# CME 295: Transformers &
# Large Language Models
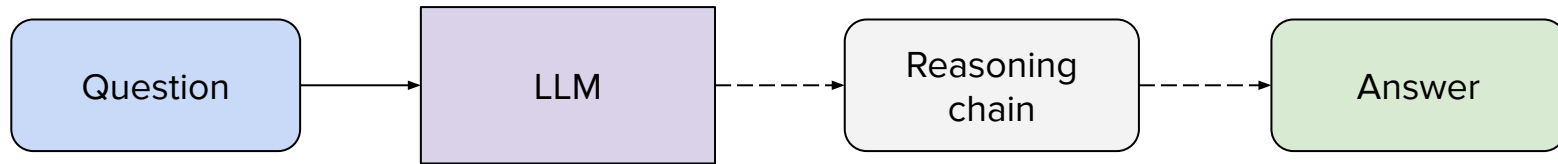
**Afshine Amidi** & **Shervine Amidi**
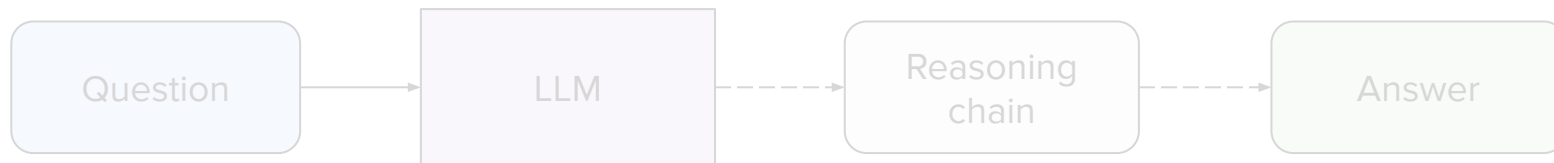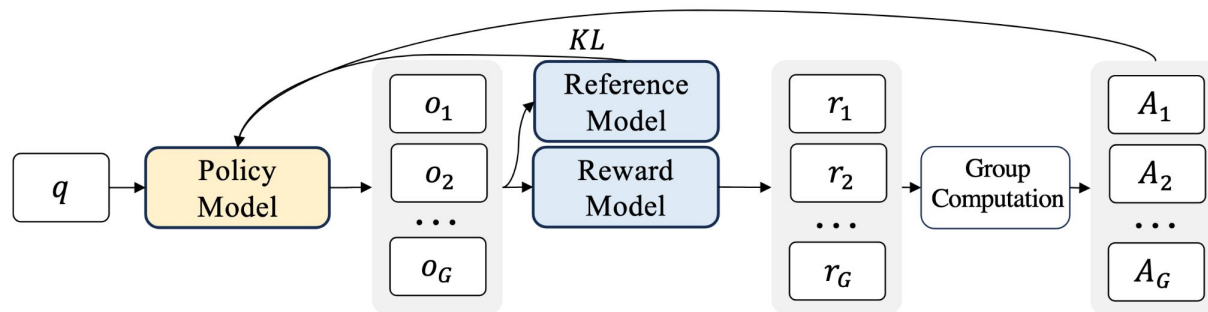
**Reasoning models**

# Recap of last episode…

**Reasoning models**



GRPO

Stanford University

ICME

# Recap of last episode...



DeepSeek-R1-Zero AIME accuracy during training

- r1-zero-pass@1
- r1-zero-cons@16
- o1-0912-pass@1
- o1-0912-cons@64



DeepSeek-R1-Zero average length per response during training

*"DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning", DeepSeek-AI, 2025.*

# Recap of last episode…



DeepSeek-R1-Zero AIME accuracy during training

DeepSeek-R1-Zero average length per response during training

Original $\quad \dfrac{1}{G}\sum_{i=1}^{G}\dfrac{\color{red}1}{\color{red}|\mathbf{o}_i|}\sum_{t=1}^{|\mathbf{o}_i|}$

*"DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning", DeepSeek-AI, 2025.*

ICME

DeepSeek-R1-Zero AIME accuracy during training


DeepSeek-R1-Zero average length per response during training

Original

$$\frac{1}{G}\sum_{i=1}^{G}\color{red}{\frac{1}{|\mathbf{o}_i|}}\color{black}\sum_{t=1}^{|\mathbf{o}_i|}$$

DAPO

$$\frac{1}{\sum_{i=1}^{G}|o_i|}\sum_{i=1}^{G}\sum_{t=1}^{|o_i|}$$

Dr. GRPO

$$\frac{1}{G}\sum_{i=1}^{G}\sum_{t=1}^{|\mathbf{o}_i|}$$

*"Understanding R1-Zero-Like Training: A Critical Perspective"*, Liu et al., 2025.
*"DAPO: An Open-Source LLM Reinforcement Learning System at Scale"*, Yu et al., 2025.

Stanford University

ICME

# Recap of last episode…

**Strengths**.

- Great at imitation or idea generation
- Amazing at generating or debugging code

**Weaknesses**.

- Limited reasoning
- Knowledge is static
- Cannot perform actions
- Hard to evaluate

ICME

**Strengths**.

- Great at imitation or idea generation
- Amazing at generating or debugging code

**Weaknesses**.

- Limited reasoning  ⟵  Focus of last lecture
- Knowledge is static
- Cannot perform actions
- Hard to evaluate

ICME

**Strengths**.

- Great at imitation or idea generation
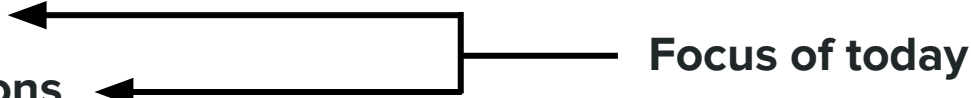- Amazing at generating or debugging code

**Weaknesses**.

- Limited reasoning
- **Knowledge is static** ←
- **Cannot perform actions** ←

  **Focus of today**

- Hard to evaluate

ICME

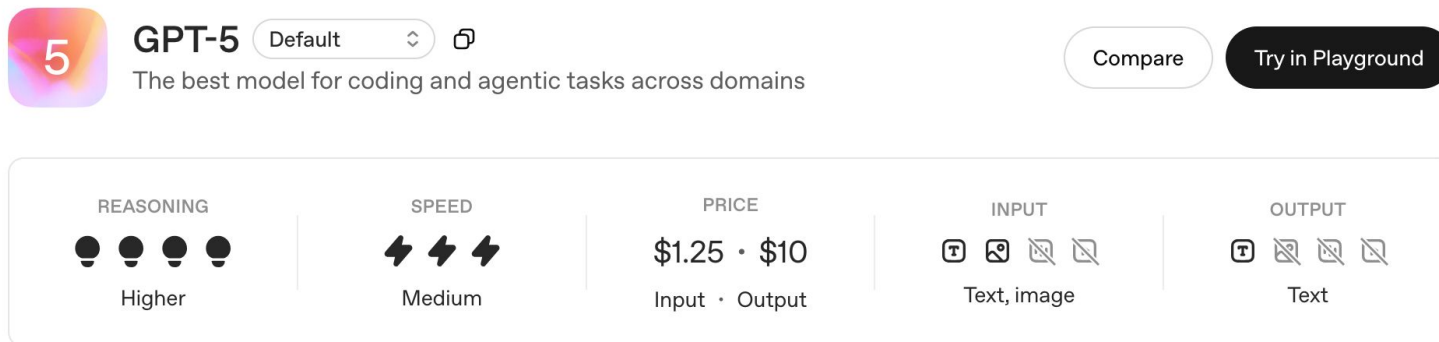# Transformers & Large Language Models

**RAG**

Tool calling

Agents

ICME

- **Knowledge** of LLM **constrained** to pretraining data

# Motivation

- **Knowledge** of LLM **constrained** to pretraining data



GPT-5 **Default**

The best model for coding and agentic tasks across domains

Compare          Try in Playground

| REASONING | SPEED | PRICE | INPUT | OUTPUT |
| --- | --- | --- | --- | --- |
| ●●●● | ⚡⚡⚡ | $1.25 · $10 | | |
| Higher | Medium | Input · Output | Text, image | Text |

GPT-5 is our flagship model for coding, reasoning, and agentic tasks across domains. Learn more in our GPT-5 usage guide.

