



# AI Agents: Mastering Complex Task Solving

## A Survey

August 18, 2023 - Machine Learning Reading Group

INTERNAL – SAP Only

# **What are AI Agents?**

“Artificial Intelligence Agent (AI Agent) is defined as a program that employs artificial intelligence techniques to perform tasks that typically require human-like intelligence”

[\[IPTU: Task Planning and Tool Usage of Large Language Model-based AI Agents\]](#)

# What are AI Agents?

“Artificial Intelligence Agent (AI Agent) is defined as a program that employs artificial intelligence techniques to perform tasks that typically require human-like intelligence”

[\[IPTU: Task Planning and Tool Usage of Large Language Model-based AI Agents\]](#)

## Are AI agents necessarily LLM-based agents?

No, but LLMs **are a paradigm shift for AI agents** because they are the **best and most generic world models** we have at our disposal.

I guess they will be replaced by some form of **multimodal models at some point**.

What

“Artificial  
intelligen  
[\[PTU: Task P](#)

Are A

No, but I  
most ge  
I guess t

rtificial  
ligence”

and  
point.

BLOG ›

## Google Research, 2022 & beyond: Robotics

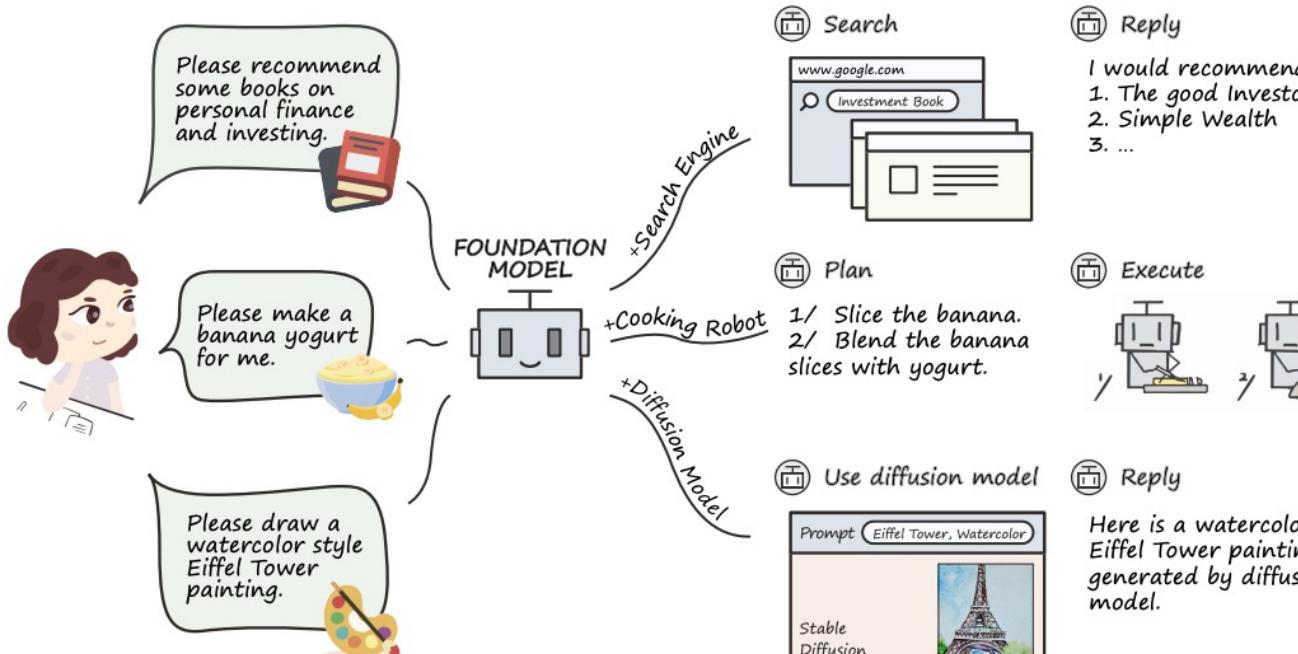
TUESDAY, FEBRUARY 14, 2023

Posted by Kendra Byrne, Senior Product Manager, and Jie Tan, Staff Research Scientist, Robotics at Google

“Advances in large models across the field of AI have spurred a leap in capabilities for robot learning. This past year, we’ve seen the sense of context and sequencing of events captured in LLMs help solve long-horizon planning for robotics and make robots easier for people to interact with and task. ...”

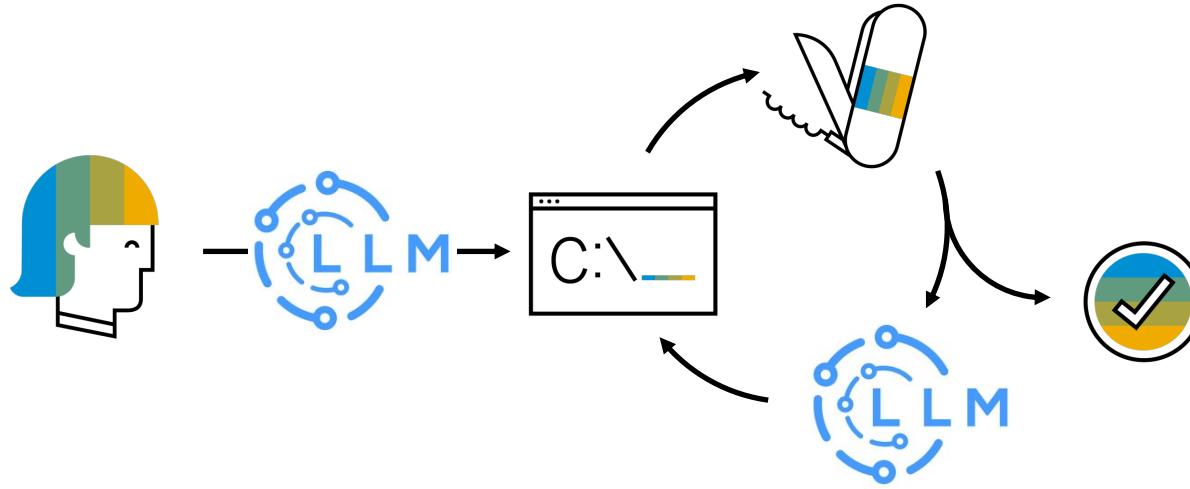
[<https://ai.googleblog.com/2023/02/google-research-2022-beyond-robotics.html>]

