

LLM Evaluation & Observability

Avinash Joshi Product Lead and Director, TruEra (India)

arer

truera

TruEra was founded in 2014 by Anupam Datta, Somesh Jha, and Tomer Geva. Some key details about the founders.

- Anupam Datta Professor of Computer Science and Electrical and Computer Engineering at Carnegie Mellon University. He served as TruEra's President and Chief Scientist.
- Somesh Jha Professor of Computer Science at the University of Wisconsin-Madison. He was the Vice President of Research at TruEra.
- Consider LLMs hallucinatory... unless proven otherwise! Tomer Geva - Previously worked at Microsoft Research TruEra's CEO.

The three foundam

mutra was acquired by SAS, the analytics software company, for an undisclosed sum. This allowed SAS to incorporate TruEra's AI auditing and risk intelligence capabilities into its software offerings for financial services, healthcare, government, and other regulated industries. While the TruEra brand was absorbed into SAS, the founders

This overlap is really murky

ML research And actively optimizes models for penalizes **Generalization**

Focus LLMs on 'General' Tasks

Generalization

Memorization

✓ Summarization

- Text Embedding
- Inference
- Planning

Leave memorization to something else



LLMs Need a Knowledge Source





Enter Retrieval Augmented Generation (RAGs)

Example: Question Answering ChatBot



But RAGs can Hallucinate too



Input	^
Who is Shayak?	
Response	^
Shayak is a computer scientist who obtained his PhD in Computer Science from Carnegie Mellon University and BTech in Comput	er
Science from the Indian Institute of Technology, Delhi. He has been building systems and leading research to make machine	
learning and big data systems more explainable, privacy compliant, and fair for the past 10 years. He is also a member of the Ban	k
of England's AI Public-Private Forum and the OECD Global Partnership on AI.	



For second meetings start here

8



Gen Al is Taking over the World

The World is Divided in Two Camps

Teams that are building cool prototypes

• Very easy to build apps that do the right thing 80% of the time.

Teams that are deploying apps in production

• Much harder to go beyond that

Three Examples from Camp #2

Education Tech

Student support for assignments

Customer Support

Summarization of call transcripts with action items Financial Services

Customer facing chatbot



What Sets these Camps Apart?

Systematic iteration to improve LLM apps in development and in production



GenAl Observability Platform Requirements



iruera

TruLens : Track and Evaluate LLM Experiments

github.com/truera/trulens





TruLens has Broad Support for LLM Apps

Broad App Support

- RAGs
- Advanced RAGs (query planning, sentence window, auto merging)
- Agents
- Fine-tuning experiments

Integration with Leading FrameworksLlamaIndexLangChainCustom Frameworks



Honest, Harmless, Helpful Evals with TruLens

A feedback function provides a score after reviewing an LLM app's inputs, outputs, intermediate results, and metadata.

Honest

...

- Answer relevance
- Embedding distance
- Context relevance
- Groundedness
- Summarization quality
- Custom evaluations

Harmless

- PII Detection
- Toxicity
- Stereotyping
- Jailbreaks
- Custom evaluations
- ...

Helpful

- Prompt sentiment
- Language mismatch
- Conciseness
- Coherence
- Custom evaluations

• ...

Track Experiments to Select the Best App Configuration

Constructing the Vector DB

- Data Selection
- Chunk Size & Chunk Overlap
- Index Distance Metric
- Selection of Embeddings

Retrieval

- Amount of Context Retrieved (top k)
- Query Planning
- Sentence window retrieval
- Auto merging retrieval

LLM

- Prompting
- Model choice
- Model Parameters (size, temperature, frequency penalty, model retries, context length etc.)



trulens ×							
	Chain1	_WikipediaQA					
board	Records	Average Latency (Seconds)	Total Cost (USD)	Total Tokens	qs_relevance	relevance	Select App
ions	5	2.4	\$0.01	6.38k	0.48	1.0	
					Low	High	
10.4.0							
	Chain2	_WikipediaQA					
	Records	Average Latency (Seconds)	Total Cost (USD)	Total Tokens	qs_relevance	relevance	Select App
	5	2.4	\$0.01	6.38k	0.48	1.0	
					Low	Migh	
	Chain3	_WikipediaQA					
	Records	Average Latency (Seconds)	Total Cost (USD)	Total Tokens	qs_relevance	relevance	Select App
	5	1.8	\$0.01	6.38k	0.47	1.0	
					Low	🖌 High	
	Chain4	WikipediaOA					
	Records	Average Latency (Seconds)	Total Cost (USD)	Total Tokens	qs_relevance	relevance	Salact Ann
	5	1	\$0	6.7k	0.47	0.82	Select App
					Low	🖌 High	
ptimal	Chain5	_WikipediaQA					
adal	Records	Average Latency (Seconds)	Total Cost (USD)	Total Tokens	qs_relevance	relevance	Select App
odel	5	1.4	Ş0.01	6.38k	0.72	0.82	
					🖌 High	Migh	

Choosing the Right Evals





LLM Evaluation & Experiment Tracking



TruEra supports requirements for effective evaluation & tracking

- Broad app support
- Reliable, comprehensive, extensible evals
- Experiment tracking to select best app config
- Easy integrations with emerging tech stack

From TruLens to TruEra Full Lifecycle AI Observability



LLM Observability across the Lifecycle



TruEra Supports GenAl Observability Requirements







Retrieval Failure



Groundedness

~ Who is Shayak?