✦ 400,000 context window

➡ 128,000 max output tokens

🗓 Sep 30, 2024 knowledge cutoff

💡 Reasoning token support

*Screenshot of OpenAI's GPT-5 model card, taken on November 9th, 2025.*

ICME

# Motivation

- **Knowledge** of LLM **constrained** to pretraining data



GPT-5 **Default**

The best model for coding and agentic tasks across domains

Compare     Try in Playground

| REASONING | SPEED | PRICE | INPUT | OUTPUT |
|---|---|---|---|---|
| Higher | Medium | $1.25 · $10<br>Input · Output | Text, image | Text |

GPT-5 is our flagship model for coding, reasoning, and agentic tasks across domains. Learn more in our GPT-5 usage guide.

- 400,000 context window
- 128,000 max output tokens
- Sep 30, 2024 knowledge cutoff
- Reasoning token support

*Screenshot of OpenAI's GPT-5 model card, taken on November 9th, 2025.*

ICME

# Motivation

- **Knowledge** of LLM **constrained** to pretraining data

- **Limited** context size

ICME

# Motivation

- Knowledge of LLM constrained to pretraining data

- **Limited** context size



REASONING · SPEED · PRICE · INPUT · OUTPUT

GPT-5 is our flagship model for coding, reasoning, and agentic tasks across domains. Learn more in our GPT-5 usage guide.

- ✦ 400,000 context window
- ➦ 128,000 max output tokens
- 🗓 Sep 30, 2024 knowledge cutoff
- 💡 Reasoning token support

Stanford University

ICME

# Motivation

- **Knowledge** of LLM **constrained** to pretraining data

- **Limited** context size

- LLM gets **distracted** by **useless** information

# Motivation

- **Knowledge** of LLM **constrained** to pretraining data

- **Limited** context size

- LLM gets **distracted** by **useless** information

*GitHub repo containing "Needle in a Haystack - Pressure Testing LLMs" tests, Kamradt, 2023.*

ICME

# Motivation

- **Knowledge** of LLM **constrained** to pretraining data

- **Limited** context size

- LLM gets **distracted** by **useless** information

- **Pricing** is **per** input/output **token**

# Motivation

- **Knowledge** of LLM **constrained** to pretraining data

- **Limited** context size

- LLM gets **distracted** by **useless** information

- **Pricing** is **per** input/output **token**

**RAG** = **R**etrieval-**A**ugmented **G**eneration

*"Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks", Lewis et al., 2020.*

**RAG** = **R**etrieval-**A**ugmented **G**eneration

**Idea**. Augment prompt with **relevant** pieces of information.

*"Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks", Lewis et al., 2020.*

ICME

**RAG** = **R**etrieval-**A**ugmented **G**eneration

**Idea**. Augment prompt with **relevant** pieces of information.

*Figure from "Super Study Guide: Transformers & Large Language Models", Amidi et al., 2024.*

**RAG** = **R**etrieval-**A**ugmented **G**eneration

**Idea**. Augment prompt with **relevant** pieces of information.

*Figure from "Super Study Guide: Transformers & Large Language Models", Amidi et al., 2024.*

**RAG** = **R**etrieval-**A**ugmented **G**eneration

**Idea**. Augment prompt with **relevant** pieces of information.

Figure from "*Super Study Guide: Transformers & Large Language Models*", Amidi et al., 2024.

1. **Retrieve** relevant document via similarity operation across the knowledge base

# Retrieve, Generate, Augment

1. **Retrieve** relevant document via similarity operation across the knowledge base

User prompt

2. **Augment** prompt with retrieved information

User prompt →

| Retrieved info |
| User prompt |

Stanford University

ICME

# Retrieve, Generate, Augment

1. **Retrieve** relevant document via similarity operation across the knowledge base

User prompt

2. **Augment** prompt with retrieved information

User prompt → Retrieved info / User prompt

3. **Generate** response

| Retrieved info |
| User prompt |

→ LLM → Response

ICME

# Retrieve, Generate, Augment

1. **R**etrieve relevant document via similarity operation across the knowledge base

2. **A**ugment prompt with retrieved information

3. **G**enerate response

# Focus on retrieval stage

**1. R**etrieve relevant document via similarity operation across the knowledge base

User prompt ◄─────────►

2. **A**ugment prompt with retrieved information

User prompt ──────► Retrieved info / User prompt

3. **G**enerate response

Retrieved info / User prompt ──────► LLM ──────► Response

Document 1

Document 2

...

**Collect**

Stanford University

ICME

Document 1    Document 2         ...

Document 1

Document 2

...

Collect                    **Divide**

Document 1  Document 2  ...

Document 1

Document 2

...

```
0.48

0.33

...

-0.51
```

**Collect**

**Divide**

**Embed**

ICME

Document 1   Document 2   ...

Document 1
Document 2
...

0.48

0.33

...

-0.51

**Collect**                    **Divide**                    **Embed**

**Hyperparameters**. Embedding size, chunk size, overlap between chunks

**Step 1 – Candidate retrieval**: Select potentially-relevant candidates

- Maximize recall
- Semantic embeddings and optionally keyword-based methods

ICME

# Retrieval overview

**Step 1 – Candidate retrieval**: Select potentially-relevant candidates

- Maximize recall
- Semantic embeddings and optionally keyword-based methods

**Step 2 – Ranking**: Give final relevance score

- Maximize precision
- Re-ranking on smaller set of candidates

ICME

**Objective**. Select potentially-relevant candidates via search across knowledge base

**Objective**. Select potentially-relevant candidates via search across knowledge base

**Method 1**. Semantic search using **embeddings-based similarity**

**Objective**. Select potentially-relevant candidates via search across knowledge base

**Method 1**. Semantic search using **embeddings-based similarity**

> Query

> Chunk

**Objective**. Select potentially-relevant candidates via search across knowledge base

**Method 1**. Semantic search using **embeddings-based similarity**

$$
\text{Query} = \begin{bmatrix} 0.48 \\ \dots \\ -0.51 \end{bmatrix}
$$

$$
\text{Chunk} = \begin{bmatrix} -0.82 \\ \dots \\ -0.43 \end{bmatrix}
$$

**Objective**. Select potentially-relevant candidates via search across knowledge base

**Method 1**. Semantic search using **embeddings-based similarity**



Stanford University

ICME

**Objective**. Select potentially-relevant candidates via search across knowledge base

**Method 1**. Semantic search using **embeddings-based similarity**

*"Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks", Reimers et al., 2019.*

ICME

**Objective**. Select potentially-relevant candidates via search across knowledge base

**Method 1**. Semantic search using **embeddings-based similarity**

"**bi-encoder**"



*"Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks", Reimers et al., 2019.*

**Objective**. Select potentially-relevant candidates via search across knowledge base

**Method 1**. Semantic search using **embeddings-based similarity**

| | [...] Huggy likes to work downstairs [...] |
|---|---|

Where is Cuddly? ⟷ [...] Where is Paris located? [...]

