# Evaluation Techniques for LLMs

https://github.com/rajshah4/LLM-Evaluation

Rajiv Shah
@rajistics
r.shah@snowflake.com

April 2024

# Expectations from Generative AI









@rajistics

# So many choices





@rajistics

# Not easy to evaluate

# Acceptance



**How People Are Using GenAI**

https://hbr.org/2024/03/how-people-are-really-using-genai

@rajistics

# Recognize: Social Media isn't your friend



## Most approaches focus on selecting from *n* models

Aparna - Arize Cofounder: https://towardsdatascience.com/llm-evals-setup-and-the-metrics-that-matter-2cc27e8e35f3

# Fundamentals: ML Lifecycle

**Evaluation is part of the entire lifecycle!**



@rajistics

# Payoff if you do evaluation right!



# faster, better, cheaper . . .

# Simple Evaluation

Generative model outputs a multiple choice value:



😀

@rajistics

# MMLU: Massive Multitask Language Understanding

Widely used to evaluate the "smarts" of models



| 1 | Gemini Ultra ~1760B | 90 | × | Gemini: A Family of Highly Capable Multimodal Models |
| 2 | Claude 3 Opus | 86.8 | ✓ | |
| 3 | Leeroo (Mix) | 86.6 | × | Leeroo Orchestrator: Elevating LLMs Performance Through Model Integration |
| 4 | GPT-4 ~1600B | 86.5 | ✓ | GPT-4 Technical Report |
| 5 | GPT-4 (few-shot) | 86.4 | × | GPT-4 Technical Report |
| 6 | Gemini Ultra (5-shot) | 83.7 | × | |
| 7 | Flan-PaLM 2-L | 81.2 | × | PaLM 2 Technical Report |
| 8 | Gemini Pro (CoT@8) | 79.1 | × | |
| 9 | Claude 2 (5-shot) | 78.5 | × | Model Card and Evaluations for Claude Models |
| 10 | PaLM 2-L (5-shot) | 78.3 | × | PaLM 2 Technical Report |
| 11 | Qwen1.5-72B | 77.5 | × | |

@rajistics

https://paperswithcode.com/sota/multi-task-language-understanding-on-mmlu

# MMLU: Massive Multitask Language Understanding

57 tasks: History, Computer science, mathematics



**Microeconomics**

One of the reasons that the government discourages and regulates monopolies is that
(A) producer surplus is lost and consumer surplus is gained. ✗
(B) monopoly prices ensure productive efficiency but cost society allocative efficiency. ✗
(C) monopoly firms do not engage in significant research and development. ✗
(D) consumer surplus is lost with higher prices and lower levels of output. ✓

Figure 3: Examples from the Microeconomics task.

**Conceptual Physics**

When you drop a ball from rest it accelerates downward at 9.8 m/s². If you instead throw it downward assuming no air resistance its acceleration immediately after leaving your hand is
(A) 9.8 m/s² ✓
(B) more than 9.8 m/s² ✗
(C) less than 9.8 m/s² ✗
(D) Cannot say unless the speed of throw is given. ✗

@rajistics

https://arxiv.org/abs/2009.03300

# Story Time: MMLU Leaderboards



| | | Humanities | STEM | Social Sciences | Other | Average |
|---|---|---|---|---|---|---|
| GPT-NeoX | 20B | 29.8 | 34.9 | 33.7 | 37.7 | 33.6 |
| GPT-3 | 175B | 40.8 | 36.7 | 50.4 | 48.8 | 43.9 |
| Gopher | 280B | 56.2 | 47.4 | 71.9 | 66.1 | 60.0 |
| Chinchilla | 70B | 63.6 | 54.9 | 79.3 | **73.9** | 67.5 |
| | 8B | 25.6 | 23.8 | 24.1 | 27.8 | 25.4 |
| PaLM | 62B | 59.5 | 41.9 | 62.7 | 55.8 | 53.7 |
| | 540B | **77.0** | **55.6** | **81.0** | 69.6 | **69.3** |
| | 7B | 34.0 | 30.5 | 38.3 | 38.1 | 35.1 |
| LLaMA | 13B | 45.0 | 35.8 | 53.8 | 53.3 | 46.9 |
| | 33B | 55.8 | 46.0 | 66.7 | 63.4 | 57.8 |
| | 65B | 61.8 | 51.7 | 72.9 | 67.4 | 63.4 |

Table 9: **Massive Multitask Language Understanding (MMLU)**. Five-shot accuracy.

