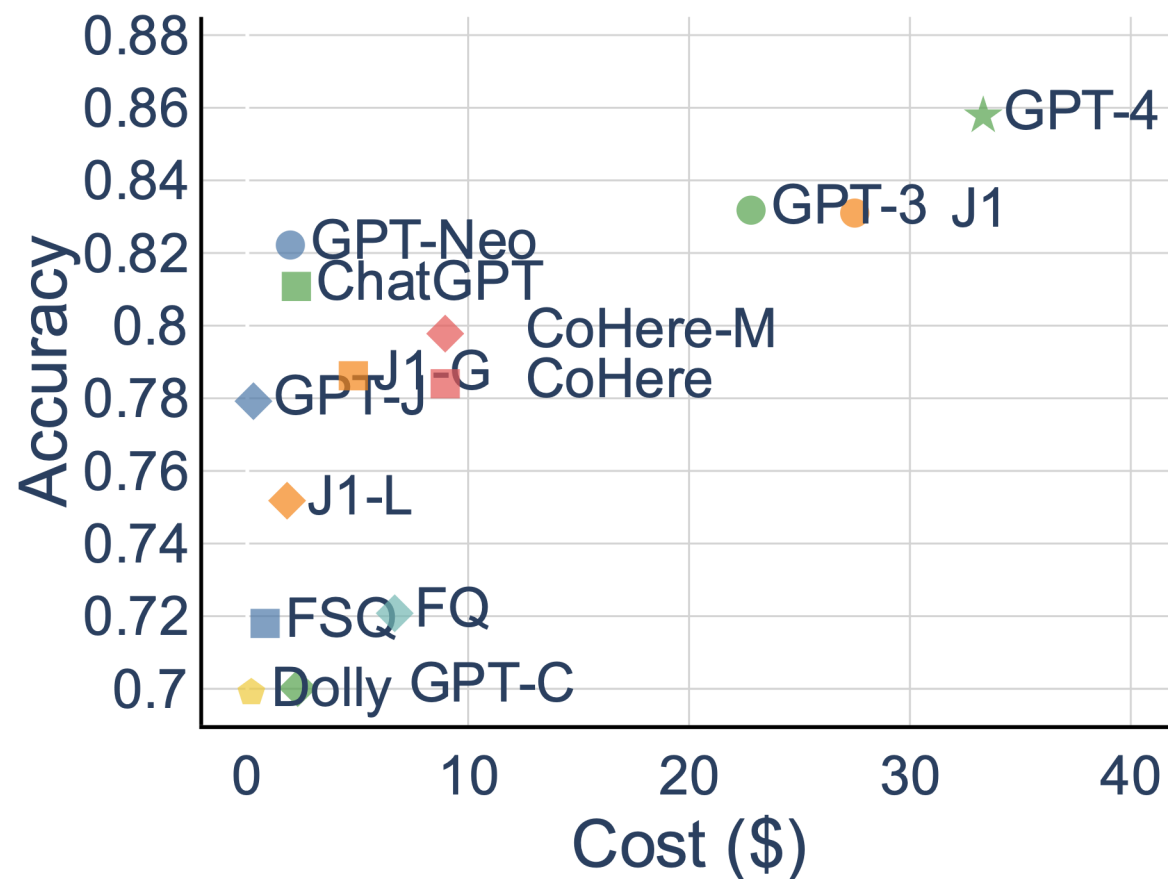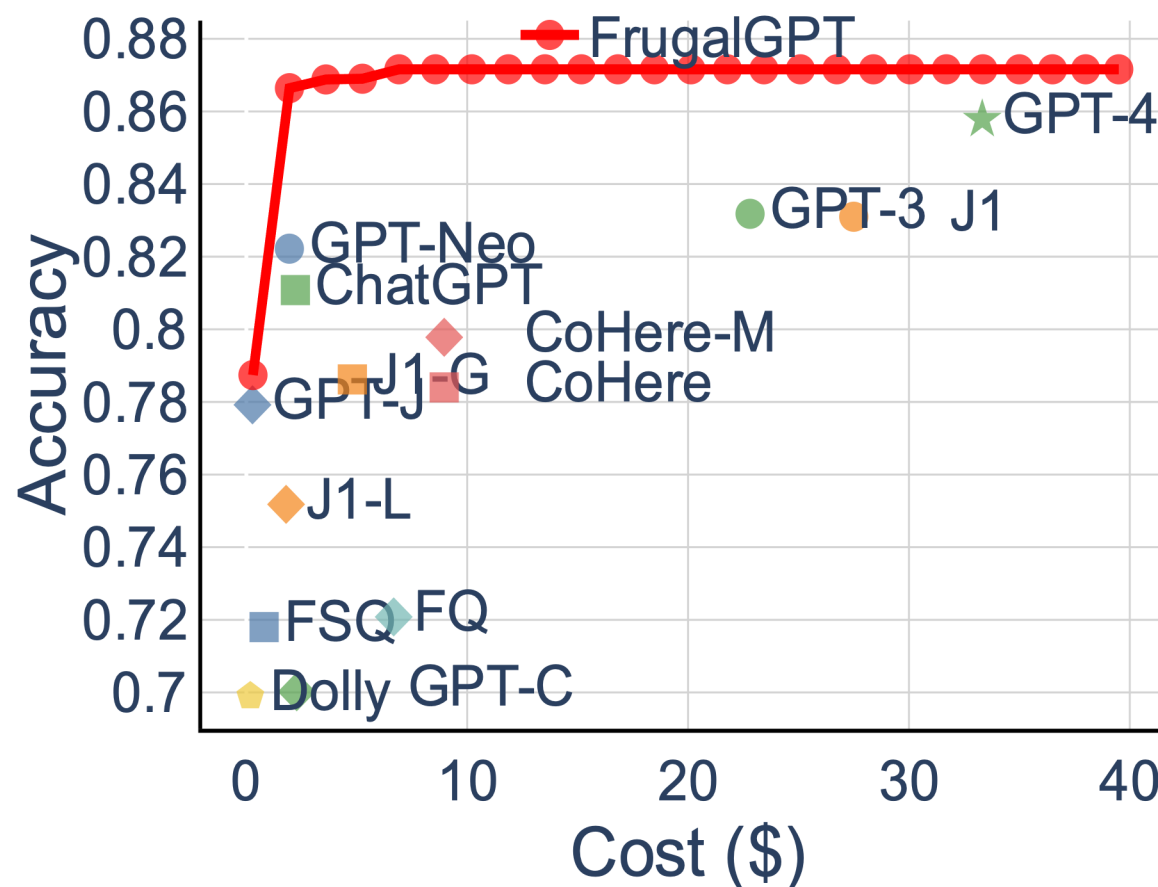
# Standard LLM usage paradigm