# “Tool Learning with Foundational Models”

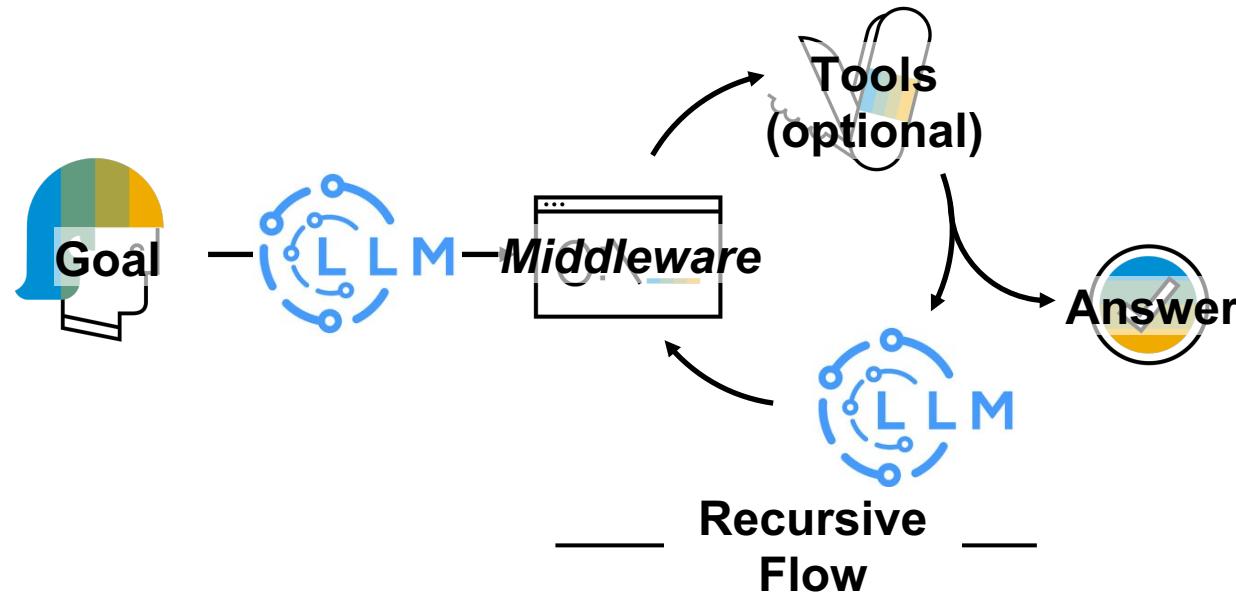


**Figure 1:** Tool learning paradigm aims to combine the strengths of specialized tools and foundation models.

# What are AI Agents?

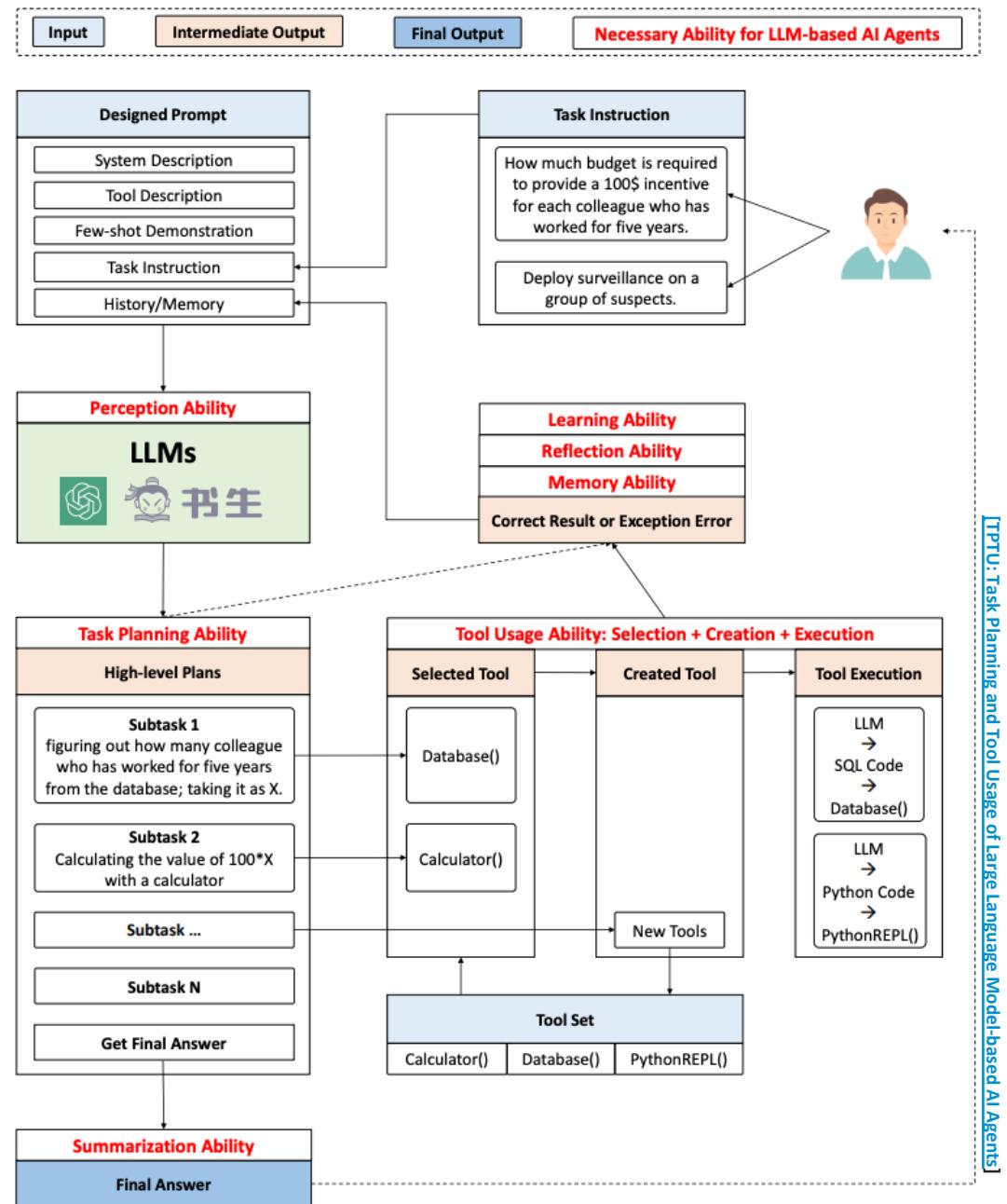


# What are AI Agents?



# What are AI Agents?

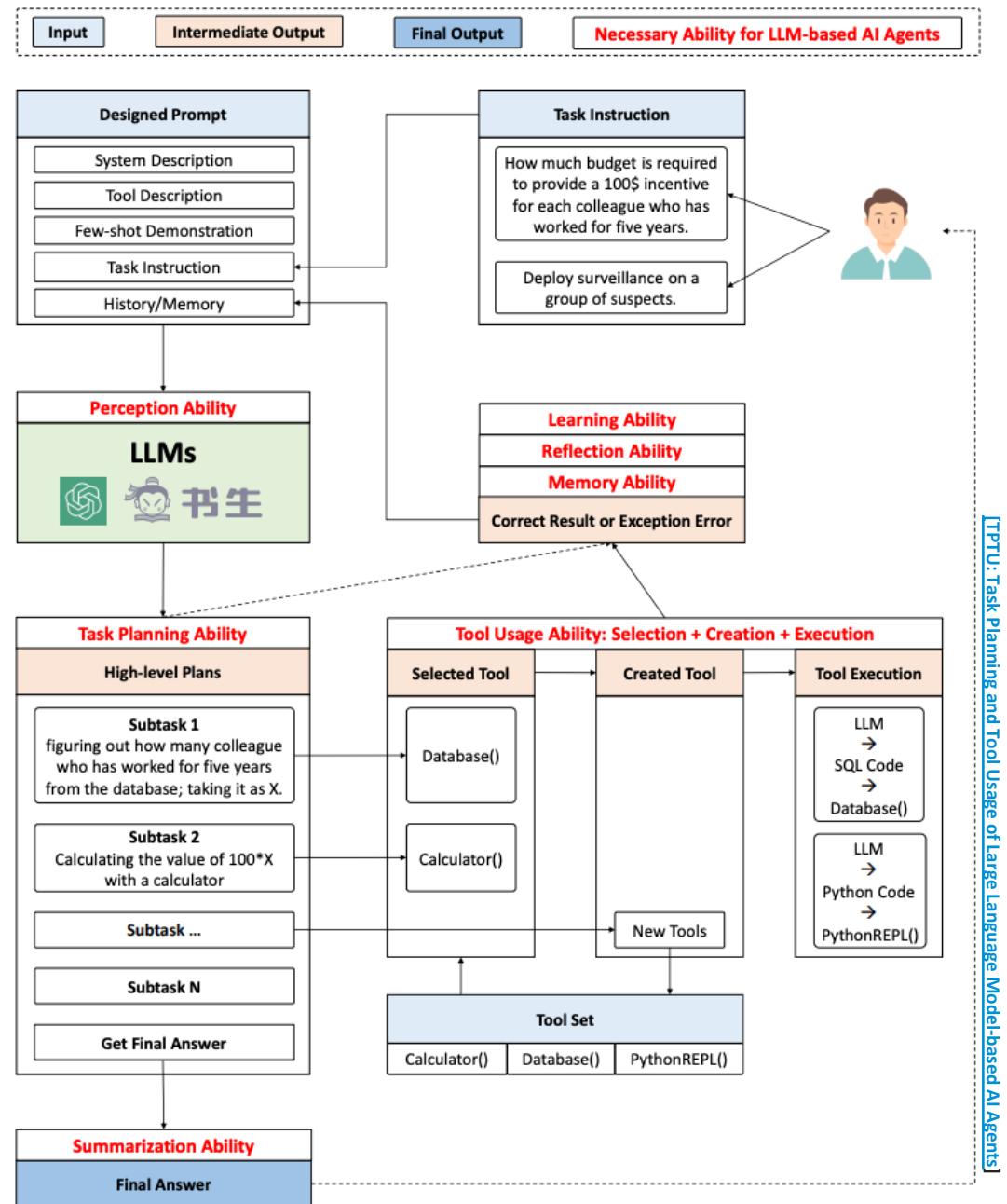
## AI Agent Framework:



# What are AI Agents?

## AI Agent Framework:

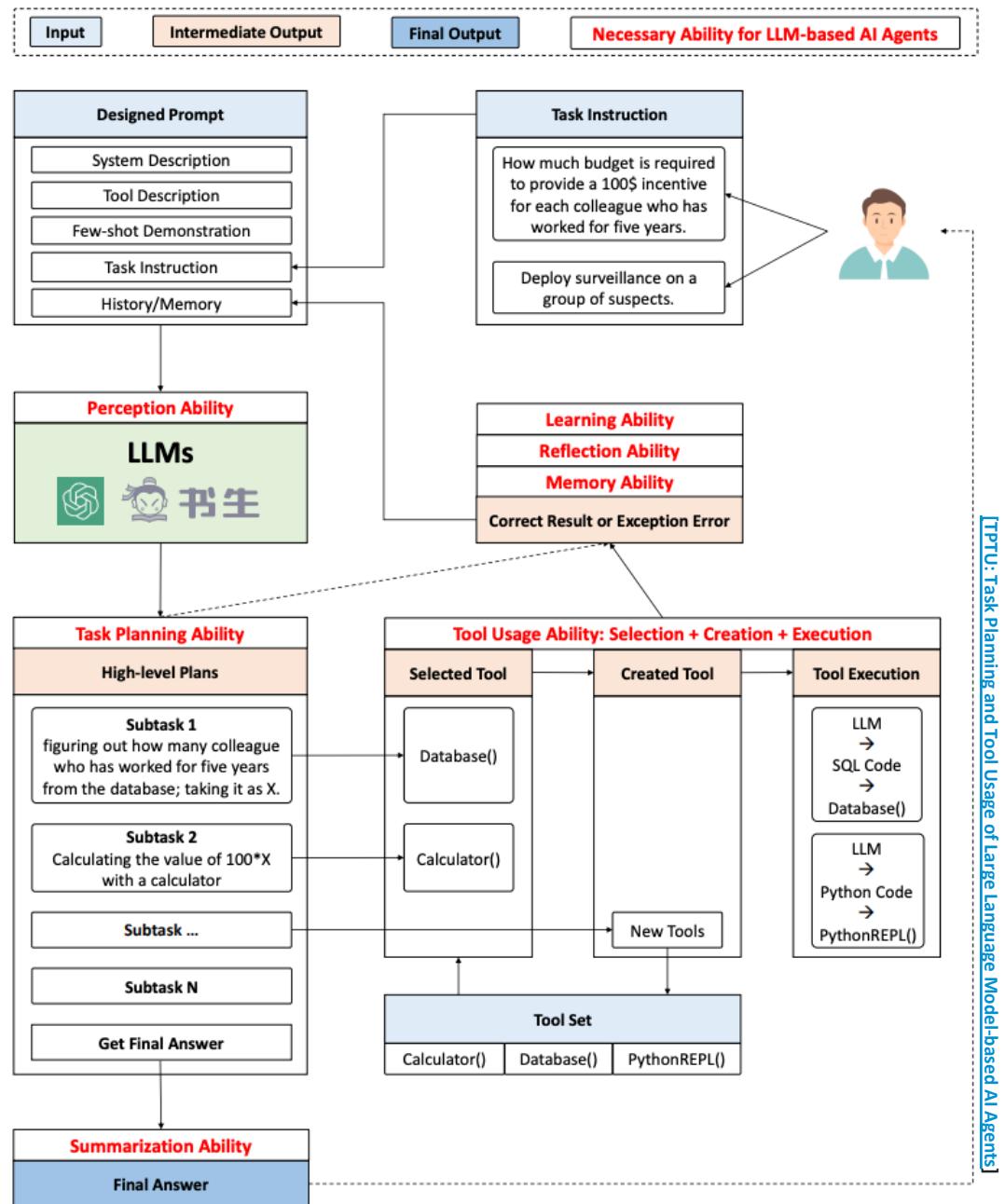
- Components:
  - Task Instruction
  - Designed Prompt
  - Tool Set
  - LLM
  - Intermediate Output/State
  - Final Answer



# What are AI Agents?

## AI Agent Framework:

- Components:
  - Task Instruction
  - Designed Prompt
  - Tool Set
  - LLM
  - Intermediate Output/State
  - Final Answer
- Abilities
  - Perception
  - Task Planning
  - Tool Usage
  - Learning/Reflection/Memory
  - Summarization



# What are AI Agents?

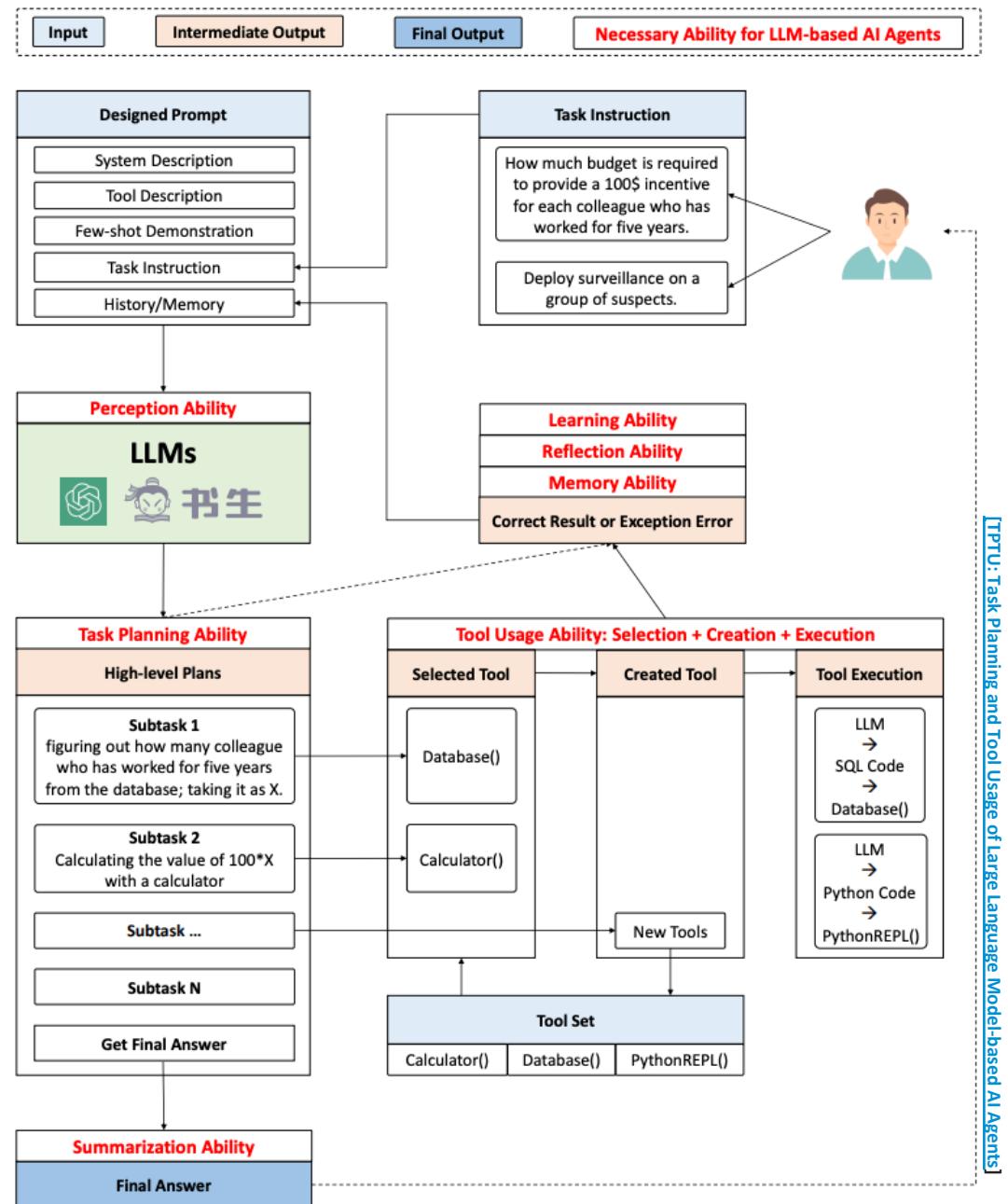
## AI Agent Framework:

- Components:

- |                               |   |
|-------------------------------|---|
| Goal                          | <ul style="list-style-type: none"> <li>Task Instruction</li> <li>Designed Prompt</li> </ul>       |
| Recursive Flow/<br>Middleware | <ul style="list-style-type: none"> <li>Tool Set</li> <li>LLM</li> </ul>                           |
| Answer                        | <ul style="list-style-type: none"> <li>Intermediate Output/State</li> <li>Final Answer</li> </ul> |

- Abilities

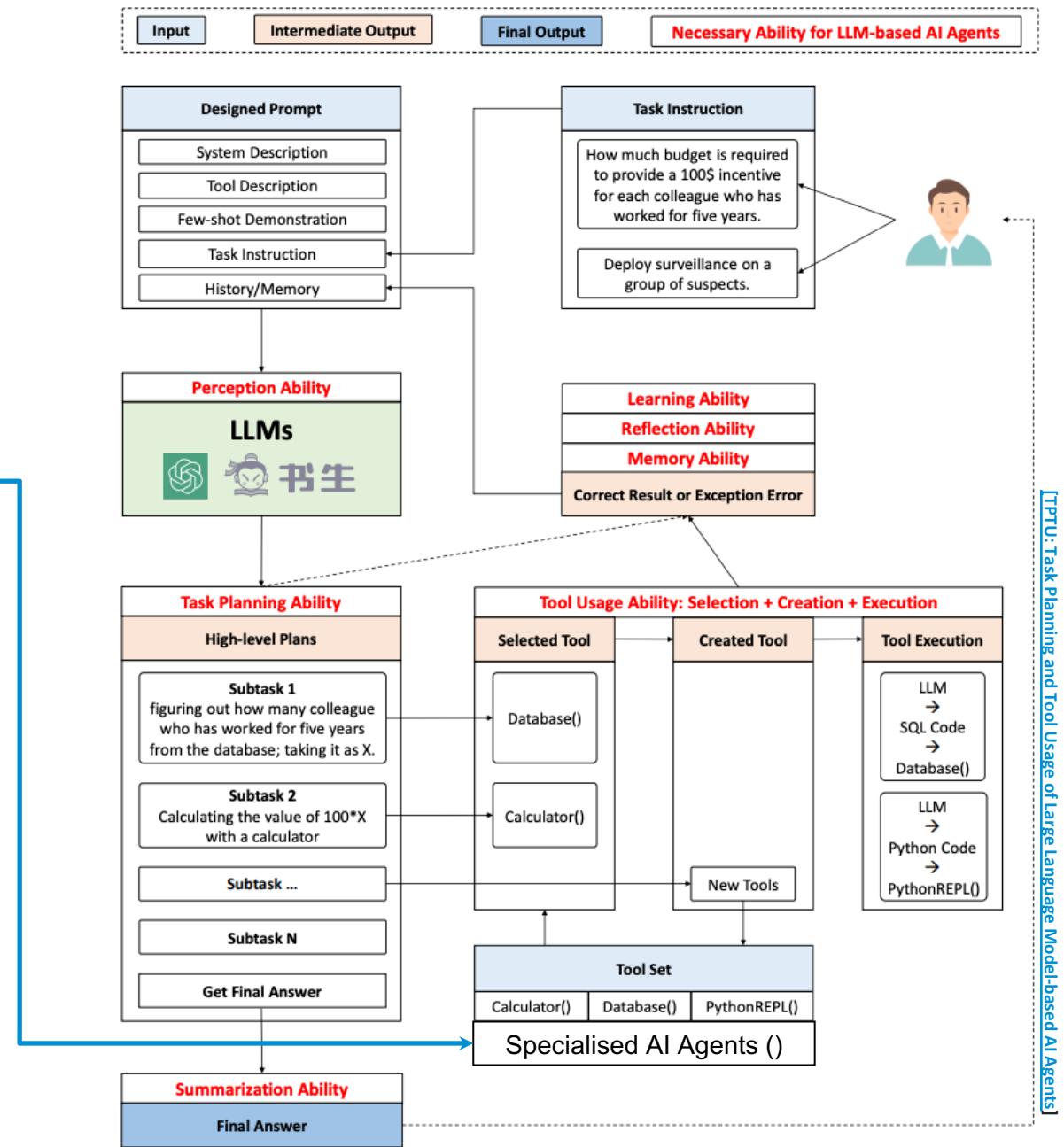
- |                               |   |
|-------------------------------|---|
| Goal                          | <ul style="list-style-type: none"> <li>Perception</li> </ul>  |
| Recursive Flow/<br>Middleware | <ul style="list-style-type: none"> <li>Task Planning</li> <li>Tool Usage</li> </ul>                 |
| Answer                        | <ul style="list-style-type: none"> <li>Learning/Reflection/Memory</li> <li>Summarization</li> </ul> |

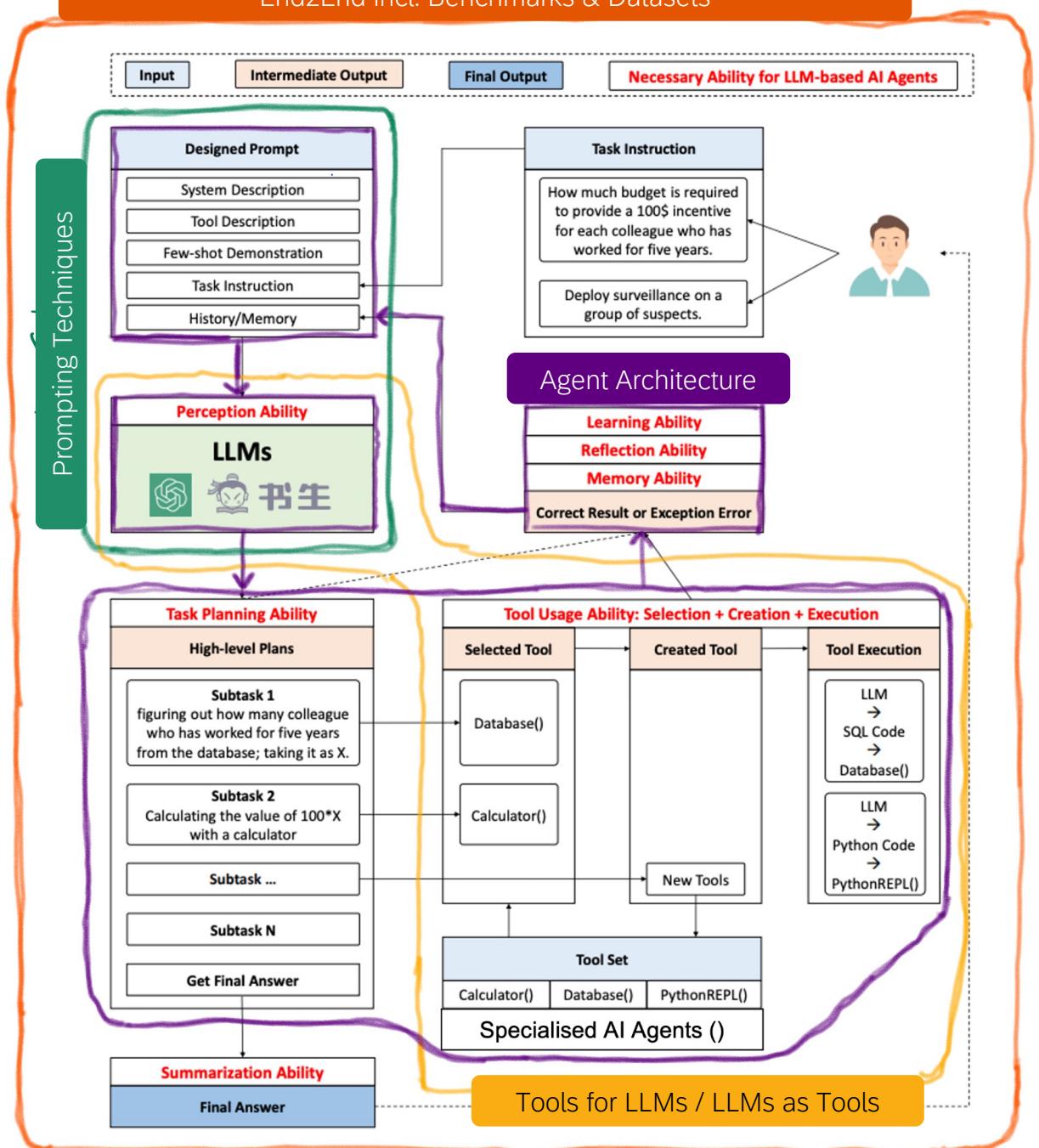


# What are AI Agents?

## AI Agent Framework:

- Components:
  - Task Instruction
  - Designed Prompt
  - Tool Set **including specialized Agents**
  - LLM
  - Intermediate Output/State
  - Final Answer
- Abilities
  - Perception
  - Task Planning
  - Tool Usage
  - Learning/Reflection/Memory
  - Summarization





# AI Agent Landscape

## Paper Categories:

End2End incl. Benchmarks & Datasets

Prompting Techniques

Agent Architecture

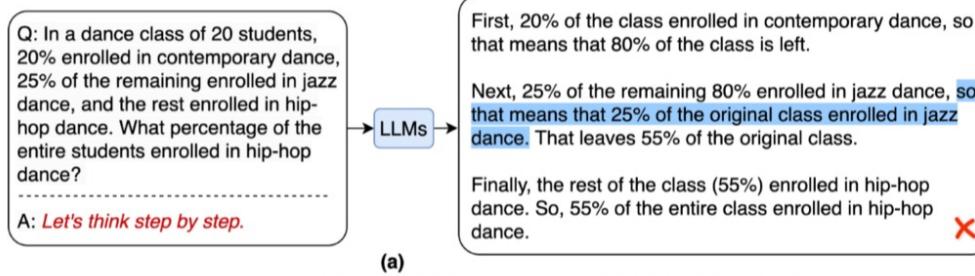
Tools for LLMs / LLMs as Tools

# „Plan-and-Solve Prompting: Improving Zero-Shot Chain-of-Thought Reasoning by Large Language Models“ [\[Link\]](#)

- 1.Singapore Management University
- 2.Southwest Jiaotong University
- 3.Singapore University of Technology and Design
- 4.East China Normal University

## Prompting Techniques

CoT:



PS:

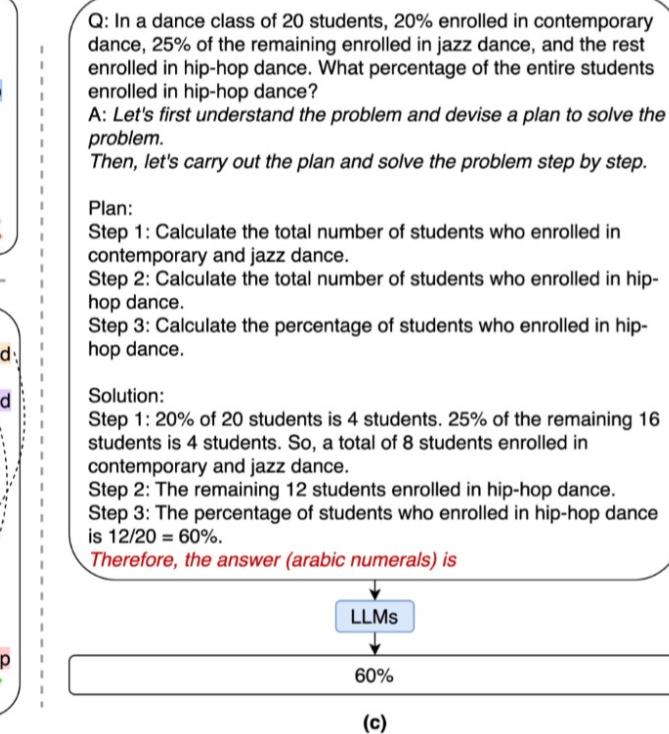
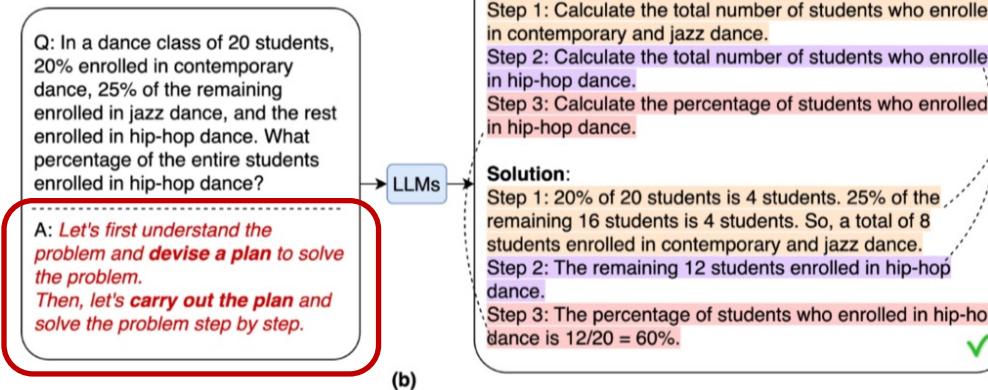


Figure 2: Example inputs and outputs of GPT-3 with (a) Zero-shot-CoT prompting, (b) Plan-and-Solve (PS) prompting, and (c) answer extraction prompting. While Zero-shot-CoT encourages LLMs to generate multi-step reasoning with “*Let's think step by step*”, it may still generate wrong reasoning steps when the problem is complex. Unlike Zero-shot-CoT, PS prompting first asks LLMs to devise a plan to solve the problem by generating a step-by-step plan and carrying out the plan to find the answer.

# „Plan-and-Solve Prompting: Improving Zero-Shot Chain-of-Thought Reasoning by Large Language Models“ [\[Link\]](#)

- 1.Singapore Management University
- 2.Southwest Jiaotong University
- 3.Singapore University of Technology and Design
- 4.East China Normal University

## Prompting Techniques

Table 2: Accuracy comparison on six math reasoning datasets. The best and second best results are boldfaced and underlined respectively.

Setting	Method (text-davinci-003)	MultiArith	GSM8K	AddSub	AQuA	SingleEq	SVAMP	Average
Zero-Shot	CoT	83.8	56.4	85.3	38.9	88.1	69.9	70.4
	PoT	<b>92.2</b>	57.0	85.1	<u>43.9</u>	<u>91.7</u>	70.8	<u>73.5</u>
	PS (ours)	87.2	<u>58.2</u>	<u>88.1</u>	42.5	89.2	72.0	72.9
	PS+ (ours)	<u>91.8</u>	<b>59.3</b>	<b>92.2</b>	<b>46.0</b>	<b>94.7</b>	<b>75.7</b>	<b>76.7</b>
Few-Shot	Manual-CoT	93.6	58.4	91.6	48.4	93.5	80.3	77.6
	Auto-CoT	95.5	57.1	90.8	41.7	92.1	78.1	75.9

Table 3: Accuracy on commonsense reasoning datasets.

Method	CSQA	StrategyQA
Few-Shot-CoT (Manual)	78.3	71.2
Zero-shot-CoT	65.2	63.8
Zero-shot-PS+ (ours)	<b>71.9</b>	<b>65.4</b>

Table 4: Accuracy on symbolic reasoning datasets.

Method	Last Letter	Coin Flip
Few-Shot-CoT (Manual)	70.6	100.0
Zero-shot-CoT	64.8	96.8
Zero-shot-PS+ (ours)	<b>75.2</b>	<b>99.6</b>

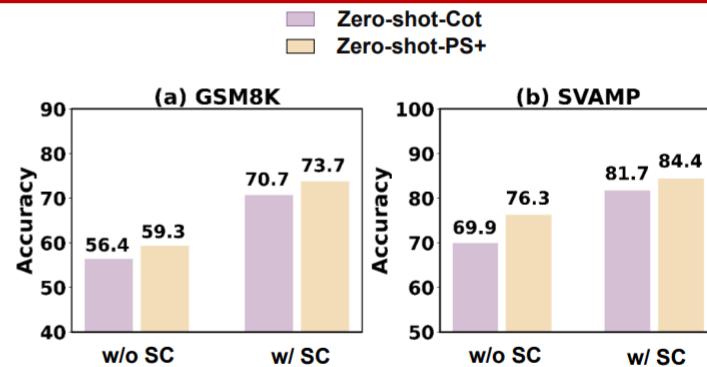


Figure 4: Results of methods with and without self-consistency (SC) on GSM8K and SVAMP.

# „Plan-and-Solve Prompting: Improving Zero-Shot Chain-of-Thought Reasoning by Large Language Models“ [\[Link\]](#)

## Prompting Techniques

Table 2: Accuracy comparison on six math reasoning datasets. The best and second best results are boldfaced and underlined.

Setting
Zero-Shot
Few-Shot

**Formulating a complete plan before solving** and giving the model a **chance to reflect** on it seems to **help the model reason significantly**.

Similar findings in:

“TPTU: Task Planning and Tool Usage of Large Language Model-based AI Agents” [\[Link\]](#)

“Reflexion: Language Agents with Verbal Reinforcement Learning” [\[Link\]](#)

“ReAct: Synergizing Reasoning and Acting in Language Models” [\[Link\]](#)

Average
70.4
<u>73.5</u>
72.9
<b>76.7</b>
77.6
75.9

Table 3: Accuracies

Method	Last Letter	Coin Flip
Few-Shot-CoT (Manual)	78.3	71.2
Zero-shot-CoT	65.2	63.8
Zero-shot-PS+ (ours)	<b>71.9</b>	<b>65.4</b>

Table 4: Accuracy on symbolic reasoning datasets.

Method	Last Letter	Coin Flip
Few-Shot-CoT (Manual)	70.6	100.0
Zero-shot-CoT	64.8	96.8
Zero-shot-PS+ (ours)	<b>75.2</b>	<b>99.6</b>

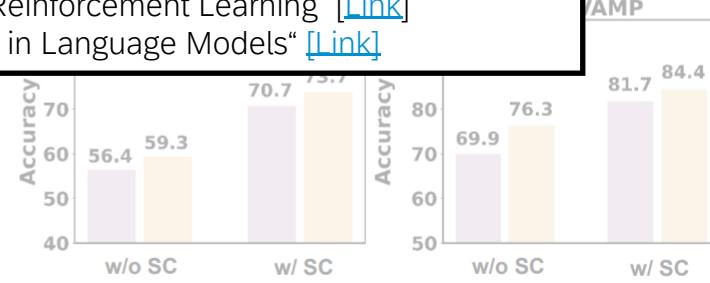


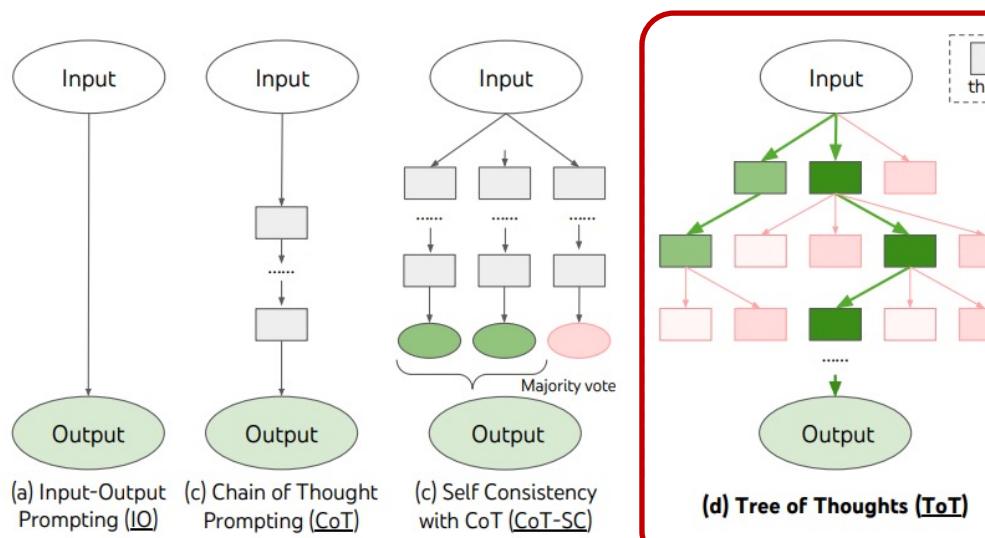
Figure 4: Results of methods with and without self-consistency (SC) on GSM8K and SVAMP.

# „Tree of Thoughts: Deliberate Problem Solving with Large Language Models“

[\[Link\]](#)

1.Princeton  
2.Google DeepMind

## Prompting Techniques



## Approach:

1. **Thought decomposition + generator**
2. **Scoring/Ordering Nodes**
3. **Search algorithm (BFS vs. DFS)**

# „Tree of Thoughts: Deliberate Problem Solving with Large Language Models“

[\[Link\]](#)

1. Princeton  
2. Google DeepMind

## Prompting Techniques

Method	Success
IO prompt	7.3%
CoT prompt	4.0%
CoT-SC ( $k=100$ )	9.0%
ToT (ours) ( $b=1$ )	45%
ToT (ours) ( $b=5$ )	<b>74%</b>
IO + Refine ( $k=10$ )	27%
IO (best of 100)	33%
CoT (best of 100)	49%

Table 2: Game of 24 Results.

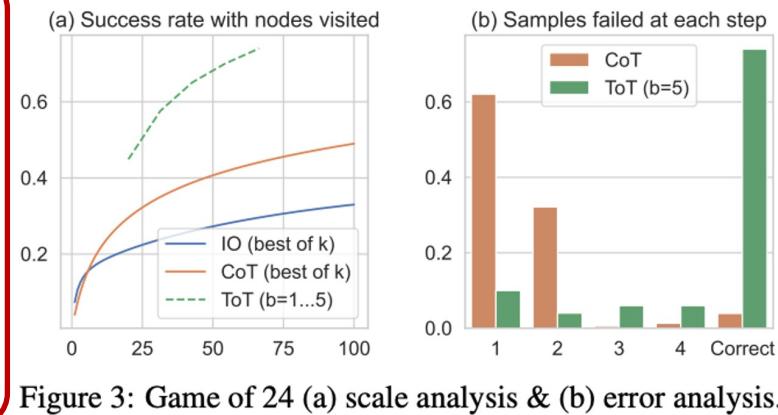
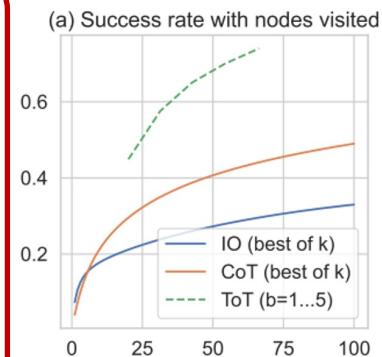


Figure 3: Game of 24 (a) scale analysis & (b) error analysis.