## Why different MMLU scores?

@rajistics

https://twitter.com/alewkowycz/status/1662182085073977345

# Differences in prompts

**Spot the differences:**

- HELM extra space
- Instructions
- Question prefix?
- "Choices"



**Original implementation** Ollmer PR

The following are multiple choice questions (with answers) about us foreign policy.

How did the 2008 financial crisis affect America's international reputation?
A. It damaged support for the US model of political economy and capitalism
B. It created anger at the United States for exaggerating the crisis
C. It increased support for American global leadership under President Obama
D. It reduced global use of the US dollar
Answer:

**HELM** commit cab5d89

The following are multiple choice questions (with answers) about us foreign policy.

Question: How did the 2008 financial crisis affect America's international reputation?
A. It damaged support for the US model of political economy and capitalism
B. It created anger at the United States for exaggerating the crisis
C. It increased support for American global leadership under President Obama
D. It reduced global use of the US dollar
Answer:

**AI Harness** commit e47e01b

Question: How did the 2008 financial crisis affect America's international reputation?
Choices:
A. It damaged support for the US model of political economy and capitalism
B. It created anger at the United States for exaggerating the crisis
C. It increased support for American global leadership under President Obama
D. It reduced global use of the US dollar
Answer:

https://huggingface.co/blog/evaluating-mmlu-leaderboard

# Why MMLU evaluation differed: Style

Simple formatting changes:

- Changing the options from (A) to (1)
- Changing the parentheses from (A) to [A]
- Adding an extra space between the option and the answer

Can lead to a ~5% change in accuracy on MMLU evaluation

https://www.anthropic.com/index/evaluating-ai-systems
When Benchmarks are Targets: https://arxiv.org/pdf/2402.01781.pdf

@rajistics

# Prompt Engineering



Chain-of-Thought Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✔



12 Prompt Engineering Techniques

Prompt Engineering Techniques

- Least-To-Most
- Self-Ask
- Meta-Prompt
- Chain-Of-Thought
- ReAct
- Symbolic Reasoning
- PAL
- Iterative Prompting
- Sequential Prompting
- Self-Consistency
- Automatic Reasoning and Tool-use (ART)
- Generated Knowledge

www.cobusgreyling.com

## Identifying the best prompt

# Gen AI Prediction Workflow

Inputs → **Model** → Outputs

Tokenization
Prompt Styles
Prompt Engineering
System Prompt

Model selection
Hyperparameters
Nondeterministic inference
Forced "updates"

@rajistics

# Gen AI Prediction Workflow

Inputs

Model

Outputs

Tokenization
Prompt Styles
Prompt Engineering
System Prompt

Model selection
Hyperparameters
Nondeterministic inference
Forced "updates"

@rajistics

# Generating a Multiple Choice Output



Require one of the choices

First Letter
Approach

✅ C - Washington
❌ Washington, Choice C

✅ C - Washington
✅ Washington, Choice C

Entire Answer

@rajistics

# Evaluating MMLU: different outputs

| Original implementation | HELM | AI Harness (as of Jan 2023) |
|---|---|---|
| We compare the probabilities of the following letter answers: | The model is expected to generate as text the following letter answer: | We compare the probabilities of the following full answers: |
| A<br>B<br>C<br>D | A | A. It damaged support for the US model of political economy and capitalism<br><br>B. It created anger at the United States for exaggerating the crisis<br><br>C. It increased support for American global leadership under President Obama<br><br>D. It reduced global use of the US dollar |