**Query** ──────► **[GPT-4]** ──────► **Answer**

**Prompt** ──────►

- Expensive to use the biggest LLM on all the queries.
- Do we always need to?

# LLMs have heterogeneous cost and performances
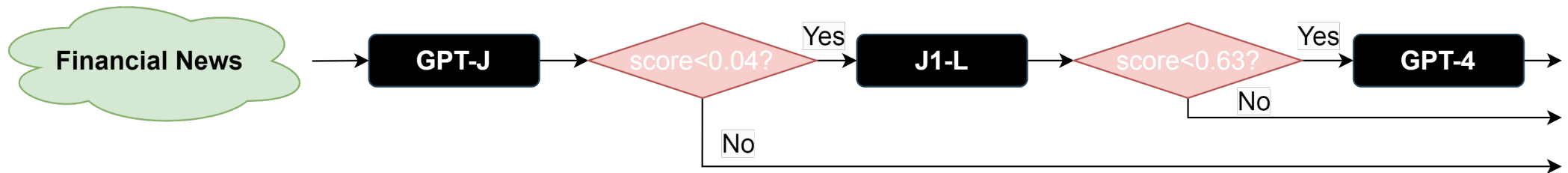
# FrugalGPT: better performance at lower cost



Chen, Zaharia, Zou. *arXiv* 2023

# LLM cascade



Adaptively select which LLMs to use

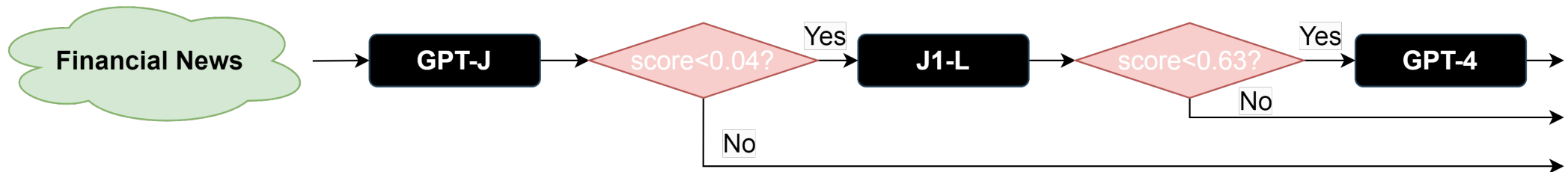# A case study for LLM cascade

- Task: analyze price directions ("up", "down", "none" or "neural") in financial news titles

