# CS224v

# Conversational Virtual Assistants with Deep Learning

# Lecture 1: Introduction

Monica Lam

# Focus of the Course

Foundation To Make LLMs Useful as a Trustworthy General Virtual Assistant

90+% Accuracy

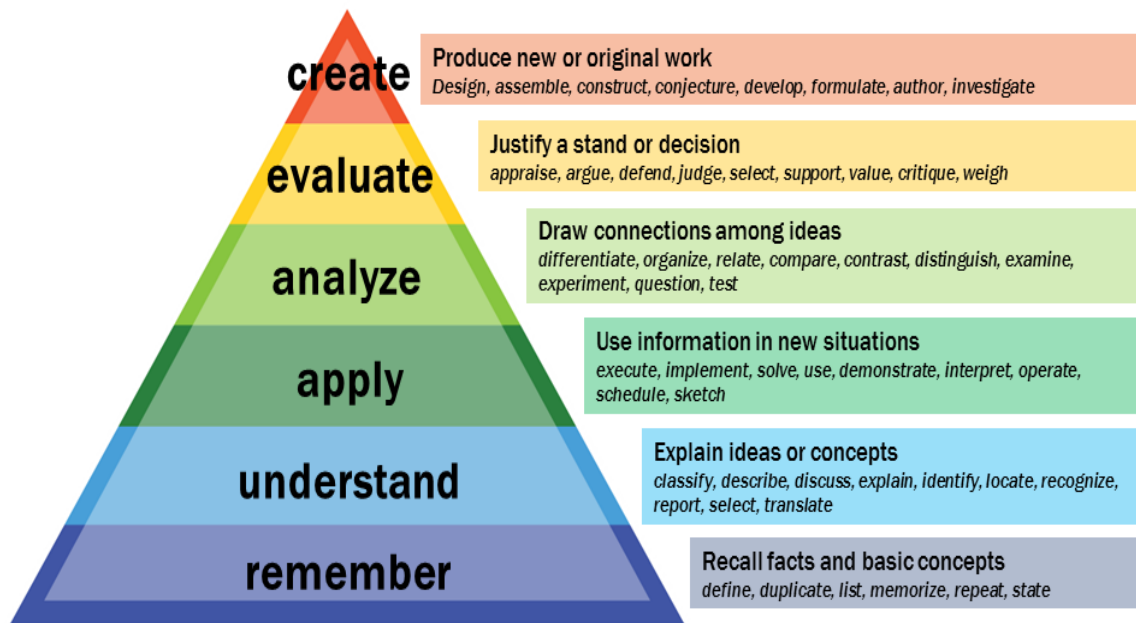Applicable to All Domains (for Non-AI Experts)

# A Project/Research Course

Actively Create and Advance New Technologies
**Adopted by Real Users**

Select Class Projects Continue as Research (2026)

Course Publications: 2023(2), 2024(5), 2025(6)

Secret to Success: **Working Software
Supports Increasingly Complex Projects**

# Goal: an AI-Based Research Assistant



Highest Level of Cognition
Subsumes all other Levels

Bloom's Taxonomy
Education Objectives in Cognition Domain

# Why an AI-Based Research Assistant?
## To Advance Knowledge Discovery

Faster

Better Results

Scalability

Broader (Breaking Down Domain & Language Barriers)
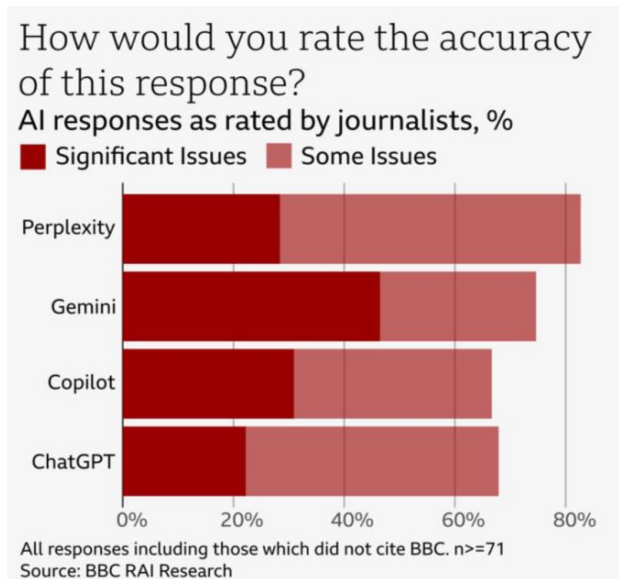
Research Accessible to More People

