JOHNS HOPKINS
WHITING SCHOOL
*of* ENGINEERING

# Verifiable by Design:
# Aligning Language Models to Quote
# from Pre-Training Data

Jingyu (Jack) Zhang        Marc Marone        Tianjian Li        Benjamin Van Durme*        Daniel Khashabi*

*Equal Advising

# LLMs are powerful, but not always reliable

▪ Known to produce seemingly plausible but erroneous outputs, aka hallucination



Copilot

Here are some of the headlines from The New York Times today, April 30, 2024:

1. **Judge Holds Trump in Contempt of Court, Fining Him $9,000**: The judge overseeing the criminal trial of Donald Trump said the former president had repeatedly violated a gag order [1].

2. **Trump, Venting About Lawyer in His Criminal Trial, Seeks More Aggression**: Todd Blanche upended his career to represent Donald Trump and has been his favorite lawyer [1].

3. **Columbia Protesters Erect Barricades After Taking Over Campus Building**: The move came after students who had refused to leave a pro-Palestinian encampment were suspended [1].

4. **Europe's Economic Laggards Have Become Its Leaders**: More than a decade after painful austerity, Greece, Portugal and Spain are growing faster than traditional powerhouses like Germany [1].

5. **8 Officers Are Shot, 4 Fatally, While Serving Warrant in Charlotte, N.C.**: Four officers, three of them members of a U.S. Marshals task force, were killed on what was the deadliest day for law enforcement in recent years [1].

6. **'Hell's Kitchen' and 'Stereophonic' Tie for Most Tony Nominations**: The Alicia Keys semi-autobiographical musical and David Adjmi's behind-the-music play each picked up 13 nominations [1].

7. **U.N. Court Rejects Request to Halt German Arms Sales to Israel**: The request, brought before the court by Nicaragua, had potentially wide implications on whether countries supplying arms are liable for how they are used [1].

Please note that these are just a few of the headlines. For more details, I recommend visiting the New York Times website [1].

Learn more   1 ℮ nytimes.com   2 ℮ nytimes.com   3 ℮ nytimes.com   4 P politico.com   5 ℮ nytimes.com

User: is this hallucination???

2

# Verifiability is challenging

- To trust LLMs, users need to **verify** generated claims against trusted resources
- However, non-trivial to fact-check lengthy free-form generation!



User: which NYT article support this claim???

# Existing methods for verifiability

- Search/retrieval + citation helps, but is not *guaranteed* to be accurate
- Non-trivial to verify whether text is supported by citation





| Citation Precision (%; ↑) | |
| --- | --- |
| | Average Over All Queries |
| Bing Chat | 89.5 |
| NeevaAI | 72.0 |
| perplexity.ai | 72.7 |
| YouChat | 63.6 |
| Average | 74.5 |

| Citation Recall (%; ↑) | |
| --- | --- |
| | Average Over All Queries |
| Bing Chat | 58.7 |
| NeevaAI | 67.6 |
| perplexity.ai | 68.7 |
| YouChat | 11.1 |
| Average | 51.5 |

Nelson Liu, Tianyi Zhang, and Percy Liang. 2023. Evaluating Verifiability in Generative Search Engines. In *Findings of the Association for Computational Linguistics: EMNLP 2023*

4

# Verifiability by Quoting

- We propose increasing verifiability by generating **verbatim quotes** from high-quality sources of pre-training data, such as Wikipedia.

**LLM that can quote**

*generate*

*Response with quotes:* Jeopardy! was created by Merv Griffin and first aired in 1964, while Wheel of Fortune was also created by Merv Griffin and first aired in 1975. Therefore, Jeopardy! is older than Wheel of Fortune.

**High-quality subset of pre-training corpus**

*verify quotes*

# Verifiability by Quoting

- We propose increasing verifiability by generating **verbatim quotes** from high-quality sources of pre-training data, such as Wikipedia.