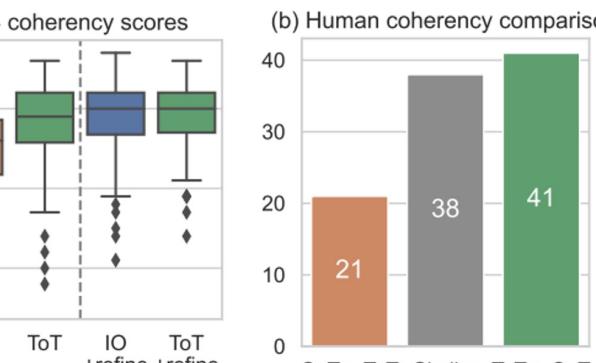
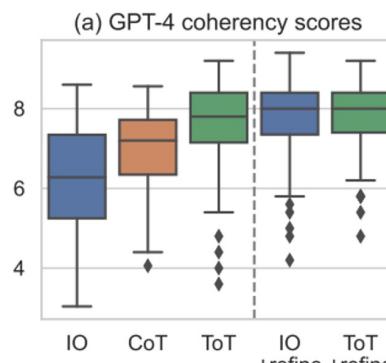


Figure 5: Creative Writing results.

Method	Success Rate (%)		
	Letter	Word	Game
IO	38.7	14	0
CoT	40.6	15.6	1
ToT (ours)	<b>78</b>	<b>60</b>	<b>20</b>
+best state	82.4	67.5	35
-prune	65.4	41.5	5
-backtrack	54.6	20	5

Table 3: Mini Crosswords results.

# „Tree of Thoughts: Deliberate Problem Solving with Large Language Models“

[\[Link\]](#)

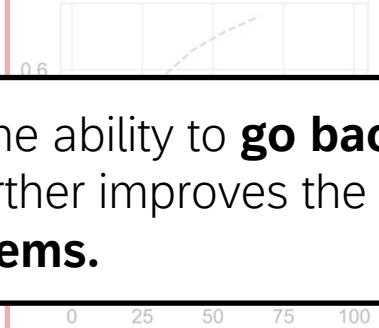
## Prompting Techniques

Method	Success
IO prompt	7.3%
CoT prompt	4.0%
ToT (ours)	78.0%
+refine	60.0%
-prune	54.6%
-backtrack	54.6%

Giving an LLM the ability to **go back and reconsider** old decisions further improves the ability to **solve complex problems**.

Table 2: Game of 24 Results.

(a) Success rate with nodes visited



(b) Samples failed at each step

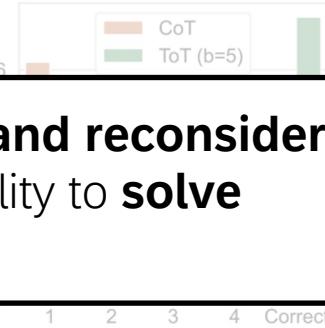


Figure 3: Game of 24 (a) scale analysis & (b) error analysis.

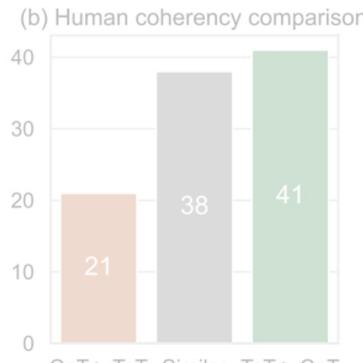
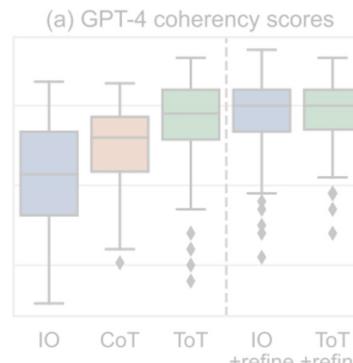


Figure 5: Creative Writing results.

Method	Success Rate (%)		
	Letter	Word	Game
IO	38.7	14	0
CoT	40.6	15.6	1
ToT (ours)	<b>78</b>	<b>60</b>	<b>20</b>
+best state	82.4	67.5	35
-prune	65.4	41.5	5
-backtrack	54.6	20	5

Table 3: Mini Crosswords results.

# „Toolformer: Language Models Can Teach Themselves to Use Tools“ [\[Link\]](#)

1.Meta AI

2.Universitat Pompeu Fabra

End2End

Tools for LLMs

[Presentation](#) by Michael in this series 17<sup>th</sup> March'23

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.

**Fine-tuning of GPT-J (6.7B) to use tools.  
Tool evaluation during generation.**

# „Tool Documentation Enables Zero-Shot Tool-Usage with Large Language Models“ [\[Link\]](#)

Tools for LLMs

Prompting Techniques

1. University of Washington
2. National Taiwan University
3. Google Cloud AI Research
4. Google Research

## Demonstration

Description: examples of questions and the tool-use plan.

- **Question:**  
Which property do these objects have in common?



### Tool-use Plan:

Text Detector → Knowledge Retriever → Solution Generator

- **Question:** [...]
- **Tool-use Plan:** [...]
- **Question:** [...]
- **Tool-use Plan:** [...]



## Documentation

Description: available tools and their functionalities.

- **Text Detector:**  
It detects the text in an image [...]
- **Knowledge Retriever:**  
It retrieves relevant knowledge [...]
- **Search Engine:**  
It searches the web for relevant info [...]
- **Image Captioner:**  
It generates a caption for an image [...]
- ...



Figure 2: Two types of knowledge for prompting LLMs for tool-use: Demonstrations (demos) and Documentations (docs). Demos consist of  $\langle$ input, output $\rangle$  pairs on input instructions and their corresponding output tool-use plans. They require manual efforts for careful curation on every new task, and the model performance can be sensitive to which demos are used [81, 12]. Many demos may also be necessary for good coverage when the number of tools scales up. On the other hand, docs provide descriptions for the tool functionality, and are usually organically available for tools.

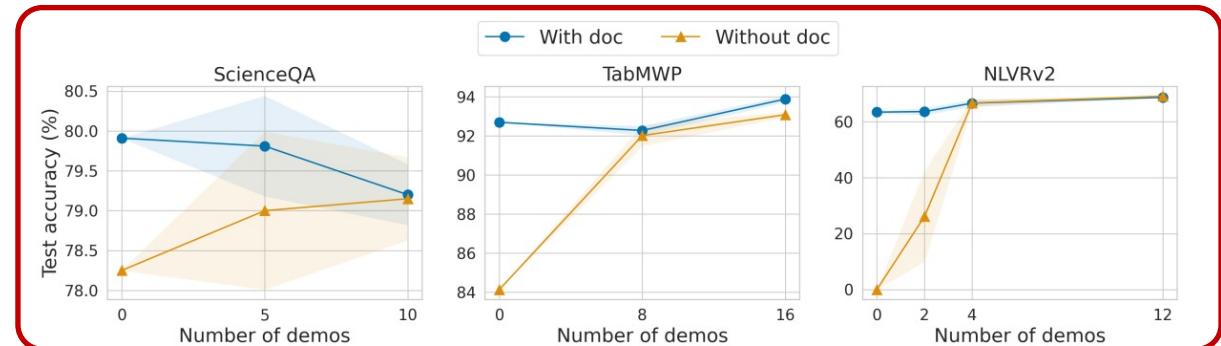


Figure 4: Tool-using performance with `gpt-3.5-turbo` on different benchmarks, which covers from language to vision modalities. We report results with and without documentation (doc) and demonstrations (demo), and their combinations. Clearly, with documentation only (upper-left blue dot) shows competitive performance across all datasets.

Compare demo-based vs documentation-based tool-augmented LLMs for complex task solving

# „Tool Documentation Enables Zero-Shot Tool-Usage with Large Language Models“ [\[Link\]](#)

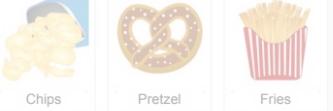
Tools for LLMs

Prompting Techniques

1. University of Washington
2. National Taiwan University
3. Google Cloud AI Research
4. Google Research

**Demonstration**

Description: examples of questions and the tool-use plan.

- Question:  
Which property do these objects have in common?  
  
Chips      Pretzel      Fries

**Tool-use Plan:**  
Text Detector → Knowledge Retriever → Solution Generator

- Question: [...]
- Tool-use Plan: [...]
- Question: [...]
- Tool-use Plan: [...]



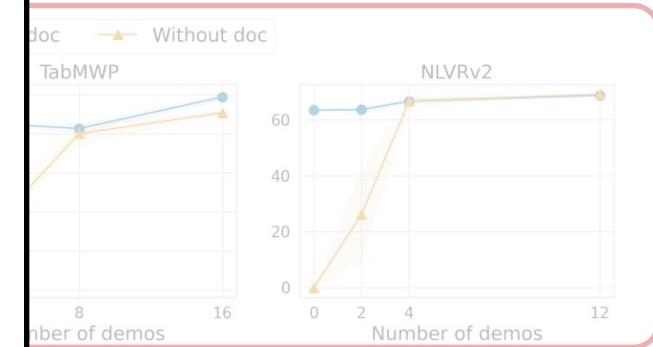
**Documentation**

**Tools** can be reasonably well controlled **by in-context learning with LLMs >100B**.

The **usefulness of few-shot examples** seems to be **limited** and sometimes even limits the model's ability to use the tools.

(In contrast to the Toolformer: In-context learning required to decouple tool evaluation and LLM generation)

Figure 2: Two types of knowledge for prompting LLMs. Documentations (docs). Demos consist of input and corresponding output tool-use plans. They require task, and the model performance can be sensitive to the number of demos. On the other hand, docs provide descriptions for the tool functionality, and are usually organically available for tools.



.5-turbo on different benchmarks, which covers results with and without documentation (doc) and . Clearly, with documentation only (upper-left blue datasets).