# Knowledge As We Know It

# Perplexity, Gemini, Copilot, ChatGPT Still Hallucinate With RAG
## (Retrieval Augmented Generation)

51% of AI answers to news problems have significant issues.

19% of AI answers which cited BBC content introduced factual errors.

13% of the quotes sourced from BBC articles were either altered or didn't actually exist in that article.

How would you rate the accuracy of this response?
AI responses as rated by journalists, %
■ Significant Issues  ■ Some Issues

All responses including those which did not cite BBC. n>=71
Source: BBC RAI Research

*Representation of BBC News content in AI Assistants. 2025.*

# The Problem

# The Enterprise AI Paradox

Enterprises are rushing to deploy AI,
but off-the-shelf LLMs are a ticking time bomb.

### Legal and Financial Liability

Air Canada forced to honor fictitious refund policy

### Reputational Damage

Lawyer sanctioned over ChatGPT-fabricated legal precedents

### Regulatory Risk

NYC's AI chatbot advised businesses to break laws

**The result: A crisis of trust that blocks AI adoption for high-value use cases.**

# The 70% Accuracy Problem

## It's Deceivingly Easy!

## Impossible to Get to Near 100% With Next-Word Prediction LLMs!

## Need of Human Filtering → Not Scalable!

# Fundamental Problems For Harder Tasks

LLMs Trained on Next-Word Prediction:

Not **Precise** Enough for Research

Cannot Perform **Complex** Tasks Reliably

Limited **Context Length**

# AUGMENT LLMs WITH
# **COMPUTATIONAL THINKING**

ALGORITHMS
DATA REPRESENTATION (PATTERN RECOGNITION)
ABSTRACTION (GENERALIZATION)
DECOMPOSITION

UNIVERSAL TO ALL DISCIPLINES

# Computational Thinking

Basic Concepts: Alan Perlis, Don Knuth (1950s)

Coined Term: Seymour Papert (1980)

Education Objective: Jeannette Wing, NSF (2006)

# Lecture Outline

1. Course Material Overview

   - Stage 1 (2022 -- 2025): Computational thinking → general research assistant
   - Stage 2 (2025 --      ): Computational thinking → scientific research assistant

2. Course overview
   - Project choice
     - Build on technology for a general assistant
     - Research the bleeding edge of a technical research assistant

   Just a high-level picture! Do not worry if you can't follow all the ideas
   (Details in the lectures)

# Genie: A Prototype of a
# Public AI Assistant to World Wide Knowledge (WWK)

Free Text

INTERNET
WIKIPEDIA — The Free Encyclopedia
SEMANTIC SCHOLAR

Knowledge Bases

WIKIDATA
FEDERAL ELECTION COMMISSION — UNITED STATES OF AMERICA
U.S. DEPARTMENT OF COMMERCE — BUREAU OF THE CENSUS
ACLED
311 SAN FRANCISCO AT YOUR SERVICE
BOSTON 311

Digitized Print & Handwriting

The African Times
THE SOUTH AFRICAN COMMERCIAL ADVERTISER

Genie

Browse

Search

Chat

Research (Storm)