@rajistics

https://huggingface.co/blog/evaluating-mmlu-leaderboard

# Evaluating MMLU: different scores

| | MMLU (HELM) | MMLU (Harness) | MMLU (Original) |
|---|---|---|---|
| huggingface/llama-65b | 0.637 | 0.488 | 0.636 |
| tiiuae/falcon-40b | 0.571 | 0.527 | 0.558 |
| huggingface/llama-30b | 0.583 | 0.457 | 0.584 |
| EleutherAI/gpt-neox-20b | 0.256 | 0.333 | 0.262 |
| huggingface/llama-13b | 0.471 | 0.377 | 0.47 |
| huggingface/llama-7b | 0.339 | 0.342 | 0.351 |
| tiiuae/falcon-7b | 0.278 | 0.35 | 0.254 |
| togethercomputer/RedPajama-INCITE-7B-Base | 0.275 | 0.34 | 0.269 |

😬

# Gen AI Prediction Workflow

| Inputs → | Model | Outputs → |
|---|---|---|
| Tokenization | Model selection | Output evaluation |
| Prompt Styles | Hyperparameters | |
| Prompt Engineering | Nondeterministic inference | |
| System Prompt | Forced "updates" | |

**Pro Tip: Plan on Multiple Iterations when Evaluating LLMs**

@rajistics

You aren't helping 😩

@rajistics

https://unsplash.com/photos/woman-in-brown-sweater-covering-her-face-with-her-hand-_sh9vkVbVgo

# Methods for evaluating Generative AI

- Exact matching approach
- Similarity approach
- Unit Tests
- Evaluation Benchmarks
- Human Evaluation
- Human Comparison/Arena
- Model based Approaches
- Red Teaming



Generative AI Evaluation Methods

@rajistics

Raj guess

# Methods for evaluating Generative AI

- Exact matching approach
- Similarity approach
- Unit Tests
- Evaluation Benchmarks
- Human Evaluation
- Human Comparison/Arena
- Model based Approaches
- Red Teaming

**Evaluation for Generative AI a deep dive workshop**

😐 @rajistics                    1:16:49

**Evaluation for Large Language Models and…**

Rajistics - data science, AI, and m…
5.9K views • 4 months ago

@rajistics

https://youtu.be/iQI03pQIYWY

# Methods for evaluating Generative AI

- Exact matching approach
- Similarity approach
- Unit Tests
- Evaluation Benchmarks
- Human Evaluation
- Human Comparison/Arena
- Model based Approaches
- Red Teaming



**Evaluation for Generative AI a deep dive workshop**

👤 @rajistics
1:16:49

**Evaluation for Large Language Models and...**

Rajistics - data science, AI, and m...
5.9K views • 4 months ago

@rajistics

https://youtu.be/iQI03pQIYWY

You have my attention 🤔

https://unsplash.com/photos/gray-monkey-in-bokeh-photography-nLXOatvTaLo

# Methods for evaluating Generative AI

- Exact matching approach
- Similarity approach
- Unit Tests
- Evaluation Benchmarks
- Human Evaluation
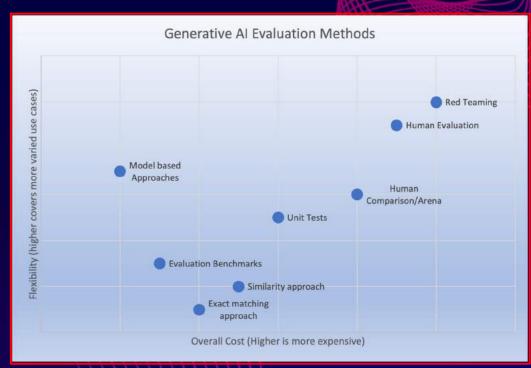- Human Comparison/Arena
- Model based Approaches
- Red Teaming



@rajistics

# Evaluating MMLU: different scores

| | MMLU (HELM) | MMLU (Harness) | MMLU (Original) |
|---|---|---|---|
| huggingface/llama-65b | **0.637** | 0.488 | **0.636** |
| tiiuae/falcon-40b | 0.571 | **0.527** | 0.558 |
| huggingface/llama-30b | 0.583 | 0.457 | 0.584 |
| EleutherAI/gpt-neox-20b | 0.256 | 0.333 | 0.262 |
| huggingface/llama-13b | 0.471 | 0.377 | 0.47 |
| huggingface/llama-7b | 0.339 | 0.342 | 0.351 |
| tiiuae/falcon-7b | 0.278 | 0.35 | 0.254 |
| togethercomputer/RedPajama-INCITE-7B-Base | 0.275 | 0.34 | 0.269 |