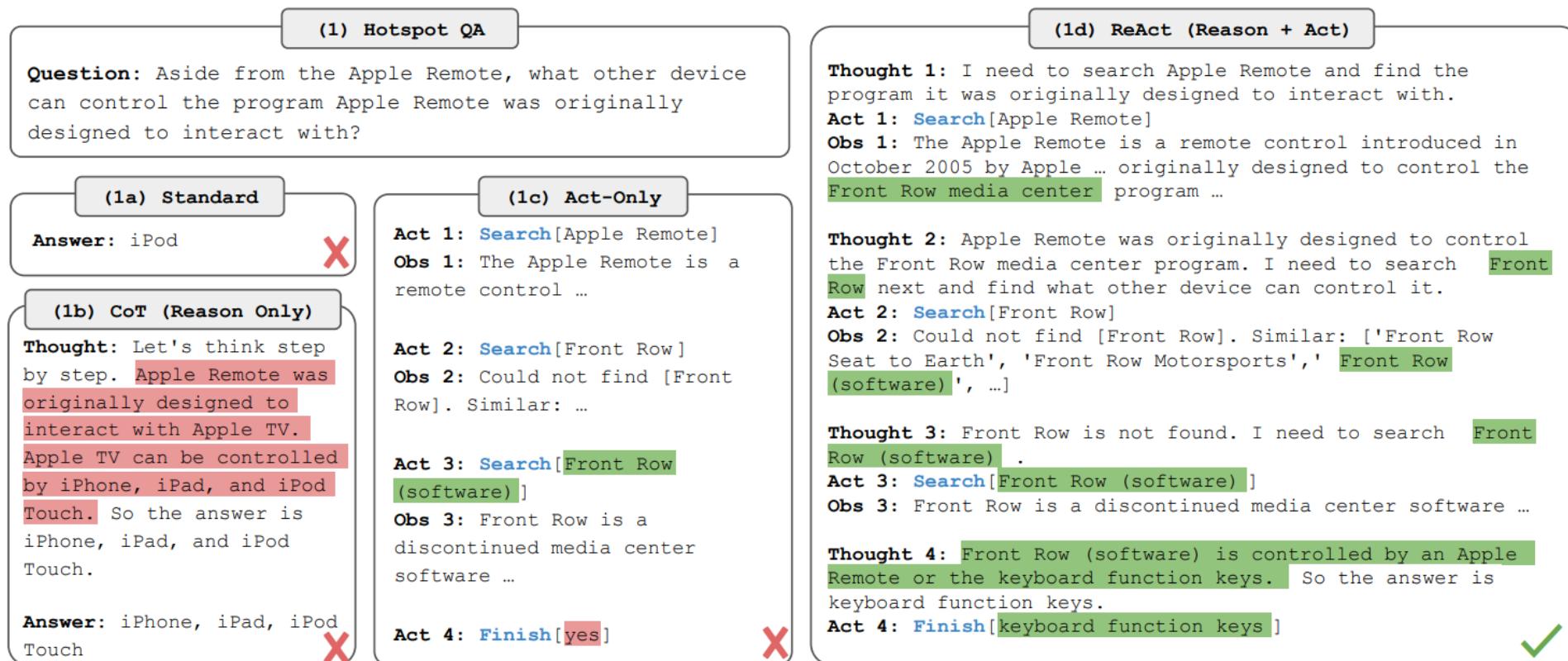
## Compare demo-based vs documentation-based tool-augmented LLMs for complex task solving

# „ReAct: Synergizing Reasoning and Acting in Language Models“ [\[Link\]](#)

Prompting Techniques

Agent Architecture

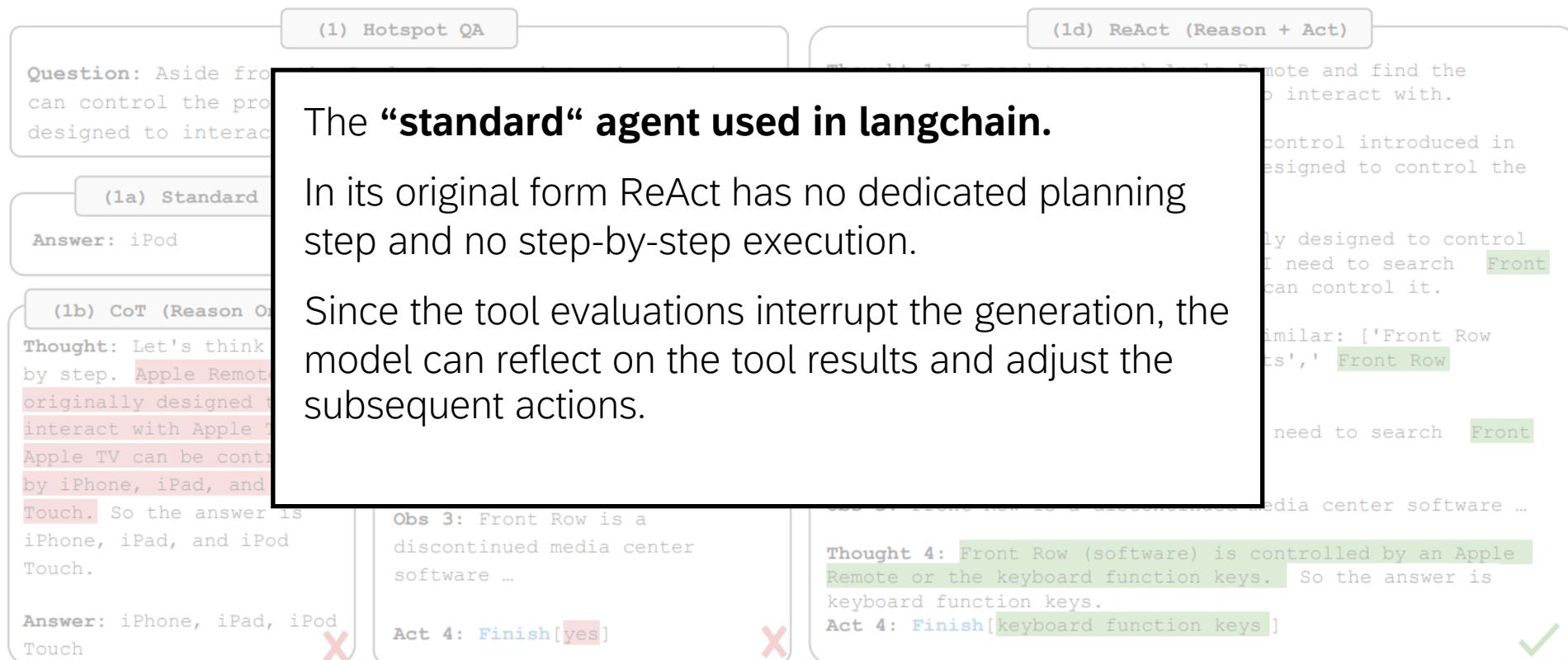
1.Princeton  
2.Google Research



# „ReAct: Synergizing Reasoning and Acting in Language Models“ [\[Link\]](#)

Prompting Techniques

Agent Architecture



# „BOLAA: Benchmarking and Orchestrating LLM-augmented Autonomous Agents“ [\[Link\]](#)

1.Salesforce

## Agent Architecture

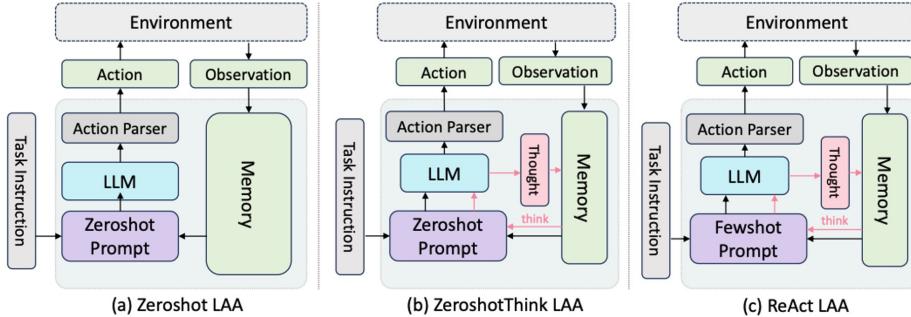


Figure 1: The LAA architectures for Zeroshot-LAA (ZS-LAA), ZeroshotThink LAA (ZST-LAA) and ReAct LAA. ZS-LAA generates actions from LLM with zeroshot prompt. ZST-LAA extends ZS-LAA with self-think. ReAct LAA advances ZST-LAA with fewshot prompt. They all resolve a given task by interacting with environment via actions to collect observations. Better view in colors.

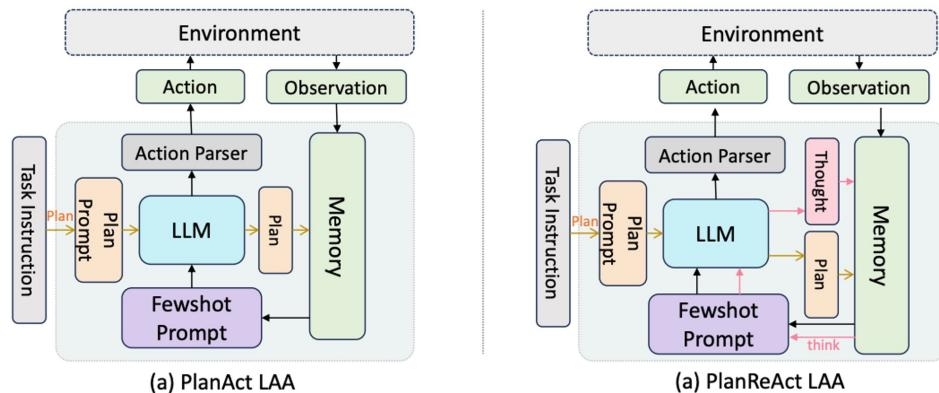


Figure 2: The LAA architectures for PlanAct LAA and PlanReAct LAA.

## End2End Evaluation

Table 3: Average reward in the HotPotQA environment. Len denotes the maximum context length. **Bold** results denote the best results in one row, *i.e.* best LAA architecture w.r.t. one LLM. Underline results denote the best performance in one column, *i.e.* best LLM regarding one LAA architecture.

LLM	Len.	LAA Architecture				
		ZS	ZST	ReAct	PlanAct	PlanReAct
fastchat-t5-3b	2k	0.0252	0.0067	0.0692	<b>0.1155</b>	0.0834
vicuna-7b	2k	<b>0.1339</b>	0.0797	0.0318	0.0868	0.0956
vicuna-13b	2k	0.1541	0.0910	<b>0.2637</b>	0.1754	0.2075
vicuna-33b	2k	0.2180	0.2223	<b>0.2602</b>	0.1333	0.2016
llama-2-7b	4k	0.0395	0.0207	<b>0.2624</b>	0.1780	0.1417
llama-2-13b	4k	0.1731	0.2313	<b>0.2521</b>	0.2192	0.2177
llama-2-70b	4k	0.2809	0.3207	<b>0.3558</b>	0.1424	0.1797
mpt-7b-instruct	8k	0.0982	0.0483	<b>0.1707</b>	0.1147	0.1195
mpt-30b-instruct	8k	0.1562	0.2141	<b>0.3261</b>	0.2224	0.2315
xgen-8k-7b-instruct	8k	0.1502	0.1244	<b>0.1937</b>	0.1116	0.1096
longchat-7b-16k	16k	0.0791	0.0672	<b>0.2161</b>	0.1296	0.0971
longchat-13b-16k	16k	0.1083	0.0562	<b>0.2387</b>	0.1623	0.1349
text-davinci-003	4k	0.3430	0.3304	<b>0.4503</b>	0.3577	0.4101
gpt-3.5-turbo	4k	<b>0.3340</b>	0.3254	0.3226	0.2762	0.3192
gpt-3.5-turbo-16k-0613	16k	<b>0.3027</b>	0.2264	0.1859	0.2113	0.2251

# „BOLAA: Benchmarking and Orchestrating LLM-augmented Autonomous Agents“ [\[Link\]](#)

1.Salesforce

## Agent Architecture

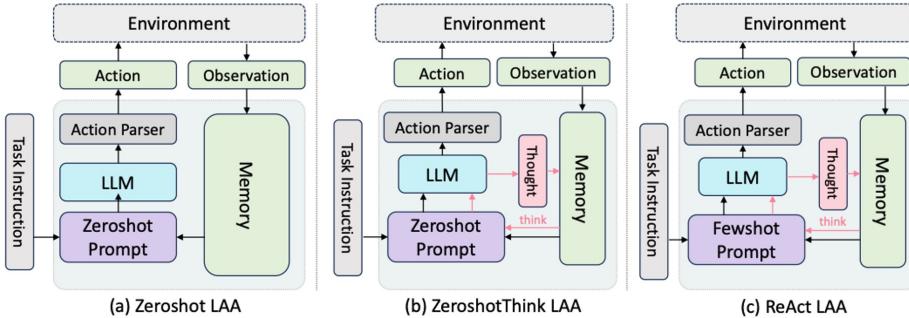


Figure 1: The LAA architectures for Zeroshot-LAA (ZS-LAA), ZeroshotThink LAA (ZST-LAA) and ReAct LAA. ZS-LAA generates actions from LLM with zeroshot prompt. ZST-LAA extends ZS-LAA with self-think. ReAct LAA advances ZST-LAA with fewshot prompt. They all resolve a given task by interacting with environment via actions to collect observations. Better view in colors.