# Live Demos: wwknowledge.org

# COMPUTATIONAL THINKING APPROACH

FORMALIZE HUMAN COGNITIVE PROCESS

PRODUCE STEP-BY-STEP INSTRUCTIONS

# 1. READING DOCUMENTS

## RAG-BASED LLMs HALLUCINATE

## ESPECIALLY WHEN RETRIEVED INFORMATION DOES NOT ANSWER THE USERS' QUESTION

# Genie: Automating Question Answering

**Ansatze** (an educated guess)

3. Ask LLM
4. Dissect into claims
5. Fact-check each claim
   - Search Wikipedia with claim
   - Filter incorrect claim

**Retrieve**

1. Search Wikipedia with query
2. Filter irrelevant info

6. Draft
7. Refine

Use multiple, easy LLM steps to consult external data

# WikiChat

- Accuracy: 97% (English)

- Best Research Award of the Year, Wikimedia Foundation

- Now in 25 languages: Breaking the language barrier!

| | | | | |
|---|---|---|---|---|
| Arabic | English | Indonesian | Persian | Serbian |
| Chinese | Finnish | Italian | Polish | Spanish |
| Czech | French | Japanese | Portuguese | Swedish |
| Danish | German | Korean | Romanian | Turkish |
| Dutch | Hebrew | Norwegian | Russian | Ukrainian |

*WikiChat: Stopping the Hallucination of Large Language Model Chatbots by Few-Shot Grounding on Wikipedia* Sina J. Semnani, Violet Z. Yao*, Heidi C. Zhang*, Monica S. Lam
In EMNLP 2023, Singapore, December 6-10, 2023.