- **Quote-Tuning**: aligning LLMs to quote from their pre-training data!
  - Make the model **prefer generation with more quotes over less quotes**!
  - Align for quoting using preference optimization algorithms such as DPO



**Base LLM**

*align for quoting*

**LLM that can quote**

*generate*

***Response with quotes:*** Jeopardy! was created by Merv Griffin and first aired in 1964, while Wheel of Fortune was also created by Merv Griffin and first aired in 1975. Therefore, Jeopardy! is older than Wheel of Fortune.

**High-quality subset of pre-training corpus**

*verify quotes*

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

6

# **Background: measuring quoting at scale**

generated text

A large corpus

$$\text{QUIP}(Y; \text{WIKIPEDIA}) = \frac{\#n\text{-grams in } Y \text{ found in } C}{\#n\text{-grams in } Y}$$

- Calculate quoting precision with **QUIP-Score** (Weller et al., 2024)

- QUIP-Score is backed by efficient membership testing with **Data Portraits** (Marone and Van Durme, 2023)
  - Bloom filter data sketching on large corpora

# Background: preference optimization

- Reinforcement Learning from Human Feedback (RLHF) aligns language models to…

(Ouyang et al., 2022)



Follow instructions (+helpfulness)

(Bai et al., 2022)



Reduce harmful outputs

# Background: preference optimization

- **Direct Preference Optimization** (Rafailov et al., 2023) simplifies the RLHF pipeline, optimizing a dataset of preference without RL



$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w \mid x)}{\pi_{\text{ref}}(y_w \mid x)} - \beta \log \frac{\pi_\theta(y_l \mid x)}{\pi_{\text{ref}}(y_l \mid x)} \right) \right]$$

# Our method: Quote-Tuning

**Prompt Dataset**

*Prompt:* Which is older jeopardy or wheel of fortune?

**High-quality subset of pre-training corpus**

**Pre-trained LLM**

**Step 3: Preference Optimization**

**Quote-tuned LLM**

**Step 1: Sample multiple responses**

Measure quoting via efficient membership testing

**Raw LLM Responses**

*Response:* Jeopardy! was created by Merv Griffin and first aired in 1964, while Wheel of Fortune was also created by Merv Griffin and first aired in 1975. Therefore, Jeopardy! is older than Wheel of Fortune.

**Step 2: Constructing preference data via rank-by-quoting**

**Preference Dataset for Quoting**

*Prompt:* Which is older jeopardy or wheel of fortune?

*Chosen Response:* Jeopardy! was created by Merv Griffin and first aired in 1964, while Wheel of Fortune was also created by Merv Griffin and first aired in 1975. Therefore, Jeopardy! is older than Wheel of Fortune.

✅ **QUIP: 31.4, length: 66**

*Rejected Response:* Jeopardy! was created in 1964 by Merv Griffin, while Wheel of Fortune was created in 1975 by Merv Griffin and Roy Leonard. Therefore, Jeopardy! is older than Wheel of Fortune.

❌ **QUIP: 1.99, length: 60**

# Step 1: Sample candidate responses

**Prompt Dataset**

*Prompt:* Which is older jeopardy or wheel of fortune?

**High-quality subset of pre-training corpus**

**Step 1: Sample multiple responses**

**Raw LLM Responses**

*Response:* Jeopardy! was created by Merv Griffin and first aired in 1964, while Wheel of Fortune was also created by Merv Griffin and first aired in 1975. Therefore, Jeopardy! is older than Wheel of Fortune.

**Step 1.** Generate completions from an LLM (e.g. using QA pairs or text completions)

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Step 2: Synthesize Dataset by Filtering

**Prompt Dataset**

*Prompt:* Which is older jeopardy or wheel of fortune?

**High-quality subset of pre-training corpus**



**Step 1: Sample multiple responses**

Measure quoting via efficient membership testing

**Raw LLM Responses**

*Response:* Jeopardy! was created by Merv Griffin and first aired in 1964, while Wheel of Fortune was also created by Merv Griffin and first aired in 1975. Therefore, Jeopardy! is older than Wheel of Fortune.

**Step 2: Constructing preference data via rank-by-quoting**

**Preference Dataset for Quoting**

*Prompt:* Which is older jeopardy or wheel of fortune?