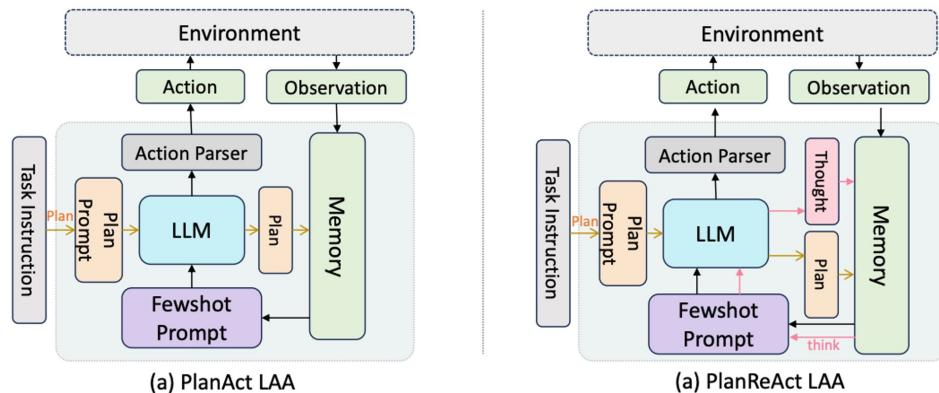


Figure 2: The LAA architectures for PlanAct LAA and PlanReAct LAA.

## End2End Evaluation

Table 1: Average reward in the WebShop environment. Len denotes the maximum context length. **Bold** results denote the best results in one row, *i.e.* best LAA architecture w.r.t. one LLM. Underline results denote the best performance in one column, *i.e.* best LLM regarding one LAA architecture.

LLM	Len.	LAA Architecture					BOLAA
		ZS	ZST	ReAct	PlanAct	PlanReAct	
fastchat-t5-3b	2k	0.3971	0.2832	0.3098	<b>0.3837</b>	0.1507	<b>0.5169</b>
vicuna-7b	2k	0.0012	0.0002	<b>0.1033</b>	0.0555	0.0674	0.0604
vicuna-13b	2k	0.0340	0.0451	0.1509	<b>0.3120</b>	0.4127	<b>0.5350</b>
vicuna-33b	2k	0.1356	0.2049	0.1887	<b>0.3692</b>	0.3125	<b>0.5612</b>
llama-2-7b	4k	0.0042	0.0068	0.1248	<b>0.3156</b>	0.2761	<b>0.4648</b>
llama-2-13b	4k	0.0662	0.0420	0.2568	<b>0.4892</b>	0.4091	0.3716
llama-2-70b	4k	0.0122	0.0080	<b>0.4426</b>	0.2979	0.3770	<b>0.5040</b>
mpt-7b-instruct	8k	0.0001	0.0001	0.0573	0.0656	<b>0.1574</b>	0.0632
mpt-30b-instruct	8k	0.1664	0.1255	0.3119	0.3060	<b>0.3198</b>	<b>0.4381</b>
xgen-8k-7b-instruct	8k	0.0001	0.0015	0.0685	<b>0.1574</b>	0.1004	<b>0.3697</b>
longchat-7b-16k	16k	0.0165	0.0171	0.069	0.0917	0.1322	<b>0.1964</b>
longchat-13b-16k	16k	0.0007	0.0007	0.2373	0.3978	<b>0.4019</b>	0.3205
text-davinci-003	4k	0.5292	0.5395	<b>0.5474</b>	0.4751	0.4912	<b>0.6341</b>
gpt-3.5-turbo	4k	0.5061	0.5057	0.5383	0.4667	<b>0.5483</b>	<b>0.6567</b>
gpt-3.5-turbo-16k	16k	<b>0.5657</b>	0.5642	0.4898	0.4565	<b>0.5607</b>	<b>0.6541</b>

# „BOLAA: Benchmarking and Orchestrating LLM-augmented Autonomous Agents“ [\[Link\]](#)



Figure 1: The LAA architectures for Zeroshot-LAA (ZS-) and ReAct LAA. ZS-LAA generates actions from LLM w/o ZS-LAA with self-think. ReAct LAA advances ZST-LAA by given task by interacting with environment via actions to complete task.

Table 1: Average reward in the WebShop environment. Len denotes the maximum context length. **Bold** results denote the best results in one row, i.e. best LAA architecture w.r.t. one LLM. Underline denotes the best LLM regarding one LAA architecture.

**There is no obvious best agent architecture.  
Open source models can't compete with the OpenAI models when used for AI agents.**

Similar findings in:

“AgentBench: Evaluating LLMs as Agents” [\[Link\]](#)

LAA Architecture	eAct	PlanAct	PlanReAct	BOLAA
mpt-30b-instruct	3098	<b>0.3837</b>	0.1507	<b>0.5169</b>
xgen-8k-7b-instruct	<b>1033</b>	0.0555	0.0674	0.0604
longchat-7b-16k	1509	<b>0.3120</b>	0.4127	<b>0.5350</b>
longchat-13b-16k	1887	<b>0.3692</b>	0.3125	<b>0.5612</b>
text-davinci-003	1248	<b>0.3156</b>	0.2761	<b>0.4648</b>
gpt-3.5-turbo	2568	<b>0.4892</b>	0.4091	0.3716
gpt-3.5-turbo-16k	4426	0.2979	0.3770	<b>0.5040</b>
	0573	0.0656	<b>0.1574</b>	0.0632
mpt-30b-instruct	8k	0.1664	0.1255	0.3198
xgen-8k-7b-instruct	8k	0.0001	0.0015	<b>0.1574</b>
longchat-7b-16k	16k	0.0165	0.0171	0.1004
longchat-13b-16k	16k	0.0007	0.0007	<b>0.4019</b>
text-davinci-003	4k	0.5292	0.5395	0.4912
gpt-3.5-turbo	4k	0.5061	0.5057	<b>0.5483</b>
gpt-3.5-turbo-16k	16k	<b>0.5657</b>	<b>0.5642</b>	<b>0.6567</b>

mpt-30b-instruct

xgen-8k-7b-instruct

longchat-7b-16k

longchat-13b-16k

text-davinci-003

gpt-3.5-turbo

gpt-3.5-turbo-16k

8k

8k

16k

16k

4k

4k

16k

0.1664

0.0001

0.0165

0.0007

0.5292

0.5061

0.5657

0.1255

0.0015

0.0171

0.0007

0.5395

0.5057

0.5642

0.3119

0.0685

0.069

0.2373

0.3978

0.4751

0.4667

0.3060

0.0917

0.3978

0.4565

0.4565

0.5607

0.3198

0.1004

0.1322

0.3205

0.4912

0.5483

0.6541

0.4381

0.1964

0.3205

0.6341

0.6567

0.6541

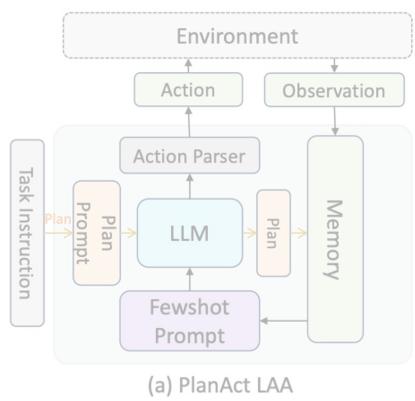


Figure 2: The LAA architectures for PlanAct LAA and PlanReAct LAA.

# „AgentBench Evaluation LLMs as AI Agents“ [\[Link\]](#)

1.Tsinghua University  
2.Ohio State University  
3.UC Berkeley

## End2End Evaluation

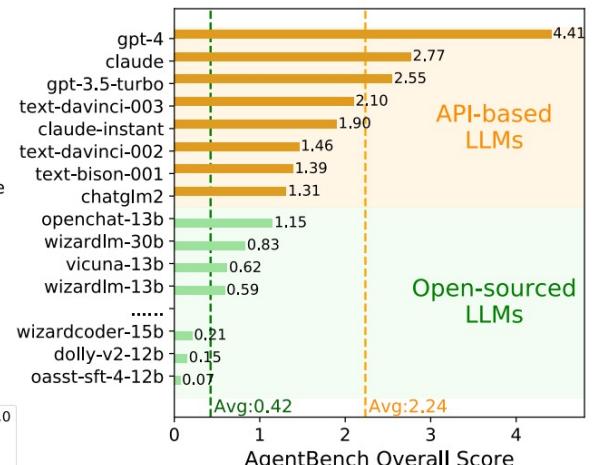
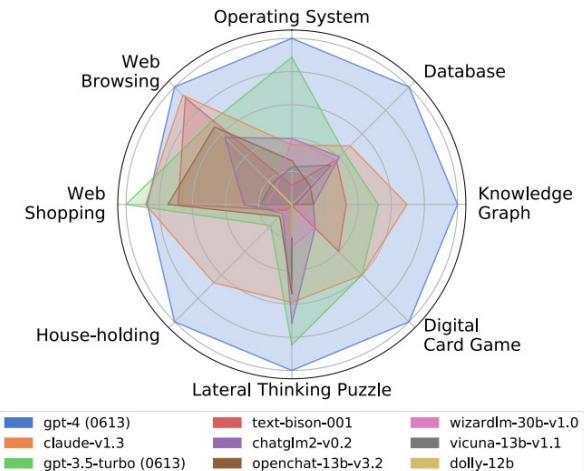
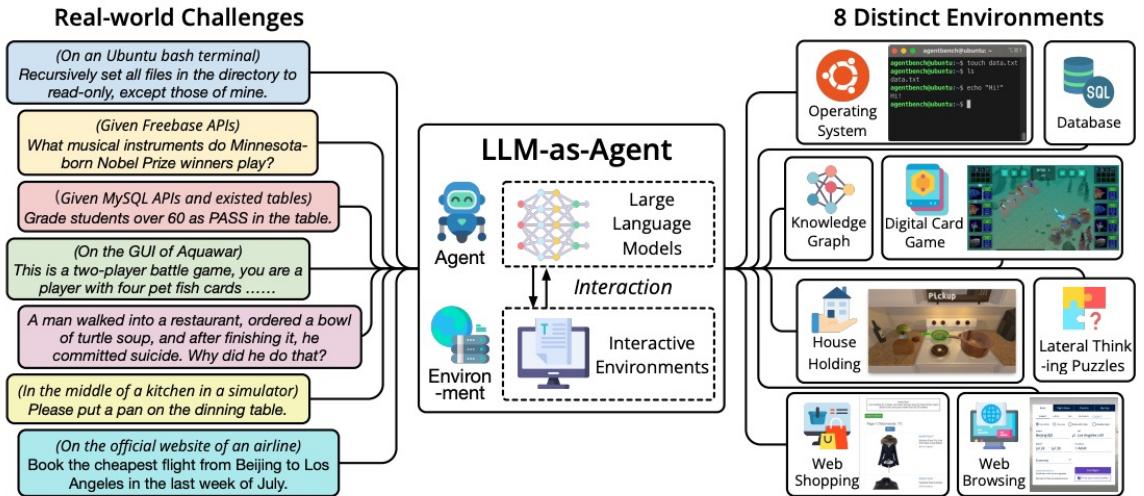


Figure 2: AgentBench is the first systematic benchmark to evaluate LLM-as-Agent on a wide array of real-world challenges and 8 distinct environments. In total, 25 LLMs are examined in its first edition.