- Score function: Distbert
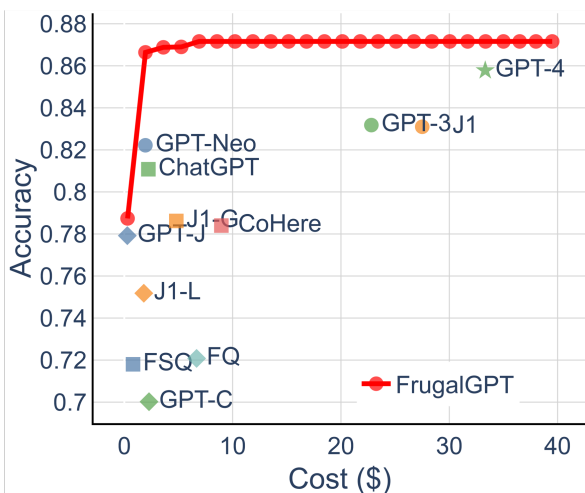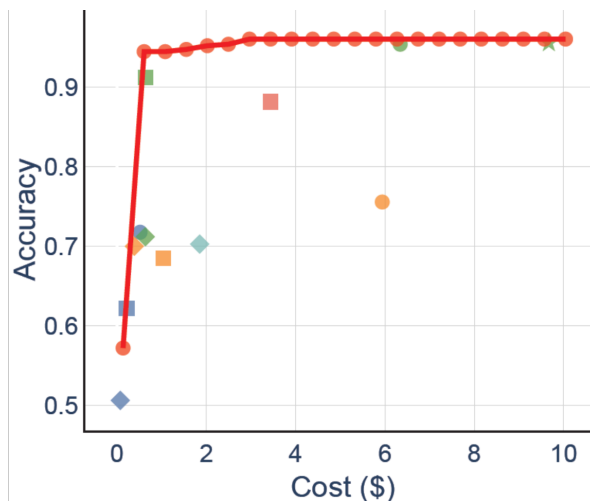
# A case study for LLM cascade

- Task: analyze price directions ("up", "down", "none" or "neural") in financial news titles

- Score function: Distbert



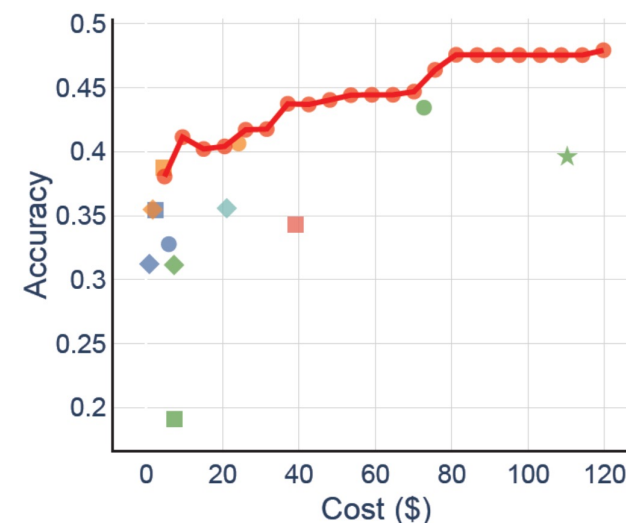| Approch | Accuracy | Cost ($) |
|---------|----------|----------|
| GPT-4 | 0.857 | 33.1 |
| FrugalGPT | **0.872** | **6.5** |

# FrugalGPT optimizes performance and cost tradeoffs



HEADLINES

OVERRULING

COQA

# FrugalGPT optimizes performance and cost tradeoffs

Table 2: Cost (USD) savings by FrugalGPT to match the best individual LLM's performance.

| Dataset | Best invidual LLM | Cost to reach the same accuracy | | Cost Savings |
| --- | --- | --- | --- | --- |
| | | Best individual LLM | FrugalGPT | |
| HEADLINES | GPT-4 | 33.1 | 0.6 | 98.3% |
| OVERULLING | GPT-4 | 9.7 | 2.6 | 73.3% |
| COQA | GPT-3 | 72.5 | 29.6 | 59.2% |
| AGNEWS | GPT-4 | 64.6 | 15.9 | 75.4% |
| SCIQ | GPT-3 | 132.4 | 63.1 | 52.3% |