[...] he was in a cuddly mood [...]

[ ...] Cuddly spends most days surrounded by books [...]

**Objective**. Select potentially-relevant candidates via search across knowledge base

**Method 2**. Search based on keywords matching via **BM25**

# Step 1: candidate retrieval

**Objective**. Select potentially-relevant candidates via search across knowledge base

**Method 2**. Search based on keywords matching via **BM25**

[...] Huggy likes to work downstairs [...]

[...] Where is Paris located? [...]

Where is Cuddly?

[...] he was in a cuddly mood [...]

[ ...] Cuddly spends most days surrounded by books [...]

# Step 1: candidate retrieval

**Objective**. Select potentially-relevant candidates via search across knowledge base

**Method 3**. Search based on **hybrid combination** of semantics and BM25

ICME

**Objective**. Select potentially-relevant candidates via search across knowledge base

**Method 3**. Search based on **hybrid combination** of semantics and BM25

[...] Huggy likes to work downstairs [...]

[...] Where is Paris located? [...]

Where is Cuddly?  ⟷

[...] he was in a cuddly mood [...]

[ ...] Cuddly spends most days surrounded by books [...]

ICME

- Mitigate **discrepancy** in nature of embeddings

*"Precise Zero-Shot Dense Retrieval without Relevance Labels"*, Gao et al., 2022.

- Mitigate **discrepancy** in nature of embeddings

Prompt

> Where is Cuddly?

*"Precise Zero-Shot Dense Retrieval without Relevance Labels", Gao et al., 2022.*

ICME

- Mitigate **discrepancy** in nature of embeddings

Prompt

Where is Cuddly?

Fake
document

Cuddly is in...

*"Precise Zero-Shot Dense Retrieval without Relevance Labels", Gao et al., 2022.*

ICME

- Mitigate **discrepancy** in nature of embeddings

Prompt

Where is Cuddly?

Chunks

Fake document

Cuddly is in...

*"Precise Zero-Shot Dense Retrieval without Relevance Labels", Gao et al., 2022.*

ICME

- Mitigate **discrepancy** in nature of embeddings

- **Contextualize** document chunks

- Mitigate **discrepancy** in nature of embeddings

- **Contextualize** document chunks



| Chunk 1 |
| --- |

⋮

| Chunk n |
| --- |

→

| Context 1 |
| --- |
| Chunk 1 |

⋮

| Context n |
| --- |
| Chunk n |

*Introducing Contextual Retrieval*, *Anthropic team, 2024.*

ICME

- Mitigate **discrepancy** in nature of embeddings

- **Contextualize** document chunks

```
<document>
{WHOLE_DOCUMENT}
</document>
Here is the chunk we want to situate within the whole document:
{CHUNK_CONTENT}


Please give a short succinct context to situate this chunk within the overall
document for the purposes of improving search retrieval of the chunk. Answer
only with the succinct context and nothing else.
```

- Mitigate **discrepancy** in nature of embeddings

- **Contextualize** document chunks

**Prompt caching**

```
<document>
{WHOLE_DOCUMENT}
</document>
Here is the chunk we want to situate within the whole document:
{CHUNK_CONTENT}


Please give a short succinct context to situate this chunk within the overall
document for the purposes of improving search retrieval of the chunk. Answer
only with the succinct context and nothing else.
```

- Mitigate **discrepancy** in nature of embeddings

- **Contextualize** document chunks

## Pricing                                          ⧉ Copy page

**Text tokens**

Prices per 1M tokens.                    | Batch | Flex | **Standard** | Priority |

| MODEL | INPUT | CACHED INPUT | OUTPUT |
|---|---|---|---|
| gpt-5 | $1.25 | $0.125 | $10.00 |
| gpt-5-mini | $0.25 | $0.025 | $2.00 |
| gpt-5-nano | $0.05 | $0.005 | $0.40 |
| gpt-5-chat-latest | $1.25 | $0.125 | $10.00 |
| gpt-5-codex | $1.25 | $0.125 | $10.00 |

# Extensions that can help with initial retrieval

- Mitigate **discrepancy** in nature of embeddings
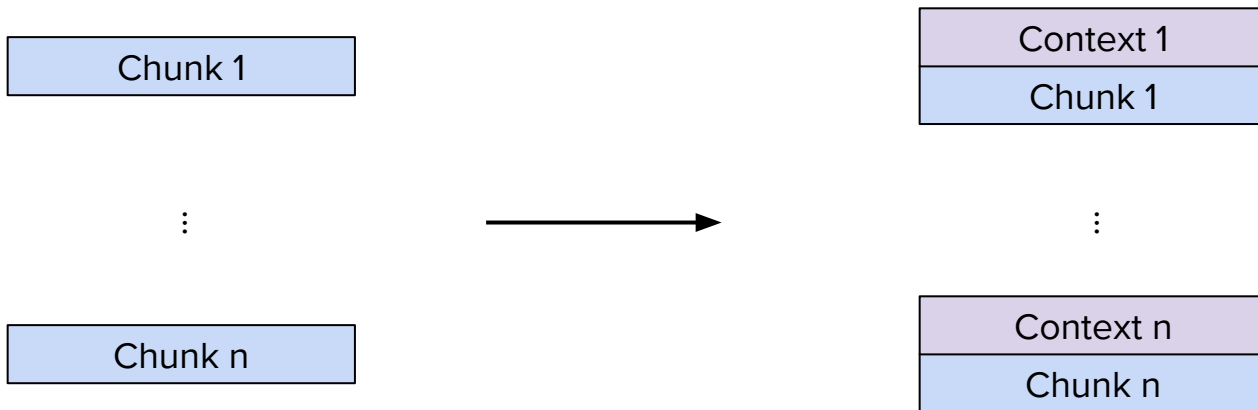
- **Contextualize** document chunks

## Pricing                                                    ⎙ Copy page

### Text tokens

Prices per 1M tokens.                          Batch    Flex   **Standard**   Priority

| MODEL | INPUT | CACHED INPUT | OUTPUT |
|---|---|---|---|
| gpt-5 | $1.25 | $0.125 | $10.00 |
| gpt-5-mini | $0.25 | $0.025 | $2.00 |
| gpt-5-nano | $0.05 | $0.005 | $0.40 |
| gpt-5-chat-latest | $1.25 | $0.125 | $10.00 |
| gpt-5-codex | $1.25 | $0.125 | $10.00 |

*OpenAI model pricing page, screenshot taken on November 10th, 2025.*                    ICME

**Objective**. Provide final relevance score using more sophisticated (re-)ranker

Query

Chunk

**Objective**. Provide final relevance score using more sophisticated (re-)ranker



*Suggested reading: "Cross-Encoders", SBERT.net*

**Objective**. Provide final relevance score using more sophisticated (re-)ranker

"**cross-encoder**"



*Suggested reading: "Cross-Encoders", SBERT.net*

**Objective**. Provide final relevance score using more sophisticated (re-)ranker

**Objective**. Provide final relevance score using more sophisticated (re-)ranker

**Setup**. Evaluate if **retrieved chunks** are **relevant**

**Setup**. Evaluate if **retrieved chunks** are **relevant**

Chunks

**Setup**. Evaluate if **retrieved chunks** are **relevant**



User prompt

Chunks

**Relevant**

# Quantify performance of retrieval

**Setup**. Evaluate if **retrieved chunks** are **relevant**



**Top k**

User prompt

**Relevant**

Chunks

# Quantify performance of retrieval

**Setup**. Evaluate if **retrieved chunks** are **relevant**

**Ranking**

**Setup**. Evaluate if **retrieved chunks** are **relevant**

**Ranking**

**Setup**. Evaluate if **retrieved chunks** are **relevant**

**Ranking**



- **N**ormalized **D**iscounted **C**umulative **G**ain at **k** (**NDCG@k**)