# THE PROBLEM IS HARDER THAN EXPECTED

3 STUDENTS, 4 MONTHS FOR DEVELOPMENT
EVALUATION WAS ALSO HARD
THE DEVIL IS IN THE DETAILS

SOLID BUILDING BLOCKS ARE CRITICAL
APPLICABLE TO GENERAL COGNITIVE PROCESSES

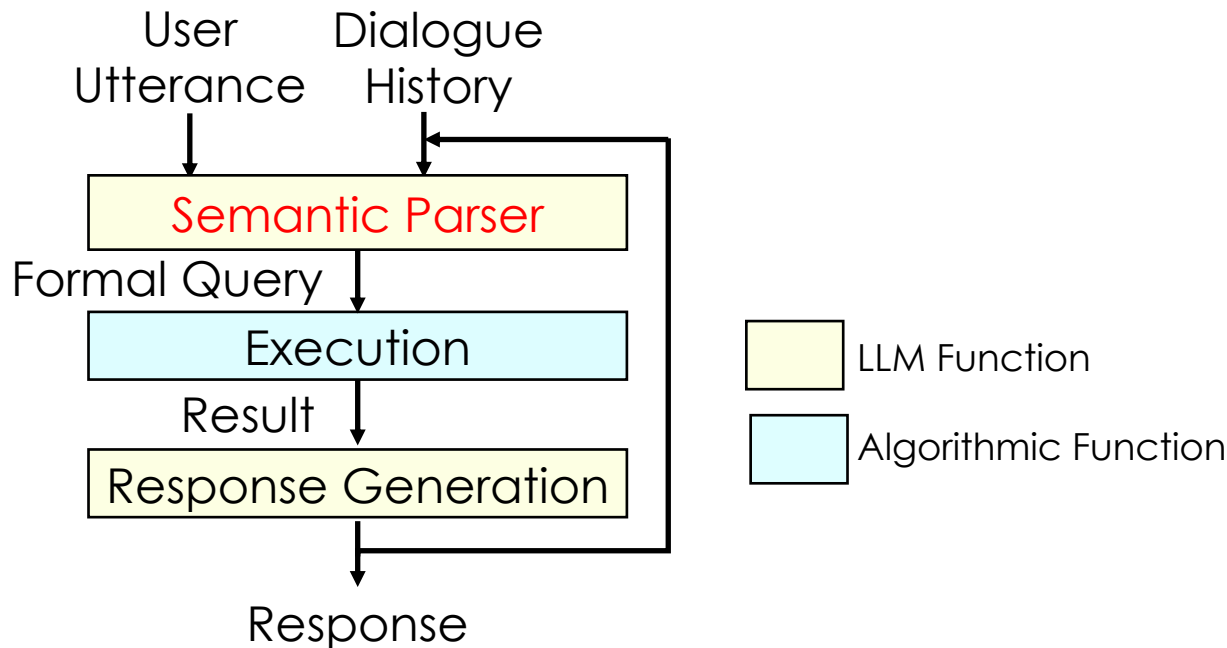# 2. Reading Data

# Example:

**WIKIDATA**

# World's Largest Live Knowledge Graph

- 15B facts, 100M entities, 10K properties, 25K contributors
- Wikidata could only be accessed by **SPARQL query language**
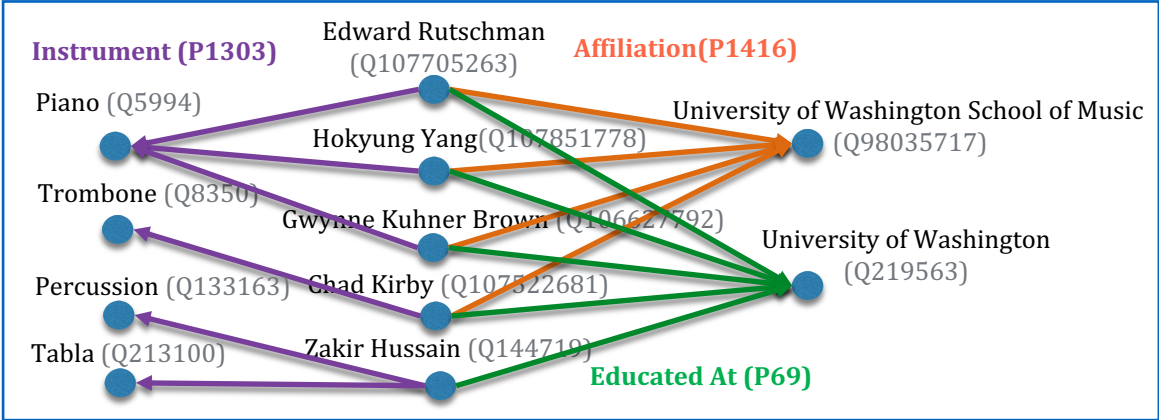
# Key Concept:
## A Semantic Parser Translates NL to Formal Semantics

What are the musical instruments played by people
who are affiliated with the U. of Washington School of Music
and have been educated at the U. of Washington, and how many people play each instrument?

Genie

```
SELECT ?instrument ?instrumentLabel (COUNT(?student) AS ?count) WHERE {
    ?student wdt:P1303 ?instrument;
        wdt:P1416 wd:Q98035717;
        wdt:P69 wd:Q219563.
    SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
}
GROUP BY ?instrument ?instrumentLabel
```



A Small Subset of Wikidata

| instrument | instrumentLabel | count |
|------------|-----------------|-------|
| Q5994      | piano           | 99    |
| Q1467960   | mbira           | 2     |
| Q8350      | trombone        | 11    |
| Q8338      | trumpet         | 8     |
| Q17172850  | voice           | 32    |
| ...        | ...             | ...   |
| Q302497    | mandolin        | 1     |
| Q187851    | recorder        | 1     |
| Q185041    | cor anglais     | 1     |
| Q83509     | piccolo         | 1     |

Result

# Break Down the Cognitive Process
## Go acquire knowledge to complete a task!

**Think**
**Act**

**Observe**
results

| User Query |
|---|

| Agent |
|---|

| SPARQL |
|---|

**Possible Actions**
**search_wikidata(**string**)**
    returns QIDs and PIDs
**get_wikidata_entry(**QID**)**
    returns Wikidata entity page
**execute_sparql(**SPARQL**)**
    returns SPARQL execution results
**get_property_examples(**PID**)**
    example of how properties are used
**stop()**