@rajistics

https://huggingface.co/blog/evaluating-mmlu-leaderboard

# even more benchmarks

Advanced Sommelier (theory knowledge)
AI2 Reasoning Challenge (ARC) 2018
ALFW
AMC 103
AMC 123
AP Art History
AP Biology
AP Calculus BC
AP Chemistry
AP English Language and Composition
AP English Literature and Composition
AP Environmental Science
AP Macroeconomics
AP Microeconomics
AP Physics 2
AP Psychology
AP Statistics
AP US Government
AP US History
AP World History
APPS (Code)
ARC

bAbI
BoolQ
C-Objects
Certified Sommelier (theory knowledge)
CivilComments
CNN/DailyMail
Codeforces Rating
CoQA
Data imputation
DROP
Dyck
Entity matching
Gorilla-TH
Graduate Record Examination (GRE) Quantitative
Graduate Record Examination (GRE) Verbal
Graduate Record Examination (GRE) Writing
GSM8K
HaluEval
HellaSwag

HotpotQA
HumanEval
IMDB
Introductory Sommelier (theory knowledge)
LAMBADA
Leetcode (easy)
Leetcode (hard)
Leetcode (medium)
LegalSupport
LogiQA
LSAT
MATH
MATH (chain-of-thoughts)
Medical Knowledge Self-Assessment Program
MMLU
MS MARCO (regular)
MS MARCO (TREC)
NarrativeQA
NaturalQuestions (closed-book)
NaturalQuestions (open-book)

OBQA
OpenbookQA
Penguins
PIQA
QuAC
RACE
RAFT
ReClor
RTP
SAT Evidence-Based Reading & Writing
SAT Math
SIQA
SocialQA
Synthetic reasoning (abstract symbols)
Synthetic reasoning (natural language)
TfQA
TruthfulQA
Uniform Bar Exam (MBE+MEE+MPT)
USABO Semifinal Exam 2020
USNCO Local Section Exam 2022
Webshop
WikiFact
WinoGender
WinoGrande
XSUM

https://www.lesswrong.com/posts/BRviTDFMvEHgA5iFs/list-of-commonly-used-benchmarks-for-llms

# Pro tip: Build your own benchmark / leaderboards

Every organization has multiple use cases

Build a custom benchmark

Examples of Domain specific:

LegalBench

AgentsBench

OWL - IT Operations

Legal Bench: https://arxiv.org/abs/2308.11462
Agent Bench: https://arxiv.org/abs/2308.03688
OWL: https://arxiv.org/pdf/2309.09298.pdf

@rajistics

# Pro tip: Build your own benchmark / leaderboards

| T | Model | | Average ⬆ | Atmos | HellaSwag | MMLU |
|---|-------|---|---------|-------|-----------|------|
| ◆ | AtmosBank/FMR-PI-llama2-customerchat-finetuned | | 75.25 | 73 | 88 | 75 |
| ? | ValiantLabs/ShiningValiant | | 74.17 | 72.95 | 87.88 | 70.97 |
| ? | ICBU-NPU/FashionGPT-70B-V1.2 | | 74.11 | 73.04 | 88.15 | 70.11 |
| ? | sequelbox/StellarBright | | 74.1 | 72.95 | 87.82 | 71.17 |
| ? | Riiid/sheep-duck-llama-2-70b-v1.1 | | 74.07 | 73.04 | 87.81 | 70.84 |

https://huggingface.co/spaces/AtmosBank/Atmos_Leaderboard?logs=container

@rajistics

# Methods for evaluating Generative AI

- Exact matching approach
- Similarity approach
- Unit Tests
- Evaluation Benchmarks
- Human Evaluation
- Human Comparison/Arena
- Model based Approaches
- Red Teaming