Response

Input

Shayak is a computer scientist who obtained his PhD in Computer Science from Carnegie Mellon University and BTech in Computer Science from the Indian Institute of Technology, Delhi. He has been building systems and leading research to make machine learning and big data systems more explainable, privacy compliant, and fair for the past 10 years. He is also a member of the Bank of England's AI Public-Private Forum and the OECD Global Partnership on AI.

Feedback

qs relevance = 0.55 question statement result Who is Shayak? When Shayak started building production grade machine learning models for algorith 0 1.0 1 Who is Shayak? When Shayak started building production grade machine learning models for algorith 1.0 2 Who is Shayak? Most recently, Shameek was Group Chief Data Officer at Standard Chartered Bank, w 0.1 3 Who is Shayak? Shameek has spent most of his career in driving responsible adoption of data analyti 0.1

~

~

~

~

Lack of Groundedness

TruEra RAG Triad Query Answer Relevance Response Context Context

Selected LLM Application: Summarize_v1

Selected Record ID: record_hash_b63e7be3dac4d43a3d79178e22710b4c

Input[Select.RecordInput]

#Person1#: Hello, Is this room service? #Person2#: Yes, May I help you? #Person1#: This is the room 1425, we asked for the room service an hour ago. #Person2#: we are very sorry to cause you a lot of inconvenience. #Person1#: What's the matter? #Person2#: We're rather busy right now. It will take another 15minutes #Person1#: Is it really going to take that long, will you rush the order? #Person2#: I'm afraid it will take 15 minutes at most. #Person1#: ha, well, we have no choice.

Response [Select.RecordOutput]

Room 1425 called room service an hour ago and they were told it will take an additional 15 minutes for their order to arrive. They are not happy about the wait but have no other option.

statement result reason 0 rson1#: Th Room 1425 called room service an hour ago and they were told it will take an addition 0 Statement Sentence: Room 1425 called room service an hour ago and they were told it will take an additional 15 minutes for their order to arrive., Supporting Evidence: This is the room 1425, we asked for the room service an hour ago. We're rather busy right now. It will take an other 15 minutes bert_score = None Statement Sentence: They are not happy about the wait, Supporting Evidence: This statement Sentence: NOTHING FOUND Score: 1 No feedback details. Statement Sentence: but have no other option., Supporting Evidence: ha, well, we have no choice. Score: 10	groundedness_measure = 0.7					
0 rson1#: Th Room 1425 called room service an hour ago and they were told it will take an addition 0 Statement Sentence: Room 1425 called room service an hour ago and they were told it will take an additional 15 minutes for their order to arrive., Supporting Evidence: This is the room 1425 wasked for the room service an hour ago. We're rather busy right now. It will take an addition use asked for the room service an hour ago. We're rather busy right now. It will take another 15 minutes Score: 10 bert_score = None Statement Sentence: The statement Sentence: The statement Sentence: They are not happy about the wait, Supporting Evidence: NOTHING FOUND Score: 1 No feedback details.			statement	resul	sult reason	
	0 Dert_sco No fee	r <mark>rson1#: Th</mark> ore = None dback de	Room 1425 called room service an hour ago and they were told it will take an additio	0.	 Statement Sentence: Room 1425 called room service an hour ago and t told it will take an additional 15 minutes for their order to arrive., Supporting Evidence: This is the room 1425, we asked for the room serv ago. We're rather busy right now. It will take another 15 minutes Score: 10 Statement Sentence: They are not happy about the wait, Supporting Evidence: NOTHING FOUND Score: 1 Statement Sentence: but have no other option., Supporting Evidence: hay well, we have no choice. Score: 10 	ney were ice an hour

Answering the Wrong Question

TruEra RAG Triad Query Answer Relevance Context Relevance Context Groundedness

Input A Which year was Hawaii's state song written?

Feedback

qs_relevance = 0.325 ^ question statement result 0 Which year was Hawaii's state song written? "Hawai` i Pono` i" is the state song of Hawaii. The words were written by King David K 1.0 1 Which year was Hawaii's state song written? The American business people made Hawaii into a republic for a short time. The new 0.1 2 Which year was Hawaii's state song written? Hawaii (sometimes spelled "Hawai'i".) is a U.S. state and the only U.S. State that is in 0.1 3 Which year was Hawaii's state song written? 1874 - Hawaii signs a treaty with the United States granting exclusive trade rights. 18 0.1

relevanc	e = 0.1		
	prompt	response	result
0	Which year was Hawaii's state song written?	Hawai`i Pono`i	0.1



Run Feedback Functions with Any Model



Connectors for 100+ LLMs



Run Evaluations on any LLM App Stack



iruera

Try Open Source TruLens





github.com/truera/trulens

Give us a star, request a feature, contribute!