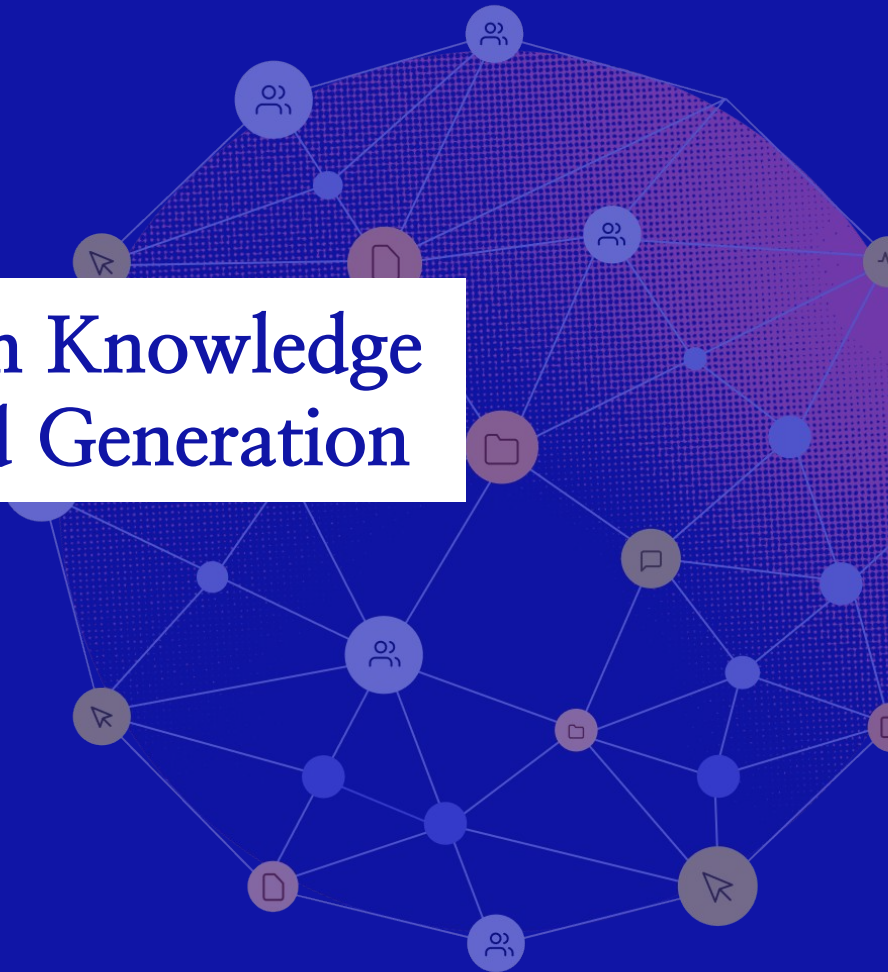


# Separating Coherence from Knowledge with Retrieval-Augmented Generation

WSAI

April 19, 2023



# Agenda

1. Introduction
2. ChatGPT: Language vs. Implicit World Knowledge
3. RAG (Retrieval Augmented Generation)
4. Hallucinations: Entailment

# Chau Tran

*Tech Lead, LLMs and Vector Search*  
Glean





# Language

# Today's LLMs

- Trained to model language
- Emergent abilities at scale:
  - Coherence
  - Reasoning
  - Implicit knowledge

The background features two large, overlapping, organic shapes. The shape on the left is a light green color with a fine halftone dot pattern. The shape on the right is a light orange color, also with a fine halftone dot pattern. Both shapes have irregular, wavy edges. The text is centered between these two shapes.

“Why is the sky…”

# Knowledge



# Knowledge

- Implicit world knowledge
  - Hallucinations
  - Factuality issues



# Knowledge

- Missing knowledge
  - Recency (ChatGPT → 2021)
  - Non-public data

# RAG

# Retrieval Augmented Generation (RAG)

- Retrieval
  - Store knowledge outside of LLM model weights
- Use generative LLMs for their strengths:
  - Reasoning
  - Synthesis



# Retrieval Augmented Generation (RAG)

- Results:
  - Knowledge-grounded
  - Factually generated output
  - LLMs as “agents” to connect to tools and knowledge

# RAG / Agent

The New England Journal of Medicine is a registered trademark of [QA("Who is the publisher of The New England Journal of Medicine?") → Massachusetts Medical Society] the MMS.