# Deployed at Wikidata Query Forum

https://www.wikidata.org/wiki/Wikidata:SPARQL_query_service/Wikidata_Query_Help
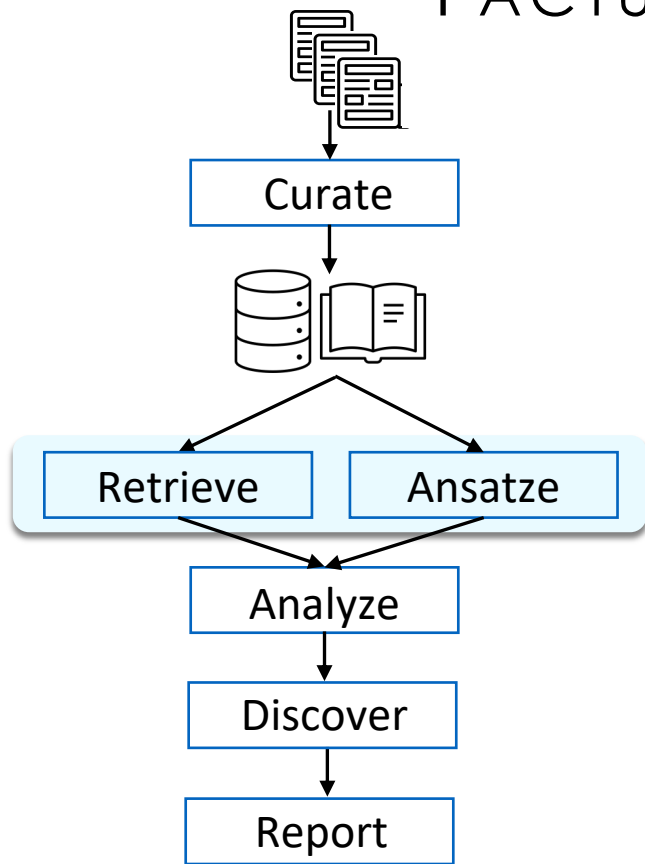
- Help users with their complex queries

  - 1700 conversations

  - In 198 samples, success rate: 78%

  - Interpretability:

    - Genie shows the query (in English) and its answer

# FACTUAL INFO RETRIEVAL

Input: Text, Databases, Knowledge Bases (Wikipedia, Wikidata, FEC ...)



**Research Pipeline**

Curate → (database/book) → Retrieve / Ansatze → Analyze → Discover → Report

| 2023/12 EMNLP | WikiChat: **Control hallucination** by grounding in **Wikipedia** <br> **- Best chat: achieving 97% factuality** <br> - Now grounded in Wikipedia in 25 languages <br> **- Best Research of the Year Award by Wikimedia Foundation, 2024** |
|---|---|
| 2024/ 6 NAACL | 1st assistant that supports **full hybrid queries on free-text + databases** <br> - Translate NL into our **novel SUQL query language** with optimizing compiler <br> **- Used by investigative journalists (FEC campaign donation), 311** |
| 2024/12 EMNLP | Best assistant for **complex queries** on **Wikidata** <br> - Agentic approach that learns about the schema on its own <br> **- Deployed on Wikidata** |
| 2025/ 7 ACL | Genie Worksheets: **Knowledge-Intensive Task-Oriented Agents** <br> - Accuracy: 80% vs 0-10% with GPT4 with function calling <br> **- Validation with industry in progress** |

# DISCOVERY OF GENERAL KNOWLEDGE



Input: Text, Databases, Knowledge Bases (Wikipedia, Wikidata, FEC ...)

| 2024/ 6 NAACL | STORM: 1st **deep research** assistant by searching the **Internet**<br>- Comprehensive Wikipedia-like articles with citations to full references<br>- 800K organic users, 1.4M articles written<br>- **Used by history students** to discover African history newspapers (1800s)<br>- **Inspired OpenAI Deep Research, Google Gemini Deep Research, Databricks Genie Deep Research Model** |
|---|---|
| 2024/12 EMNLP | 1st assistant that lets users **interactively direct the deep research**<br>- Uses a novel round-table of experts to explore unknown unknowns |