Stanford University

ICME

# Quantify performance of retrieval

**Setup**. Evaluate if **retrieved chunks** are **relevant**

**Ranking**



- **N**ormalized **D**iscounted **C**umulative **G**ain at **k** (**NDCG@k**)

$$\text{DCG@}k = \sum_{i=1}^{k} \frac{\text{rel}_i}{\log_2(i+1)}$$

with $\text{rel}_i \in \{0, 1\}$

ICME

**Setup**. Evaluate if **retrieved chunks** are **relevant**

**Ranking**

1

⋮

k

- - - - - - - - - - - - - -

k+1

⋮

n

- **N**ormalized **D**iscounted **C**umulative **G**ain at **k** (**NDCG@k**)

$$\text{DCG@}k = \sum_{i=1}^{k} \frac{\text{rel}_i}{\log_2(i+1)}$$

with $\text{rel}_i \in \{0, 1\}$

$$\text{NDCG@}k = \frac{\text{DCG@}k}{\text{IDCG@}k}$$

$\text{DCG@}k$  if ranking was perfect

ICME

**Setup**. Evaluate if **retrieved chunks** are **relevant**

**Ranking**

- Normalized Discounted Cumulative Gain at **k** (**NDCG@k**)

- **R**eciprocal **R**ank at **k** (**RR@k**)

$$RR = \frac{1}{rank}$$

rank of the first relevant chunk

# Quantify performance of retrieval

**Setup**. Evaluate if **retrieved chunks** are **relevant**

**Ranking**



- Normalized Discounted Cumulative Gain at **k** (**NDCG@k**)

- Reciprocal Rank at **k** (**RR@k**)

- **Recall at k**

Stanford University

ICME

**Setup**. Evaluate if **retrieved chunks** are **relevant**

**Ranking**



- Normalized Discounted Cumulative Gain at **k** (**NDCG@k**)

- Reciprocal Rank at **k** (**RR@k**)

- **Recall at k**

$$\text{Recall@}k = \frac{|\text{relevant in top } k|}{|\text{relevant}|}$$

ICME

# Quantify performance of retrieval

**Setup**. Evaluate if **retrieved chunks** are **relevant**

**Ranking**



- Normalized Discounted Cumulative Gain at k (NDCG@k)

- Reciprocal Rank at k (RR@k)

- Recall at k

- **Precision at k**

false

# Quantify performance of retrieval

**Setup**. Evaluate if **retrieved chunks** are **relevant**

**Ranking**



- Normalized Discounted Cumulative Gain at **k** (**NDCG@k**)

- Reciprocal Rank at **k** (**RR@k**)

- Recall at k

- **Precision at k**

$$\text{Precision@}k = \frac{|\text{relevant in top } k|}{k}$$

ICME

# Transformers & Large Language Models

RAG

**Tool calling**

Agents

**Unstructured**

RAG

## Motivation

| ID | Field | … |
|---|---|---|
| 123 | Obs | … |
| … | … | … |

**Unstructured**

RAG

**Structured**

Stanford University

ICME

# Idea

| ID | Field | … |
|----|-------|---|
| 123 | Obs | … |
| … | … | … |

**Function**

```python
def get_data(
  id, field, ...
):
    # Logic.
    ...
    return result
```

"**Tool calling** [...] allows autonomous systems to complete complex tasks by dynamically accessing and [may act] upon external resources."

"**Tool calling** […] allows autonomous systems to **complete** complex **tasks** by dynamically accessing and [may act] upon **external resources**."

# Real-life example

LLM as we've known it so far

```
Find a bear near me!
```

→ LLM →

```
Sorry, I don't know
which bears are near
you.
```

ICME

# Real-life example

**LLM as we've known it so far**

Find a bear near me! → LLM → Sorry, I don't know which bears are near you.

**LLM with tools**

```
<function API>
Find a bear near me!
```
→ LLM → ?

Stanford University

ICME

# Real-life example

LLM as we've known it so far

Find a bear near me! → LLM → Sorry, I don't know which bears are near you.

LLM with tools

<function API>
Find a bear near me! → LLM → ?

ICME

# Real-life example

## find_teddy_bear.py

```python
from dataclasses import dataclass
from geopy.distance import geodesic
import requests

@dataclass
class TeddyBearInfo:
    name: str
    distance_meters: float
    mood: str
    message: str

def find_teddy_bear(location: tuple[float, float]) -> TeddyBearInfo:
    """
    Finds the nearest teddy bear to the given GPS coordinates.

    Parameters:
        location: A (latitude, longitude) pair representing the user's current
            location.

    Returns:
        TeddyBearInfo: Information about the nearest teddy bear found.
    """
    # Call API to get the closest teddy bear
    user_lat, user_lon = location
    api_url = "https://api.to.teddy.bears.com/v1/closest"
    try:
        response = requests.get(
            api_url,
            params={"latitude": user_lat, "longitude": user_lon},
            timeout=5
        )
        response.raise_for_status()
        closest_teddy_bear = response.json()
    except requests.RequestException as e:
        raise RuntimeError(f"Failed to fetch teddy bear data from API: {e}")
```

```python
    # Extract coordinates from API response
    bear_lat, bear_lon = closest_teddy_bear["coords"]

    # Compute distance to the bear using geopy (returns distance in meters)
    distance = geodesic((user_lat, user_lon), (bear_lat, bear_lon)).meters

    return TeddyBearInfo(
        name=closest_teddy_bear["name"],
        distance_meters=round(distance, 2),
        mood=closest_teddy_bear["mood"],
        message=f"{closest_teddy_bear['name']} is {closest_teddy_bear['mood']} "
        f"and only {round(distance, 2)} meters away!"
    )
```

ICME

# Real-life example

**find_teddy_bear.py**

```python
from dataclasses import dataclass
from geopy.distance import geodesic
import requests

@dataclass
class TeddyBearInfo:
    name: str
    distance_meters: float
    mood: str
    message: str

def find_teddy_bear(location: tuple[float, float]) -> TeddyBearInfo:
    """
    Finds the nearest teddy bear to the given GPS coordinates.

    Parameters:
        location: A (latitude, longitude) pair representing the user's current
            location.

    Returns:
        TeddyBearInfo: Information about the nearest teddy bear found.
    """
    # Call API to get the closest teddy bear
    user_lat, user_lon = location
    api_url = "https://api.to.teddy.bears.com/v1/closest"
    try:
        response = requests.get(
            api_url,
            params={"latitude": user_lat, "longitude": user_lon},
            timeout=5
        )
        response.raise_for_status()
        closest_teddy_bear = response.json()
    except requests.RequestException as e:
        raise RuntimeError(f"Failed to fetch teddy bear data from API: {e}")
```

```python
    # Extract coordinates from API response
    bear_lat, bear_lon = closest_teddy_bear["coords"]

    # Compute distance to the bear using geopy (returns distance in meters)
    distance = geodesic((user_lat, user_lon), (bear_lat, bear_lon)).meters

    return TeddyBearInfo(
        name=closest_teddy_bear["name"],
        distance_meters=round(distance, 2),
        mood=closest_teddy_bear["mood"],
        message=f"{closest_teddy_bear['name']} is {closest_teddy_bear['mood']} "
        f"and only {round(distance, 2)} meters away!"
    )
```