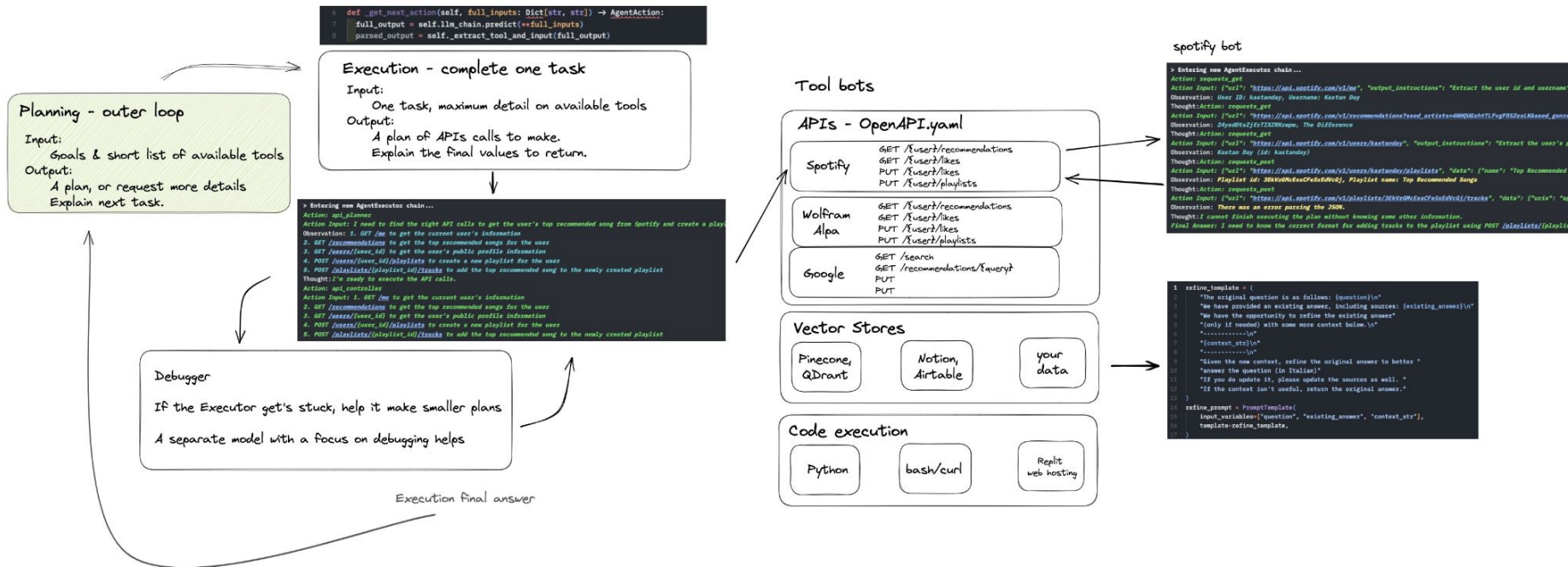
Out of 1400 participants, 400 (or [Calculator(400 / 1400) → 0.29] 29%) passed the test.

The name derives from "la tortuga", the Spanish word for [MT("tortuga") → turtle] turtle.

The Brown Act is California's law [WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.] that requires legislative bodies, like city councils, to hold their meetings open to the public.

Figure 1: Exemplary predictions of Toolformer. The model autonomously decides to call different APIs (from top to bottom: a question answering system, a calculator, a machine translation system, and a Wikipedia search engine) to obtain information that is useful for completing a piece of text.

# RAG / Agent





# Domain-Adapted Retrieval

- Retrieving the right knowledge is a search problem
- Hybrid search (SOTA)
  - Keywords & terms
  - Dense (embeddings)

Fine-tuned for specific  
corpora & knowledge

asset allocation model

Input text



Embedding  
model



-0.015

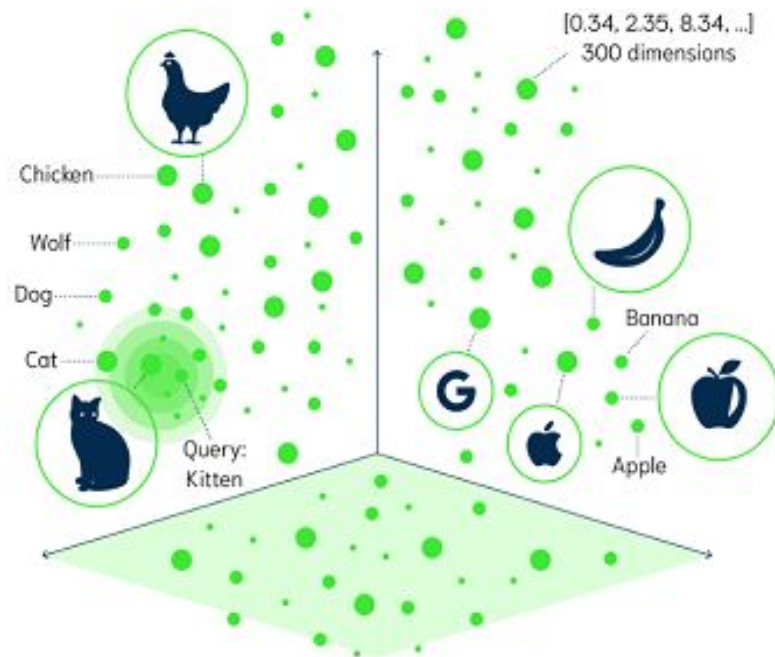
0.012

-0.014

...

-0.013

Numerical representation





# Entailment as a Hallucination Defense

- For LLM input X (relevant knowledge), is generated output Y entailed by X?
- Note: requires X to be source of knowledge (RAG)

## TEXT

- *Eyeing the huge market potential, currently led by Google, Yahoo took over search company Overture Services Inc last year.*

## HYPOTHESIS

- *Yahoo bought Overture.*

## ENTAILMENT

- **TRUE**

# Q&A

If we don't get to your question live, we'll still make sure to get you an answer.

A member of our team will follow up with you after today's event.

glean

# Thank you!

Interested in more information?

👉 [glean.com/get-a-demo](https://glean.com/get-a-demo)