**New: DATASTORM discovers insights by combining Storm + DataTalk**

# Genie: A Prototype of a
# Public AI Assistant to World Wide Knowledge (WWK)

Demos: wwknowledge.org

**Free Text**

INTERNET

WIKIPEDIA
The Free Encyclopedia

SEMANTIC SCHOLAR

**Knowledge Bases**

WIKIDATA

FEDERAL ELECTION COMMISSION
UNITED STATES OF AMERICA

U.S. DEPARTMENT OF COMMERCE
BUREAU OF THE CENSUS

ACLED

311
SAN FRANCISCO
AT YOUR SERVICE

BOSTON 311

**Digitized Print & Handwriting**

The African Times.

THE SOUTH AFRICAN COMMERCIAL ADVERTISER.

**Genie**

Browse

Search

Chat

Research (Storm)

NLP Research Can Create Useful Tools

Usage Also Helps Generate Research Ideas

# Lecture Outline

1. Course Material Overview

   - Stage 1 (2022 -- 2025): Computational thinking → general research assistant
   - Stage 2 (2025 --        ): Computational thinking → scientific research assistant

2. Course overview
   - Project choice
     - Build on technology for a general assistant
     - Research the bleeding edge of a technical research assistant

# Scientific Research
## Requires Changing the Fundamentals of NLP
## With Computation Training (CT)

1. Information Retrieval (Embedding Similarity)
   → Formal Semantics

2. Reasoning: Chain of Thought
   → Computational Thinking (CT)

3. Training: Unsupervised Next-word Prediction Training
   Fine-tuned with Chain of Thought
   → Fine-tuned with CT Training

# Knowledge in Large Sets of Long Documents

**Analysis** combines info from every doc

Decision, Stats, Trend, ...

| | Resumes<br>Paper Reviews<br>Proposals<br>Financial Reports | Pubmed Articles<br>Human Genome<br>Drug/Protein/Disease<br>Event News | News<br>Social Media<br>Interviews<br>Medical Rec. | Requirements<br>Regulations<br>Clinical Trials<br>Invoice Match |
|---|---|---|---|---|
| | Ranking (Top n) of ALL documents | Accumulation | Stats/trends | Satisfiability |
| | Finalists | Knowledge Bases | Discover Knowledge | Compliance to Constraints |

# An Example: Clinical Trials

- Motivation
  - Patients cannot navigate clinical trial websites
  - Clinical trials fail for a lack of finding patients

- Challenges: matching millions of trials with millions of patients
  - Informal Retrieval (IR)
    - Similarity → false positives; false negatives
  - Matching
    - LLM "reasoning" → inaccurate

  <span style="color:red">Missing one appropriate trial is fatal!</span>

# Approach for Effective Clinical Trial Matching

1. Representation
   - Translate free-text trials/patient records to formal semantics
     - SMT (Satisfiability Modulo Theories)
     - UMLS (Universal Medical Language System) with 3M terms
2. Retrieval algorithm
   - Store trials in DBs with columns representing simpler predicates
   - Retrieve with DB queries
3. Match algorithm: SMT theorem prover

# Curating Knowledge

**Semantic Parsing from Free-Text to:**

- **Databases**:

  for analyzing large & long document sets

  - Sets of papers, articles, reports, transcripts

- **Knowledge graphs**: a semantic web

  - NSF Proto-OKN project

- **SMT: constraint satisfaction**

  - Degree satisfaction, regulation compliance, invoice matching



**Full Research Pipeline**

# Knowledge Curation—Important but Hard Example:



ACLED: Armed Conflict Location & Event Data

- Best dataset of qualitative coding (Free-text → DB)
- Human annotators, with rounds of reviews
- Over 200 territories, 80 different languages