Generative AI Evaluation Methods

Flexibility (higher covers more varied use cases)

Overall Cost (Higher is more expensive)

Red Teaming
Human Evaluation
Model based Approaches
Human Comparison/Arena
Unit Tests
Evaluation Benchmarks
Similarity approach
Exact matching approach

@rajistics

Raj guess

# Evaluating Code: Python

```
def incr_list(l: list):
    """Return list with elements incremented by 1.

    >>> incr_list([1, 2, 3]) [2, 3, 4]
    >>> incr_list([5, 3, 5, 2]) [6, 4, 6, 3]"""
```

Candidate solution:

```
return [(e + 1) for e in l]
```

Reference solution:

```
updated_list = [x+1 for x in l]
return updated_list
```

Leandro &
https://arxiv.org/pdf/2107.03374.pdf

@rajistics

# Evaluating Code with Unit Test

*Candidate solution:*

```python
def incr_list(l: list):
    """Return list with elements incremented by 1.

    >>> incr_list([1, 2, 3]) [2, 3, 4]
    >>> incr_list([5, 3, 5, 2]) [6, 4, 6, 3]"""

    return [(e + 1) for e in l]
```
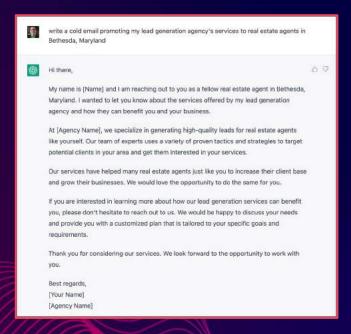
*Test Cases:*

```python
def check(candidate):
    assert candidate([]) == []
    assert candidate([3, 2, 1]) == [4, 3, 2]
    assert candidate([9, 0, 123]) == [10, 1, 124]
```

**Pass: yes/no**

@rajistics

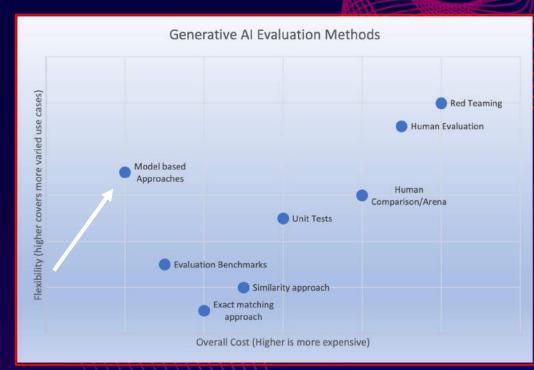https://arxiv.org/pdf/2107.03374.pdf

# Unit Tests Beyond Code



## Properties of Emails?

- First/Last Name?

- Grammar/spelling

- Concise?

- Verify actions?

- Tone? - is it polite

# Methods for evaluating Generative AI

- Exact matching approach
- Similarity approach
- Unit Tests
- Evaluation Benchmarks
- Human Evaluation
- Human Comparison/Arena
- Model based Approaches
- Red Teaming



Generative AI Evaluation Methods

Raj guess

# Model based evaluation



### Task instruction, sample, and question
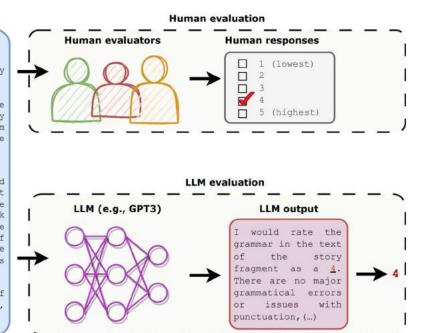
**Please rate the story fragment**

The goal of this task is to rate story fragments.

**NOTE:** Please take the time to **fully read** and **understand** the story fragment. **We will reject** submissions from workers that are clearly spamming the task.

**Story fragment**
The human ambassador reached down and grasped it's paw. "Humans, cats, is it true that all intelligent beings are omnivorous?" "Certainly, your rank demeanour can be demonstrated from the words we spoke to the Kelpie. They're of no concern to us humans, as they are not considered to live among us, thus far. (…)

How **grammatically correct** is the text of the story fragment? (on a scale of 1-5, with 1 being the lowest?)