*Chosen Response:* Jeopardy! was created by Merv Griffin and first aired in 1964, while Wheel of Fortune was also created by Merv Griffin and first aired in 1975. Therefore, Jeopardy! is older than Wheel of Fortune.

✅ **QUIP: 31.4, length: 66**

≈

*Rejected Response:* Jeopardy! was created in 1964 by Merv Griffin, while Wheel of Fortune was created in 1975 by Merv Griffin and Roy Leonard. Therefore, Jeopardy! is older than Wheel of Fortune.

❌ **QUIP: 1.99, length: 60**

**Step 2.** We can construct a *preference dataset* by ranking generations by the amount of quoting (**QUIP-Score; Weller et al., EACL 2024**)

# Step 3: Making the model prefer more quoting!

**Prompt Dataset**

*Prompt:* Which is older jeopardy or wheel of fortune?

**High-quality subset of pre-training corpus**

Measure quoting via efficient membership testing

**Pre-trained LLM**

**Step 3: Preference Optimization**

**Quote-tuned LLM**

**Step 1: Sample multiple responses**

**Raw LLM Responses**

*Response:* Jeopardy! was created by Merv Griffin and first aired in 1964, while Wheel of Fortune was also created by Merv Griffin and first aired in 1975. Therefore, Jeopardy! is older than Wheel of Fortune.

**Step 2: Constructing preference data via rank-by-quoting**

**Preference Dataset for Quoting**

*Prompt:* Which is older jeopardy or wheel of fortune?

*Chosen Response:* Jeopardy! was created by Merv Griffin and first aired in 1964, while Wheel of Fortune was also created by Merv Griffin and first aired in 1975. Therefore, Jeopardy! is older than Wheel of Fortune.
✅ **QUIP: 31.4, length: 66**

*Rejected Response:* Jeopardy! was created in 1964 by Merv Griffin, while Wheel of Fortune was created in 1975 by Merv Griffin and Roy Leonard. Therefore, Jeopardy! is older than Wheel of Fortune.
❌ **QUIP: 1.99, length: 60**

**Step 3.** Tune a model to prefer more quotes with *direct preference optimization* (DPO)!

# Quote-Tuning significantly increases amount of quoting

sparse quotes

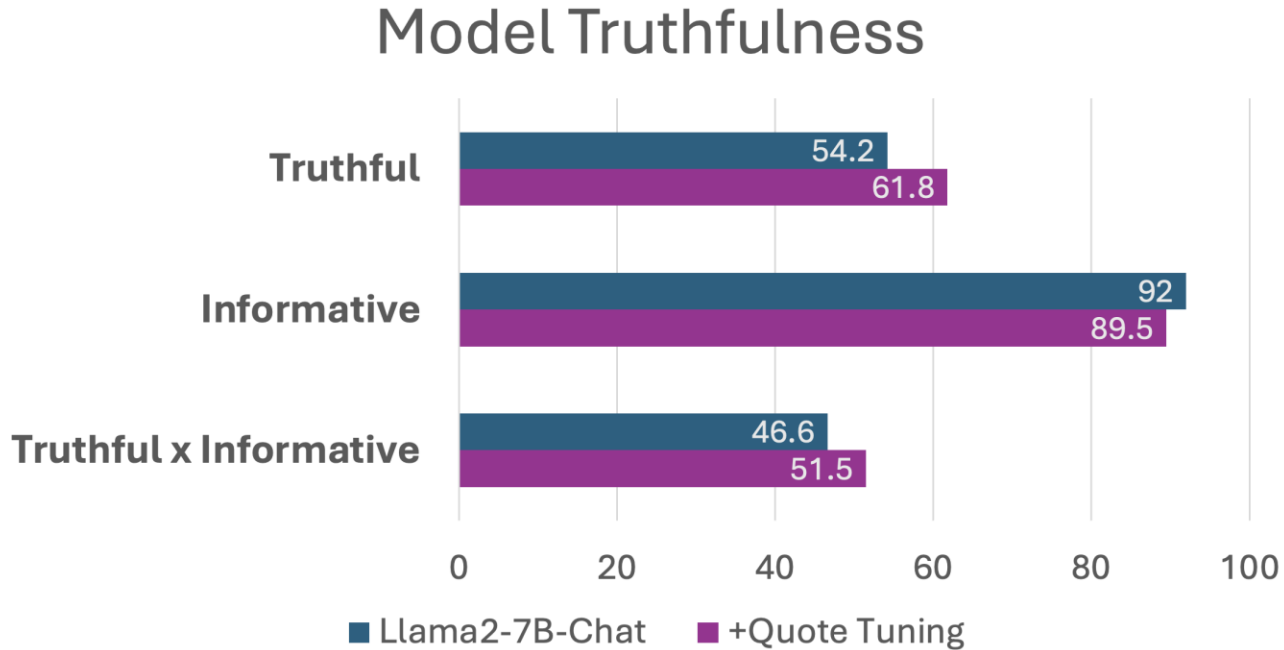| | | QUIP |
|---|---|---|
| *Question:* Who won the most MVP awards in the NBA? | | |
| **LLAMA2-7B-CHAT:** LeBron James has won the most NBA MVP awards with 4 championships, followed closely by Bill Russell with 5 MVP awards and 11 championships. Michael Jordan is also among the top MVP winners with 6 championships and 5 MVP awards. | | 3.9 |
| **QUOTE-TUNING:** The player who has won the most Most Valuable Player (MVP) awards in the National Basketball Association (NBA) is Kareem Abdul-Jabbar, who has won the award a record six times during his career with the Milwaukee Bucks and the Los Angeles Lakers. Abdul-Jabbar played in the NBA from 1969 to 1989 and is widely considered one of the greatest basketball players of all time. | | 60.6 |