has a descriptive,
well-documented API

ICME

# Real-life example

## find_teddy_bear.py

```python
from dataclasses import dataclass
from geopy.distance import geodesic
import requests

@dataclass
class TeddyBearInfo:
    name: str
    distance_meters: float
    mood: str
    message: str

def find_teddy_bear(location: tuple[float, float]) -> TeddyBearInfo:
    """
    Finds the nearest teddy bear to the given GPS coordinates.

    Parameters:
        location: A (latitude, longitude) pair representing the user's current
            location.

    Returns:
        TeddyBearInfo: Information about the nearest teddy bear found.
    """
    # Call API to get the closest teddy bear
    user_lat, user_lon = location
    api_url = "https://api.to.teddy.bears.com/v1/closest"
    try:
        response = requests.get(
            api_url,
            params={"latitude": user_lat, "longitude": user_lon},
            timeout=5
        )
        response.raise_for_status()
        closest_teddy_bear = response.json()
    except requests.RequestException as e:
        raise RuntimeError(f"Failed to fetch teddy bear data from API: {e}")
```

```python
    # Extract coordinates from API response
    bear_lat, bear_lon = closest_teddy_bear["coords"]

    # Compute distance to the bear using geopy (returns distance in meters)
    distance = geodesic((user_lat, user_lon), (bear_lat, bear_lon)).meters

    return TeddyBearInfo(
        name=closest_teddy_bear["name"],
        distance_meters=round(distance, 2),
        mood=closest_teddy_bear["mood"],
        message=f"{closest_teddy_bear['name']} is {closest_teddy_bear['mood']} "
        f"and only {round(distance, 2)} meters away!"
    )
```

(optional) has some backend call

Stanford University

ICME

# Real-life example

**find_teddy_bear.py**

```python
from dataclasses import dataclass
from geopy.distance import geodesic
import requests

@dataclass
class TeddyBearInfo:
    name: str
    distance_meters: float
    mood: str
    message: str

def find_teddy_bear(location: tuple[float, float]) -> TeddyBearInfo:
    """
    Finds the nearest teddy bear to the given GPS coordinates.

    Parameters:
        location: A (latitude, longitude) pair representing the user's current
            location.

    Returns:
        TeddyBearInfo: Information about the nearest teddy bear found.
    """
    # Call API to get the closest teddy bear
    user_lat, user_lon = location
    api_url = "https://api.to.teddy.bears.com/v1/closest"
    try:
        response = requests.get(
            api_url,
            params={"latitude": user_lat, "longitude": user_lon},
            timeout=5
        )
        response.raise_for_status()
        closest_teddy_bear = response.json()
    except requests.RequestException as e:
        raise RuntimeError(f"Failed to fetch teddy bear data from API: {e}")
```

```python
    # Extract coordinates from API response
    bear_lat, bear_lon = closest_teddy_bear["coords"]

    # Compute distance to the bear using geopy (returns distance in meters)
    distance = geodesic((user_lat, user_lon), (bear_lat, bear_lon)).meters

    return TeddyBearInfo(
        name=closest_teddy_bear["name"],
        distance_meters=round(distance, 2),
        mood=closest_teddy_bear["mood"],
        message=f"{closest_teddy_bear['name']} is {closest_teddy_bear['mood']} "
        f"and only {round(distance, 2)} meters away!"
    )
```

returns some info

1. Let **LLM** find **argument** for **relevant function** call

```
<function API>

Find a bear near me!
```

→ LLM →

```
location = (37.42,-122.17)
with find_teddy_bear()
```

1. Let **LLM** find **argument** for **relevant function** call

```
<function API>

Find a bear near me!
```

→ LLM →

```
location = (37.42,-122.17)
with find_teddy_bear()
```

2. Make **function call**

```
find_teddy_bear(location)
```

→ Backend →

```
{"name": "Teddy", …}
```

ICME

# How it works

## 1. Let **LLM** find **argument** for **relevant function** call

```
<function API>

Find a bear near me!
```
→ LLM →
```
location = (37.42,-122.17)
with find_teddy_bear()
```

## 2. Make **function call**

```
find_teddy_bear(location)
```
→ Backend →
```
{"name": "Teddy", ...}
```

## 3. Let **LLM deduce** conclusion **based on results**

```
{"name": "Teddy", ...}
```
→ LLM →
```
response
```

ICME

**Method 1**: via training

**Method 1**: via training



```
<function API>

Find a bear near me!
```

```
location = (37.42,-122.17)
with find_teddy_bear()
```

Tool prediction

# Teach a model to use a tool

**Method 1**: via training

```
<function API>

Find a bear near me!
```

↓

```
location = (37.42,-122.17)
with find_teddy_bear()
```

Tool prediction

```
<function API>

Find a bear near me!
```

LLM

```
location = (37.42,-122.17)
with find_teddy_bear()
```

Backend

```
{"name": "Teddy", …}
```

↓

```
response
```

Response generation

ICME

# Teach a model to use a tool

**Method 1**: via training

```
<function API>

Find a bear near me!
```

↓

```
location = (37.42,-122.17)
with find_teddy_bear()
```

Tool prediction

conversation history so far

```
<function API>

Find a bear near me!
```

LLM

```
location = (37.42,-122.17)
with find_teddy_bear()
```

Backend

```
{"name": "Teddy", …}
```

↓

```
response
```

Response generation

desired prediction

ICME

# Teach a model to use a tool

**Method 1**: via training

```
<function API>

Find a bear in Paris!
```

```
location = (48.86,2.35)
with find_teddy_bear()
```

Tool prediction

```
<function API>

Find a bear in Paris!
```

LLM

```
location = (48.86,2.35)
with find_teddy_bear()
```

Backend

```
{"name": "Ours", …}
```

response

Response generation

ICME

# Teach a model to use a tool

**Method 1**: via training

**Method 2**: via prompting

```
<function API> + <detailed explanation on how to use it>

Find a bear near me!
```

ICME

**Method 1**: via training

**Method 2**: via prompting

how to write such a description?

```
<function API> + <detailed explanation on how to use it>

Find a bear near me!
```

Stanford University

# Teach a model to use a tool

**Method 1**: via training

**Method 2**: via prompting

how to write such a description?

```
<function API> + <detailed explanation on how to use it>

Find a bear near me!
```

**One way:** Use SFT pairs as evaluation + use a powerful reasoning model to write it for you!

ICME

# Examples of common use cases

**Information**
- Web/database search
- Weather, stocks, and any other tracker
- Codebase

**Computation**
- Calculator

...and many more!

- Code execution (often in Python)

**Action**
- Send emails/messages and other in-computer action
- anything else within the domain of an assistant

ICME

# In practice: many tools

```python
def find_teddy_bear(
 location: tuple[float,
                 float]
) -> TeddyBearInfo:
    # ...
```

```python
def hug_teddy_bear(
 recipient: TeddyBear,
 intensity: str = "warm"
) -> HugResponse:
    # ...
```

```python
def check_teddy_mood(
 name: TeddyBear
) -> TeddyMood:
    # ...
```

```python
def send_teddy_gift(
 recipient: TeddyBear,
 gift: str = "poetry book"
) -> GiftResponse:
    # ...
```

```python
def schedule_playdate(
 host: TeddyBear,
 guest: TeddyBear,
 time: datetime
) -> Confirmation:
    # ...
```

```python
def send_message(
 recipient: TeddyBear,
 message: str
) -> MessageResponse:
    # ...
```

ICME

# In practice: many tools

```
<function API>

Find a bear near me!
```
→ LLM →
```
location = (37.42,-122.17)
with find_teddy_bear()
```

```
<list of function APIs>

Find a bear near me!
```
→ LLM →
```
location = (37.42,-122.17)
with find_teddy_bear()
```

**Benefits**.

- LLMs just became way more useful!