# 🍋 LEMONADE: A Large Multilingual Expert-Annotated Abstractive Event Dataset for the Real World

**Sina J. Semnani**[1]    **Pingyue Zhang**[2]    **Wanyue Zhai**[1]    **Haozhuo Li**[1]
**Ryan Beauchamp**[1]    **Trey Billing**[3]    **Katayoun Kishi**[3]    **Manling Li**[2]    **Monica S. Lam**[1]
[1]Stanford University  [2]Northwestern University  [3]ACLED

In Findings of ACL, 2025

## ACADEMIC FORMULATION FALLS SHORT

EXTRACTIVE → ABSTRACTIVE

CANONICALIZATION OF ENTITIES

NEEDED: RESEARCH ON ENTITIES IN EXTERNAL LITERATURE

SOTA, BUT NOT ACCURATE ENOUGH

# Goal: Create A Universal Semantic Parser

*Free Text, Sample Annotations ➔ Formal Representation*

1. Ontology: Automatic standardization of terminology

2. Translation to formal semantics. Challenges:
   - Classification among many entities
   - Acquisition of missing knowledge
   - Long documents exceeding LLM contexts
   - Consistency checking
   - Refinement with annotated samples

3. Information retrieval optimizations

# 3 Case Studies of Scientific Assistants

All the technology developed is generally applicable to any domain
You can join these new projects from the ground up.

1.  **Clinical Trial Matching** (Mayo Clinic, Stanford Medicine) using **SMT**
    *   Off-line matching & a conversational agent for patients looking for trials

2.  **Cancer drug resistance**: Data Synthesis with literature in RNA-Seq Analysis (Stanford Oncology)
    *   Analysis of experimental results
    *   Extraction of **knowledge graph** of pubmed articles
    *   Synthesizing data with literature

    Similar use case: **Repurposing drugs for rare diseases** (Everycure, Rare Genomics Inst.)

3.  **Sustainability Modeling** (Stanford Cornerstone Initiative)
    *   Synthesizing model from many different raw **databases** with judgement

# Scientific Research
## Requires Changing the Fundamentals of NLP
## With Computation Training (CT)

1. Information Retrieval (Embedding Similarity)
   → Formal Semantics

2. **Reasoning: Chain of Thought**
   → Computational Thinking (CT)

3. Training: Unsupervised Next-word Prediction Training
   Fine-tuned with Chain of Thought
   → Fine-tuned with CT Training

# LLMs Lack Computational Thinking (CT)
Example: Composition

| Question | OpenAI-03 Model |
|---|---|
| Who is the wife of Benjamin Harrison? | Caroline Harrison |
| Who is the grandfather of Caroline Harrison? | George Scott |
| Who's the grandfather of the wife of Benjamin Harrison? | Dr. John Witherspoon |

## Idea

Do not use LLMs to answer complex questions directly

We propose:

1. Use LLM to perform simple functions

2. Use an Algorithmic Engine to invoke LLM functions and compose them

# LLMs Lack Computational Thinking (CT)

Example: Composition

Question                                             OpenAI-O3 Model

Who is the wife of Benjamin Harrison?                 Caroline Harrison

Who is the grandfather of Caroline Harrison?          George Scott

Who's the grandfather of the wife of Benjamin Harrison?



**Computation Thinking**

Input

*Algorithmic* Recursion & Decomposition Framework

Problem statement

Primitive?