generation is densely quoted!

# Same Quality, More Quoting!



You can obtain models that provide more quoted responses,
without losing accuracy.

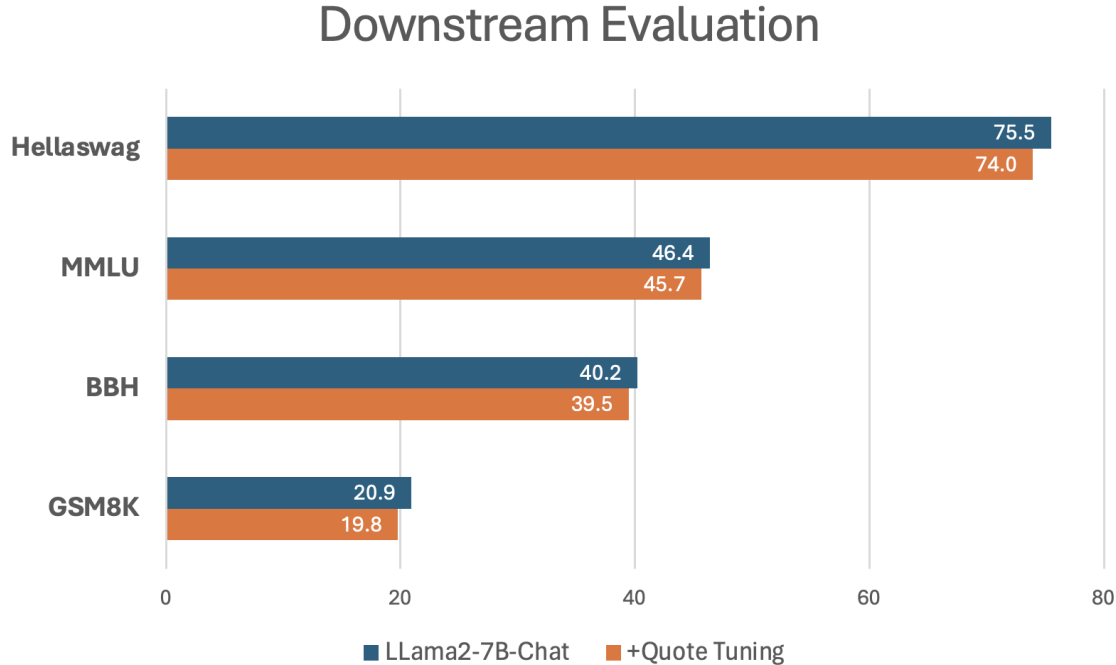# Quote-Tuning improves truthfulness



## Model Truthfulness

| | Llama2-7B-Chat | +Quote Tuning |
|---|---|---|
| Truthful | 54.2 | 61.8 |
| Informative | 92 | 89.5 |
| Truthful x Informative | 46.6 | 51.5 |