- They can also interact with the real world

- Overcomes "knowledge cutoff" limitation

**Benefits**.

- LLMs just became way more useful!

- They can also interact with the real world

- Overcomes "knowledge cutoff" limitation

**Challenges**.

- More tools = decrease performance

- Finite context length: not scalable

- Many tools to define. Lots of work.

# Tool selection

**Goal**. Both reduce latency and improve performance

*"Automatic Tool Selection to Reduce Large Language Model Latency", Robert et al., 2024.*

ICME

# Tool selection

**Goal**. Both reduce latency and improve performance

(1)    ┌─────────────────────────┐         ┌──────────┐
       │ `Find a bear near me!`   │  ───▶   │  Router  │  ───▶   List of **selected** tools
       └─────────────────────────┘         └──────────┘

*"Automatic Tool Selection to Reduce Large Language Model Latency", Robert et al., 2024.*

# Tool selection

**Goal**. Both reduce latency and improve performance



① Find a bear near me! → Router → List of **selected** tools

② <**selected** function APIs>
Find a bear near me! → LLM → ...

*"Automatic Tool Selection to Reduce Large Language Model Latency", Robert et al., 2024.*

ICME

**Goal**. Avoid duplication of tool implementations

# **MCP** = **M**odel **C**ontext **P**rotocol

**Idea**. Connect tools/data to LLMs in a standard way

*"Introducing the Model Context Protocol", Anthropic, 2024.*

**MCP** = **M**odel **C**ontext **P**rotocol

**Idea**. Connect tools/data to LLMs in a standard way



*"Introducing the Model Context Protocol", Anthropic, 2024.*

# Standardization: MCP



MCP host

MCP client

User

MCP server

Tools  Prompts  Resources

ICME

# Standardization: MCP



Claude Desktop

MCP client

"Recommend a new poetry book to my teddy bear"

Book provider
MCP server

Find

Recommend

find_title

recommend_taste

Personal collection

Top books

*"About MCP > Architecture overview", Anthropic, 2024.*

ICME

# Transformers & Large Language Models

RAG

Tool calling

**Agents**

ICME

"An **agent** is a system that autonomously pursues goals and completes tasks on a user's behalf."

ICME

# Definition

"An **agent** is a system that **autonomously pursues goals** and **completes tasks** on a user's behalf."

Stanford University

ICME

# High-level idea

**Traditional**

Question → LLM → Answer

**Reasoning**

Question → LLM ⇢ Reasoning ⇢ Answer

ICME

# High-level idea

**Traditional**

Question → LLM → Answer

**Reasoning**

Question → LLM ⤏ Reasoning ⤏ Answer

**Agent**

Question → LLM ⤏ Calls ⤏ LLM ⤏ ... ⤏ Answer

**ReAct** = **Re**ason + **Act**

*"ReAct: Synergizing Reasoning and Acting in Language Models", Yao et al., 2022.*

**Input**

`My teddy bear is cold.`
`Please do something.`

Examples:

- Manually entered (e.g. user question)

- External event (e.g. metric going beyond a threshold)

ICME

# ReAct in action

The user's teddy bear is cold, which may be due to the current temperature of the room, which is currently unknown.

Input ┈┈┈┈┈┈┈┈┈⟶ **Observe**

My teddy bear is cold.
Please do something.

- **Synthesize** previous actions + explicitly **state** what is **currently known** including own knowledge

- **Reasoning**-heavy step to figure out **what is needed**

ICME

# ReAct in action

The user's teddy bear is cold, which may be due to the current temperature of the room, which is currently unknown.

Input - - - - - - - - - - - - - -> Observe

My teddy bear is cold.
Please do something.

Plan

Determine the temperature of the room.

Detail **what tasks** need to be accomplished and **what tools** to call

# ReAct in action

The user's teddy bear is cold, which may be due to the current temperature of the room, which is currently unknown.

Input ---------→ Observe

My teddy bear is cold.
Please do something.

Act ←-------- Plan

Determine the temperature of the room.

- Perform an action via an **API**

- **Look** for info in a **database** of documents

ICME

The user's teddy bear is cold, which may be due to the current temperature of the room, which is currently unknown.

```
Input  ┄┄┄┄┄┄┄┄┄┄┄→  Observe
```

My teddy bear is cold.
Please do something.

```
Act  ←───────────  Plan
```

Determine the temperature of the room.

`get_current_room_temperature()`

- Perform an action via an **API**

- **Look** for info in a **database** of documents

ICME

The user's teddy bear is cold, which may be due to the current temperature of the room, which is currently unknown.

Input

**Observe**

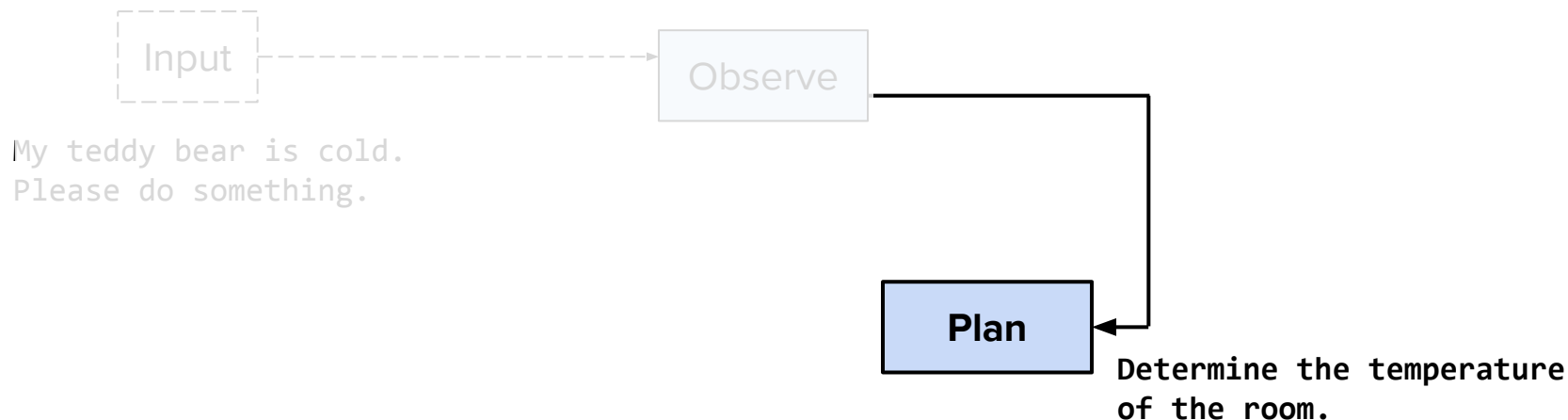My teddy bear is cold.
Please do something.

Act

Plan

Determine the temperature of the room.

get_current_room_temperature()

# ReAct in action

The temperature in the room is currently 65F. This is about 5F less than an average temperature. We need to increase the room temperature.

Input ⇢ **Observe**

My teddy bear is cold.
Please do something.

Act        Plan

Determine the temperature of the room.

get_current_room_temperature()

Stanford University

# ReAct in action

The temperature in the room is currently 65F. This is about 5F less than an average temperature. We need to increase the room temperature.

Input

Observe

My teddy bear is cold. Please do something.

Act

Plan

Determine the temperature of the room.

get_current_room_temperature()

Stanford University

ICME

# ReAct in action

The temperature in the room is currently 65F. This is about 5F less than an average temperature. We need to increase the room temperature.

Input ----> Observe

My teddy bear is cold.
Please do something.

Act <---- **Plan**

get_current_room_temperature()

Increase the temperature by 5F.