Solve-Primitive
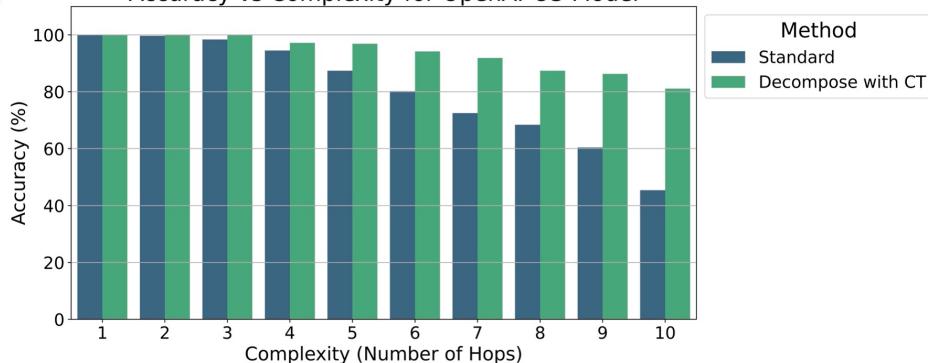
Decompose

output

LLM Call



Accuracy vs Complexity for OpenAI-O3 Model

Method
Standard
Decompose with CT

Dataset: LLM knows every hop of the question

# Potential Applications

- Semantic Parser

  This just in: Beats SOTA on Spider-2 Text-to-SQL

- Data Synthesis

- Coding

- Scientific and Engineering Workflows

# Scientific Research
## Requires Changing the Fundamentals of NLP With Computation Training (CT)

1. Information Retrieval (Embedding Similarity)
   → Formal Semantics

2. Reasoning: Chain of Thought
   → Computational Thinking (CT)

3. **Training: Unsupervised Next-word Prediction Training Fine-tuned with Chain of Thought**
   → Fine-tuned with CT Training

# Training CT

- Why?
  - To improve ansatze (educated guess) to aid knowledge discovery
    - Example: Ask LLM "What drug may cure this rare disease?" Literature search on the answer.
- How?
  - Fine-tune with training samples:
    (Natural language question → Formal Query; Answer)
  - Teach it formal semantics representation
  - Improve its latent representation of knowledge in LLMs

- Expected result: improved LLMs recall and ansatzes

# The CS224V Course

# Course Objectives

Foundation To Make LLMs Useful as a Trustworthy General Virtual Assistant

Hands-on Project/Research Experience Building on the State-of-the-Art Tools Or Advancing the State-of-the-Art

# Course Design

1. 2 homeworks to bring everybody onboard with SOTA tools
2. Lectures on techniques of LLM-based conversational agents
3. Supervised quarter-long project
   - Create a new application on existing tools
   - Enhance existing tools (with an application)
   - Develop new tools  (with an application)

# Purpose of the 2 Assignments

Prepare you for your project proposal

1. How to create an autonomous LLM-based research agent?
2. How to create a conversational agent using Genie Worksheet
   - Performs tasks; Answers question
   - Does not hallucinate

# Project Apprenticeship

- Assistance with project selection: Hardest part in research!
  - We suggest over 20 research-oriented projects on the website
  - Student-initiated projects are hugely welcome

- Weekly group mentorship meeting
  - We want to make you succeed!

# Project Mentorship

All homeworks and projects are to be done in pairs

- Week 4: Project proposal, with a weekly plan
- Weeks 5-10 (excluding Thanksgiving break):
  - Submit a written weekend update (every Monday)
  - Group meeting with mentors during the week
- Week 11: Poster presentation (Dec 4)
- Final project report due Dec 9, 2025.

# Course Schedule at a Glance

| Dates | Lectures / Homeworks | Projects |
|---|---|---|
| 9/22 - 10/ 6 | Introduction;<br>Autonomous Research Agents (HW1)<br>Task-oriented agents (HW2)<br>Grounding on free text | Research Project Ideas |
| 10/ 8 - 10/20 | | Student-initiated ideas<br>Project discussions<br>Project proposals (2) |
| 10/22 - 11/3 | Knowledge retrieval: databases, knowledge graphs, hybrid data, long documents | Weekly meetings with mentor |
| 11/ 5  - 11/12 | Reasoning with computational thinking;<br>Formal semantics;  NLP building blocks | Weekly meetings with mentor |
| 11/17 – 11/19 | Misc: Multimodal apps; Training LLMs | Weekly meetings with mentor |
| 11/25 - 11/27 | *Thanksgiving* | |
| 12/3 | | Final project posters<br>(3:00-5:40) |

# This Course

| | Grade |
|---|---|
| Participation | 15% |
| Assignment | 25% |
| Final Project | 60% |

Participation includes

- Class attendance and participation
- Ed discussion
- Meetings with project mentors