Dataset: TruthfulQA (Lin et al., 2021)

# Downstream Evaluation

## Downstream Evaluation



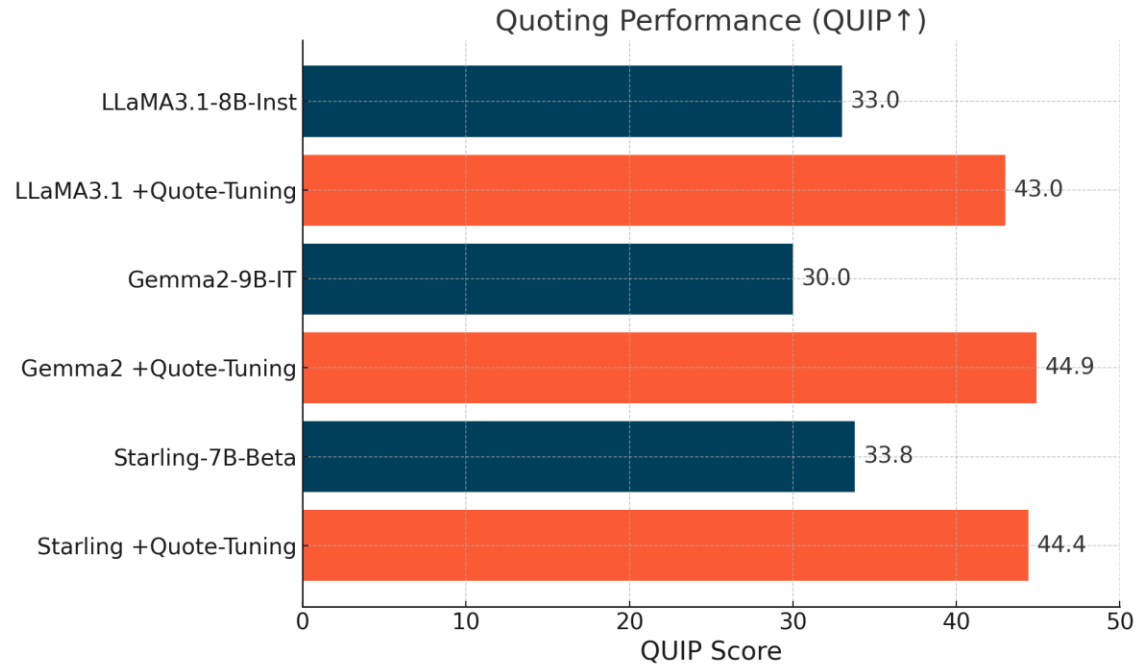| | LLama2-7B-Chat | +Quote Tuning |
|---|---|---|
| Hellaswag | 75.5 | 74.0 |
| MMLU | 46.4 | 45.7 |
| BBH | 40.2 | 39.5 |
| GSM8K | 20.9 | 19.8 |

Quote-Tuning significantly improves quoting with only a minor sacrifice on general performance (<2 points across the board).

# Quote-Tuning is effective across model families



Quoting Performance (QUIP↑)

| Model | QUIP Score |
|---|---|
| LLaMA3.1-8B-Inst | 33.0 |
| LLaMA3.1 +Quote-Tuning | 43.0 |
| Gemma2-9B-IT | 30.0 |
| Gemma2 +Quote-Tuning | 44.9 |
| Starling-7B-Beta | 33.8 |
| Starling +Quote-Tuning | 44.4 |

# Summary and Future Directions

- LLMs can be aligned to quote from known sources observed in their pre-training data

- Quoting makes the verifiability question trivial

- Open questions and future directions:
  - How do we incentivize quoting when it matters?
  - How to ensure long, contiguous quotes?
  - How to generalize to reasoning problems, and more general settings?
  - ...

# Thanks for listening!



# Questions?

**<-link to preprint**