Stanford University

ICME

The temperature in the room is
currently 65F. This is about 5F less
than an average temperature. We need
to increase the room temperature.

```
Input
```

```
Observe
```

My teddy bear is cold.
Please do something.

```
Act
```

```
Plan
```

Increase the temperature
by 5F.

get_current_room_temperature()

The temperature in the room is
currently 65F. This is about 5F less
than an average temperature. We need
to increase the room temperature.

Input

Observe

My teddy bear is cold.
Please do something.

Act

Plan

Increase the temperature
by 5F.

increase_temperature(value=5)

ICME
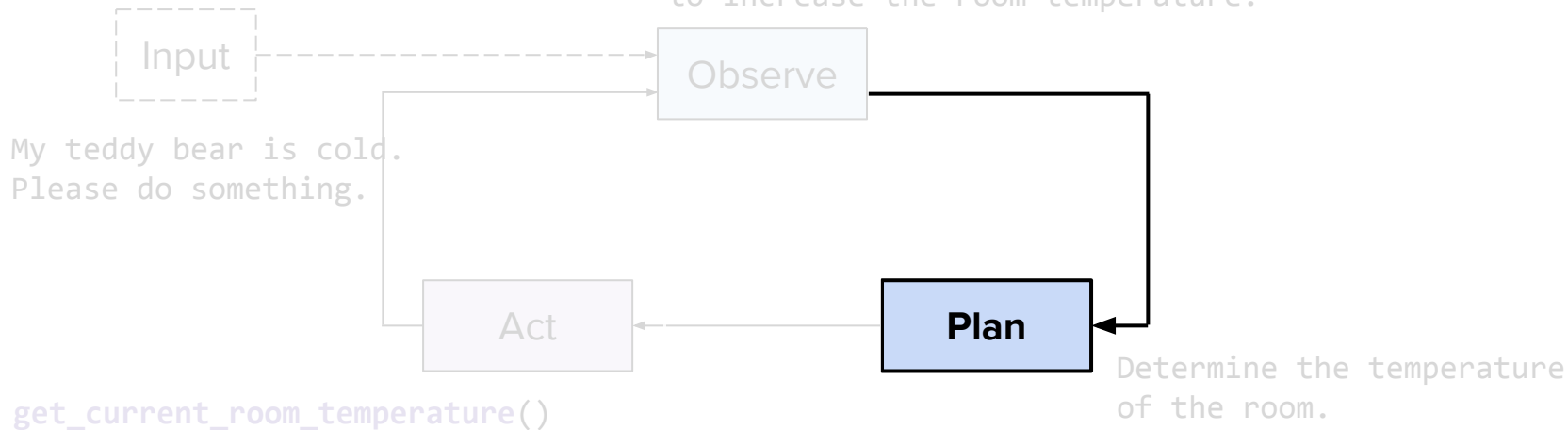
The temperature in the room is
currently 65F. This is about 5F less
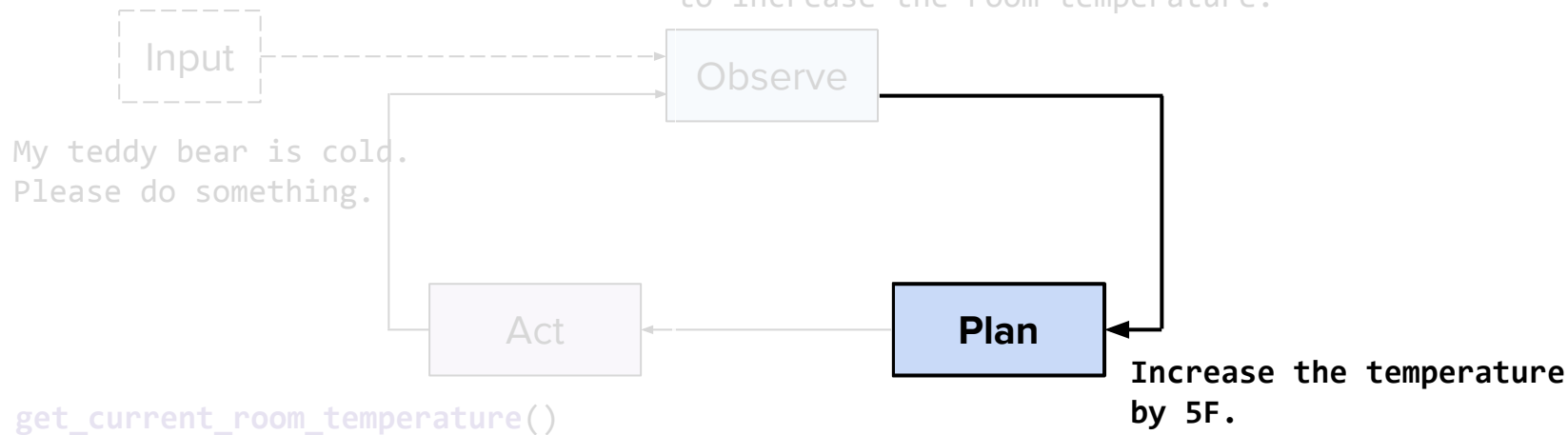than an average temperature. We need
to increase the room temperature.

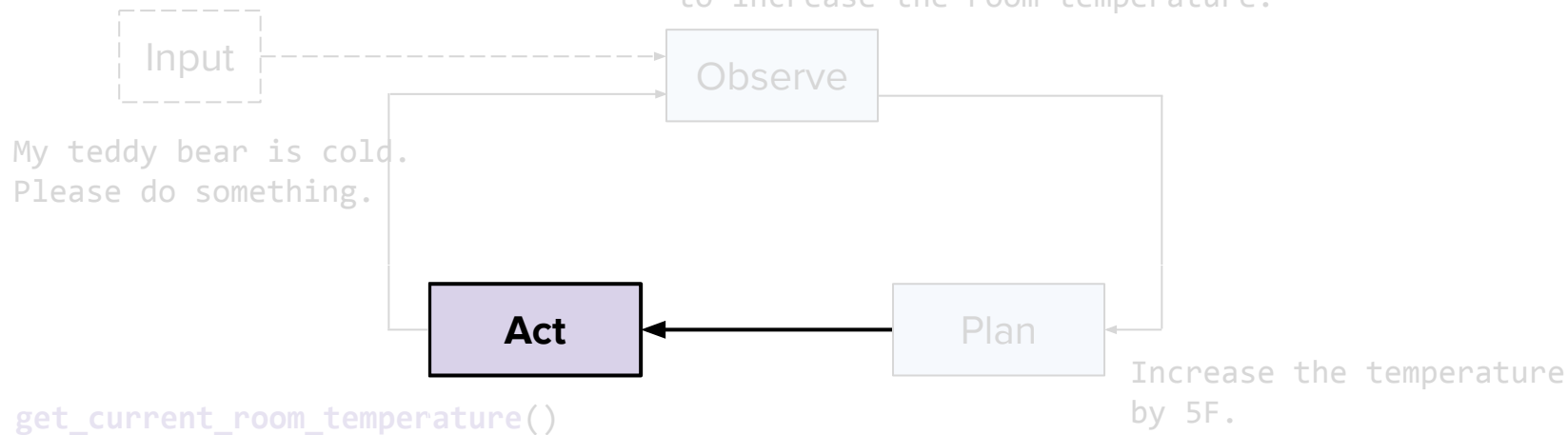Input - - - - - - - - - - - → **Observe**

My teddy bear is cold.
Please do something.

Act ← Plan

increase_temperature(value=5)

Increase the temperature
by 5F.