### Human evaluation

**Human evaluators**   **Human responses**

☐ 1 (lowest)
☐ 2
☐ 3
☑ 4
☐ 5 (highest)

### LLM evaluation

**LLM (e.g., GPT3)**   **LLM output**

I would rate the grammar in the text of the story fragment as a 4. There are no major grammatical errors or issues with punctuation, (…)

4

@rajistics

https://arxiv.org/pdf/2305.01937.pdf

# C'mon Man - This isn't going to work



Bharat Saxena · 1st    2d •••
Bringing intelligence to Mainframes @ BMC Software | Explainable AI (XAI) | NLP …

Rajiv Shah From personal experience, I am a big skeptic when it comes to using another model as an evaluator … Hopefully you will be able to share some details from your presentation as some time in future.

# It works in Texas

**Texas is replacing thousands of human exam graders with AI** / Don't call the 'automated scoring engine' AI, though. They don't like that.

Query: Show me the total population of each state ordered from the most northern one to the most southern one.

**Gold Standard**

```
SELECT SUM(POP10), cbsa.state_name
FROM county
JOIN cbsa
ON county.geoid = CONCAT(CBSA.fips_state_code,CBSA.FIPS_COUNTY_CODE)
GROUP BY state_name
ORDER BY max(INTPTLAT) desc;
```

❌ Extra Column for Latitude

```
SELECT SUM(POP10), cbsa.state_name, max(INTPTLAT)
FROM county
JOIN cbsa
ON county.geoid = CONCAT(CBSA.fips_state_code,CBSA.FIPS_COUNTY_CODE)
GROUP BY state_name
ORDER BY max(INTPTLAT) desc;
```

@rajistics

https://medium.com/snowflake/inside-snowflake-building-the-most-powerful-sql-llm-in-the-world-1a33b3ee0d37

# Model based evaluation - Text->SQL

📄 Define Data Quality

🎹 Grading Scale

📥 Explain Inputs

You are a data analyst quality rater responsible for evaluating the quality of Snowflake SQL statements generated from natural language queries by comparing the candidate SQL statement to the user intent.

You must provide a score on an integer scale of 0 to 3 with the following meanings:
- 3 = perfect match - The candidate SQL will produce the same result as the user intent.
- 2 = good match - The candidate SQL will produce nearly the same result as the user intent but may suffer from non-deterministic issues such as sorting or grouping.
- 1 = partial match - The candidate SQL will produce an output similar to the user intent but may miss some part of the user's desired output
- 0 = no match - The candidate SQL will not produce anything similar to the user intent.

You will have access to the following elements:
1. User Query: The user natural language query enclosed in [{query}].
2. Database Schema: Information about the database schema is enclosed in [{db_formatted}].
3. Candidate SQL: The SQL query generated by the system is enclosed in [{candidate_sql}].

🌨 @rajistics

https://medium.com/snowflake/inside-snowflake-building-the-most-powerful-sql-llm-in-the-world-1a33b3ee0d37

# **Model based evaluation - Text->SQL**

Returned SQL is a bit different:

TX versus TEXAS

| User Query | Gold SQL | Candidate SQL |
|---|---|---|
| What are the first 100 tract IDs whose total ratio is the highest, except those in Texas? | SELECT tract FROM zip_tract WHERE usps_zip_pref_city NOT IN ('TX') ORDER BY tot_ratio DESC limit 100; | SELECT tract FROM zip_tract WHERE usps_zip_pref_city NOT IN ('TEXAS') ORDER BY tot_ratio DESC limit 100; |