Stanford University

ICME

# ReAct in action

The thermostat is now set to 70F. This
should be warm enough.

```
Input
```

My teddy bear is cold.
Please do something.

**Observe**

```
Act
```

```
Plan
```

Increase the temperature
by 5F.

increase_temperature(value=5)

Stanford University

ICME

# ReAct in action

The thermostat is now set to 70F. This
should be warm enough.

Input - - - → Observe ═ ═ ═ ═ ═ → **Output**

My teddy bear is cold.
Please do something.

Act ← Plan

increase_temperature(value=5)

Increase the temperature
by 5F.

Stanford University

# ReAct in action

The thermostat is now set to 70F. This
should be warm enough.

Input - - - - - - - - - - - → Observe - - - - - - - - → **Output**

My teddy bear is cold.
Please do something.

**The thermostat is now
set to 70F. Your
teddy bear soon feel
warmer.**

Act ← Plan

Increase the temperature
by 5F.

increase_temperature(value=5)

ICME

Input

My teddy bear is cold.
Please do something.

Thermostat agent

Output

The thermostat is now
set to 70F. Your
teddy bear soon feel
warmer.

ICME

# There can be agents for many things

Occupancy agent

Energy management agent

Thermostat agent

Air quality agent

Stanford University

ICME

# How can agents communicate with each other?



Occupancy agent

?

Energy management agent

?

Thermostat agent

?

Air quality agent

ICME

**A2A** = **A**gent**2A**gent

*"Announcing the Agent2Agent Protocol (A2A)"*, *Google, 2025.*

ICME

# Standardization: A2A

Thermostat agent    =

### AgentSkill

```
id='optimize_
energy'
description=
'Reduces energy
consumption'
examples=['Save
energy
overnight']
```

### AgentCard

```
name='thermostat
agent'
url='http://path.
to.agent'
version='1.0.0'
skills=[
 optimize_energy,
 prepare_home,
]
```

### AgentExecutor

```
async def
execute(self, ...):
  # Logic.

async def
cancel(self, ...):
  # Logic.
```

*"Agent2Agent (A2A) Protocol Official Specification", Google, 2025.*

ICME

**Risks**.

- Potential for harm in the real world

- Example: data exfiltration

ICME

# Safety

**Risks**.

- Potential for harm in the real world

- Example: data exfiltration

**Remediations**.

- Training steps

- Inference safeguards

- Benchmarks, e.g. Agent-SafetyBench

*"Towards Tool Use Alignment of Large Language Models", Chen et al., 2024.*

ICME

**Risks**.

- Potential for harm in the real world

- Example: data exfiltration

**Remediations**.

- Training steps

- Inference safeguards

- Benchmarks, e.g. Agent-SafetyBench

**...very important topic!**

ICME

# Just yesterday in the news

*"Disrupting the first reported AI-orchestrated cyber espionage campaign", Anthropic, 2025.*

ICME

- **Hallucination** is a (big) problem

# Closing thoughts

- **Hallucination** is a (big) problem

- **Reasoning** abilities are a bottleneck
    - Finetuning helps, but hard
    - New capabilities are very welcome

ICME

Closing thoughts

- **Hallucination** is a (big) problem

- **Reasoning** abilities are a bottleneck

  - Finetuning helps, but hard

  - New capabilities are very welcome

- **Evaluation** is challenging

Stanford University
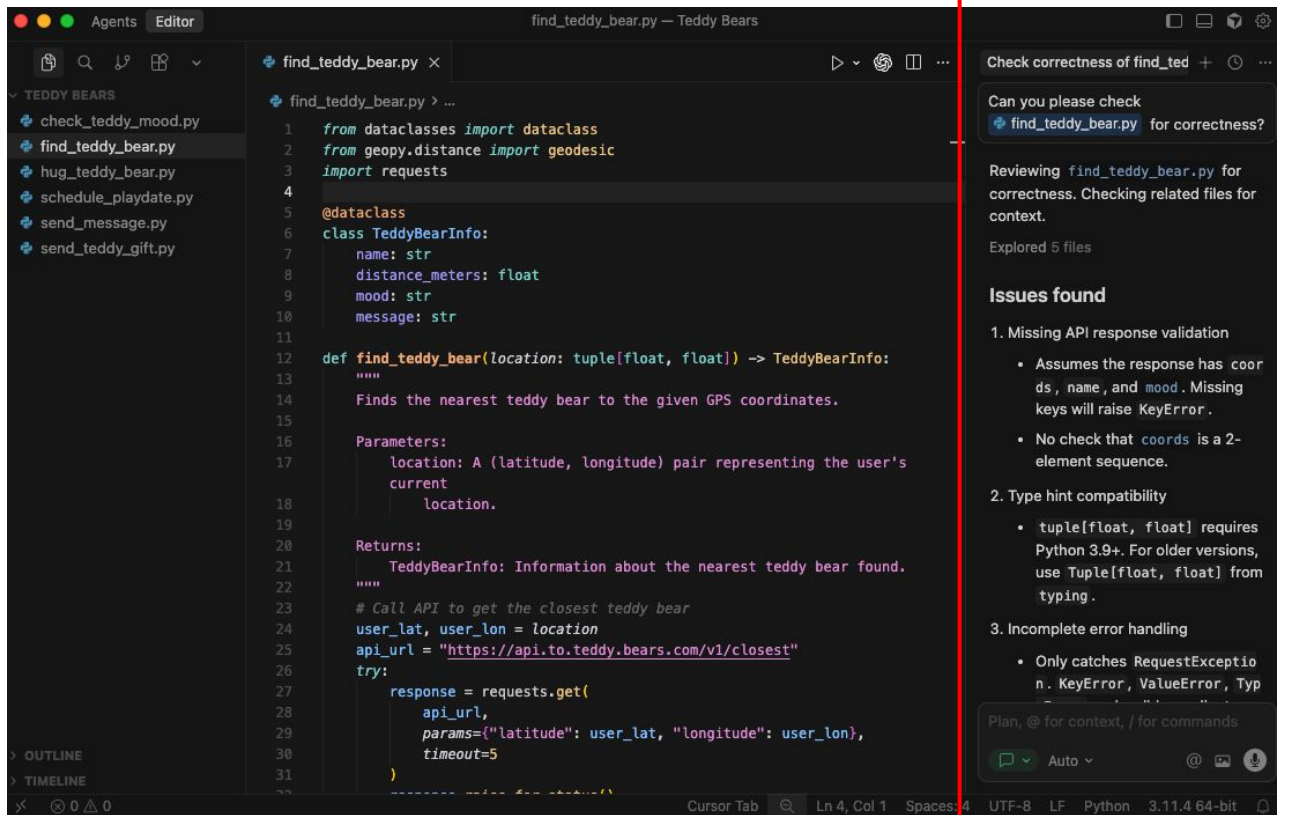
ICME

# Closing thoughts

- **Hallucination** is a (big) problem

- **Reasoning** abilities are a bottleneck

  - Finetuning helps, but hard

  - New capabilities are very welcome

- **Evaluation** is challenging

- Good to **start simple**, then **iterate** and **progressively scale up**

- Good to **start** with **capable models**, **optimize** on **size later**

ICME

- **Hallucination** is a (big) problem

- **Reasoning** abilities are a bottleneck

    - Finetuning helps, but hard

    - New capabilities are very welcome

- **Evaluation** is challenging

- Good to **start simple**, then **iterate** and **progressively scale up**

- Good to **start** with **capable models**, **optimize** on **size later**

- **Transparency** / **observability** helps with user trust and debuggability

Stanford University

ICME

# Bonus: AI agents in your daily life

Personal favorite use case: coding!

agent

Thank you for your attention!