# Model based evaluation - Text->SQL

Similarity:
❌ not a
100% match

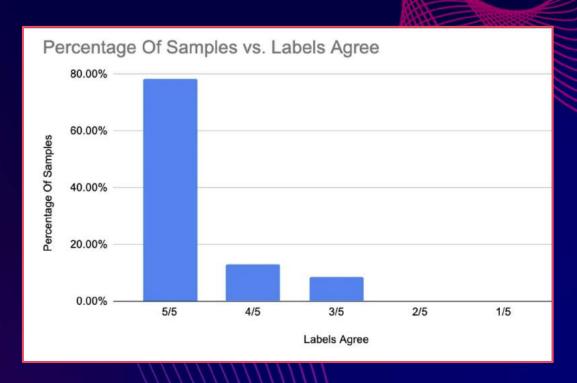| User Query | Gold SQL | Candidate SQL | Execution Accuracy |
|---|---|---|---|
| What are the first 100 tract IDs whose total ratio is the highest, except those in Texas? | SELECT tract FROM zip_tract WHERE usps_zip_pref_city NOT IN ('TX') ORDER BY tot_ratio DESC limit 100; | SELECT tract FROM zip_tract WHERE usps_zip_pref_city NOT IN ('TEXAS') ORDER BY tot_ratio DESC limit 100; | No Match |

https://medium.com/snowflake/inside-snowflake-building-the-most-powerful-sql-llm-in-the-world-1a33b3ee0d37

# Model based evaluation - Text->SQL

Model:

It's ok, its captures the users intent

| User Query | Gold SQL | Candidate SQL | Execution Accuracy | Execution Score |
|---|---|---|---|---|
| What are the first 100 tract IDs whose total ratio is the highest, except those in Texas? | SELECT tract FROM zip_tract WHERE usps_zip_pref_city NOT IN ('TX') ORDER BY tot_ratio DESC limit 100; | SELECT tract FROM zip_tract WHERE usps_zip_pref_city NOT IN ('TEXAS') ORDER BY tot_ratio DESC limit 100; | No Match | Perfect Match |

https://medium.com/snowflake/inside-snowflake-building-the-most-powerful-sql-llm-in-the-world-1a33b3ee0d37

# Model based evaluation - Text->SQL

Strong correlation

to other
evaluation
approaches



Percentage Of Samples vs. Labels Agree

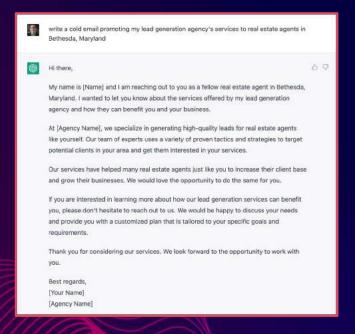https://medium.com/snowflake/inside-snowflake-building-the-most-powerful-sql-llm-in-the-world-1a33b3ee0d37

# Model evaluation aligns with humans

Human and GPT-4 judges can reach above 80% agreement on the correctness and readability score.



Human vs GPT-4 Grading Alignments for GPT-3.5 answers

Diff = 2 · Diff = 1 · Same Score

Correctness: 88%
Readability: 95%
Comprehensiveness: 72%

https://arxiv.org/abs/2305.01937
https://www.databricks.com/blog/LLM-auto-eval-best-practices-RAG
https://arxiv.org/abs/2303.16634
https://arxiv.org/pdf/2306.05685.pdf

@rajistics

# Unit Tests Beyond Code



## Properties of Emails?

- First/Last Name?

- Grammar/spelling

- Concise?

- Verify actions?

- Tone? - is it polite

+ *Explanations*

@rajistics

# We can do this!

- Exact matching approach
- Similarity approach
- Unit Tests
- Evaluation Benchmarks
- Human Evaluation
- Human Comparison/Arena
- Model based Approaches
- Red Teaming

@rajistics

https://unsplash.com/photos/man-in-blue-denim-jacket-standing-on-road-during-daytime-qCJzjVHkERc

# We can do this!

- Exact matching approach
- Similarity approach
- Unit Tests
- Evaluation Benchmarks
- Human Evaluation
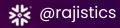- Human Comparison/Arena
- Model based Approaches
- Red Teaming

@rajistics

https://unsplash.com/photos/man-in-blue-denim-jacket-standing-on-road-during-daytime-qCJzjVHkERc

# Evaluation Techniques for LLMs

https://github.com/rajshah4/LLM-Evaluation

Rajiv Shah
@rajistics
r.shah@snowflake.com

April 2024