Figure 1: An overview of LLMs on AgentBench. While LLMs begin to manifest their proficiency in LLM-as-Agent, gaps between models and the distance towards practical usability are significant.

# „AgentBench Evaluation LLMs as AI Agents“ [\[Link\]](#)

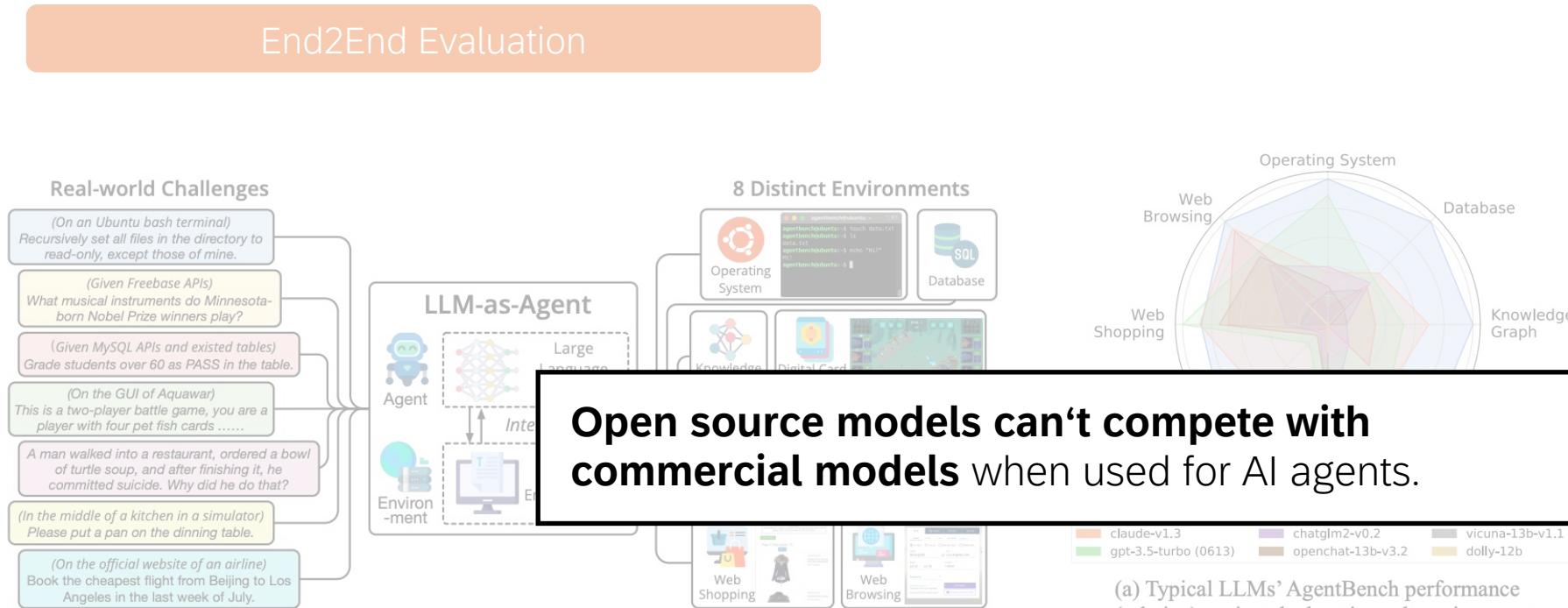


Figure 2: AgentBench is the first systematic benchmark to evaluate LLM-as-Agent on a wide array of real-world challenges and 8 distinct environments. In total, 25 LLMs are examined in its first edition.

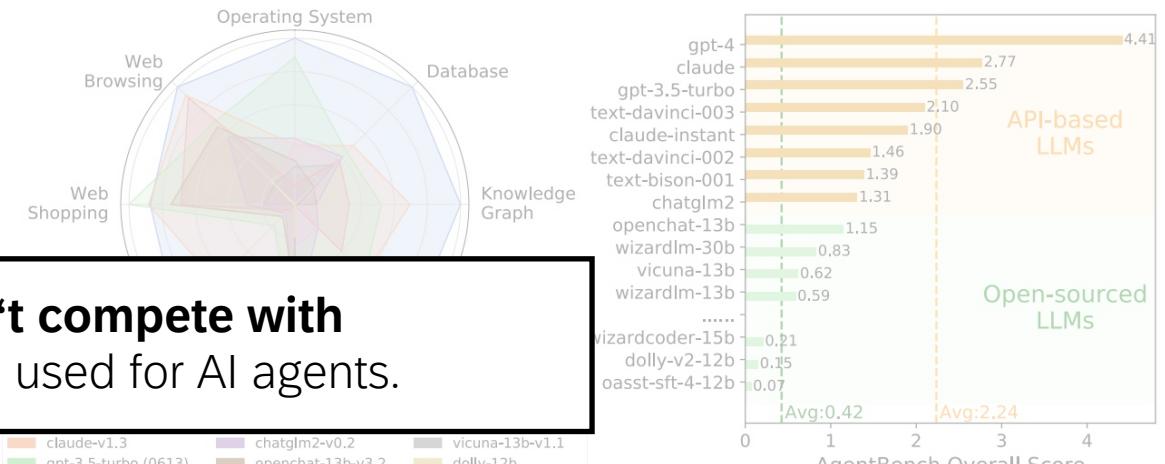


Figure 1: An overview of LLMs on AgentBench. While LLMs begin to manifest their proficiency in LLM-as-Agent, gaps between models and the distance towards practical usability are significant.

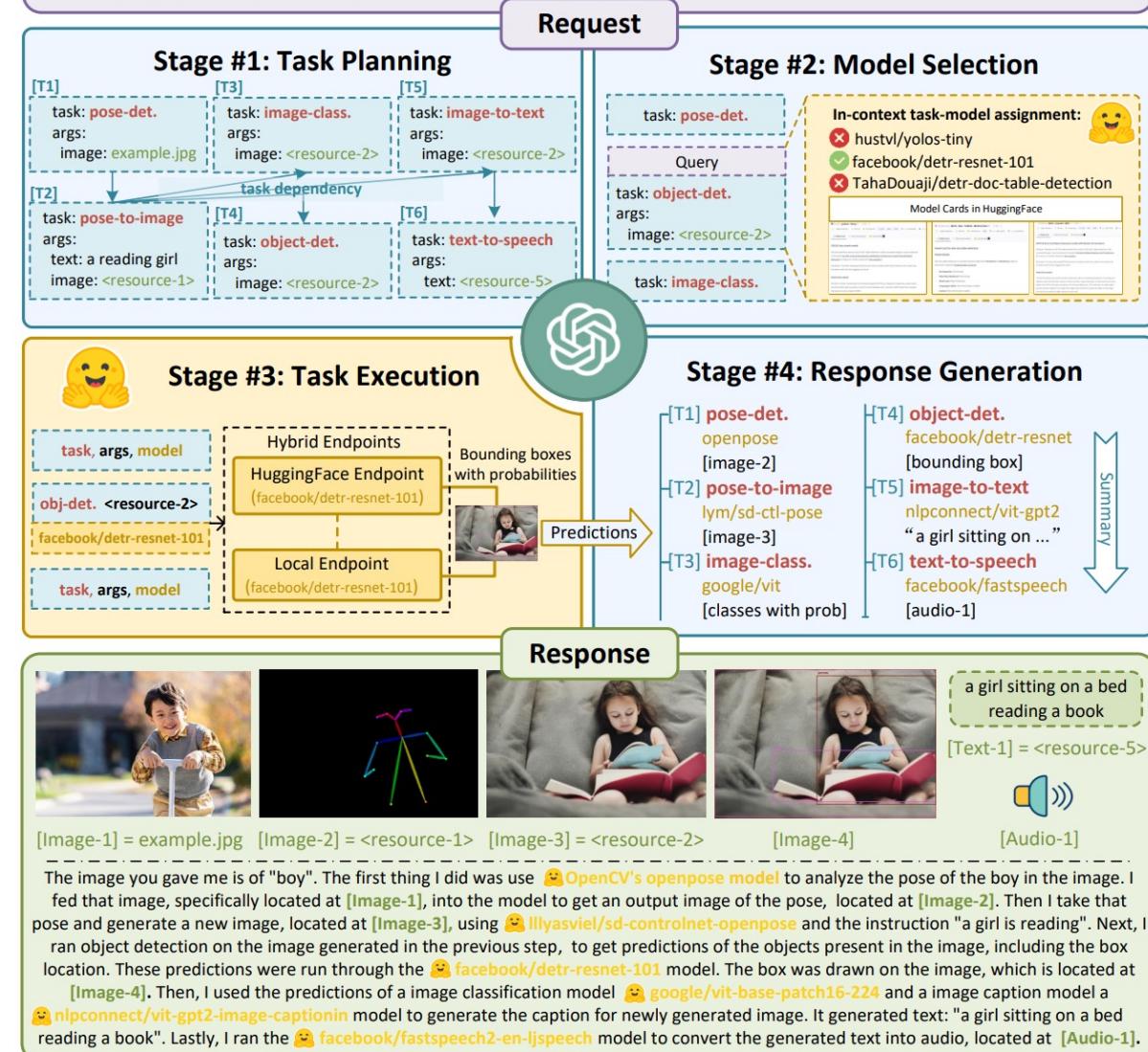
# „HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in Hugging Face“ [\[Link\]](#)

## Tools for LLMs / LLMs as Tools

1.Zhejiang University

2.Microsoft Research Asia

Please generate an image where a girl is reading a book, and her pose is the same as the boy in the image example.jpg, then please describe the new image with your voice.



# „HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in Hugging Face“ [\[Link\]](#)

Tools for LLMs / LLMs as Tools

**LLMs** can be used as **controllers for other AI models**.

Their **reasoning and problem-solving abilities** work in areas **beyond NLP**.

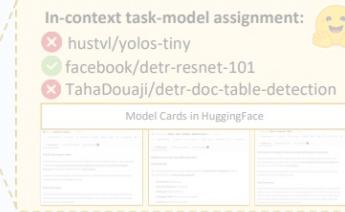
Please generate an image where a girl is reading a book, and her pose is the same as the boy in the image example.jpg, then please describe the new image with your voice.

Request

Stage #1: Task Planning

[T1]	task: pose-det. args: image: example.jpg	[T3]	task: image-class. args: image: <resource-2>	[T5]	task: image-to-text args: image: <resource-2>
					task: pose-det. args: image: <resource-2>

Stage #2: Model Selection



#4: Response Generation



Response



[Image-1] = example.jpg    [Image-2] = <resource-1>    [Image-3] = <resource-2>    [Image-4]    [Audio-1]

The image you gave me is of "boy". The first thing I did was use [OpenCV's openpose model](#) to analyze the pose of the boy in the image. I fed that image, specifically located at [Image-1], into the model to get an output image of the pose, located at [Image-2]. Then I take that pose and generate a new image, located at [Image-3], using [ilyasvie/sd-controlnet-openpose](#) and the instruction "a girl is reading". Next, I ran object detection on the image generated in the previous step, to get predictions of the objects present in the image, including the box location. These predictions were run through the [facebook/detr-resnet-101](#) model. The box was drawn on the image, which is located at [Image-4]. Then, I used the predictions of a image classification model [google/vit-base-patch16-224](#) and a image caption model a [nlpconnect/vit-gpt2-image-captionin](#) model to generate the caption for newly generated image. It generated text: "a girl sitting on a bed reading a book". Lastly, I ran the [facebook/fastspeech2-en-ljspeech](#) model to convert the generated text into audio, located at [Audio-1].

# „Voyager: An Open-Ended Embodied Agent with Large Language Models“

[\[Link\]](#)

1.NVIDIA  
2.Caltech  
3.UT Austin  
4.Stanford  
5.ASU

## Agent Architecture

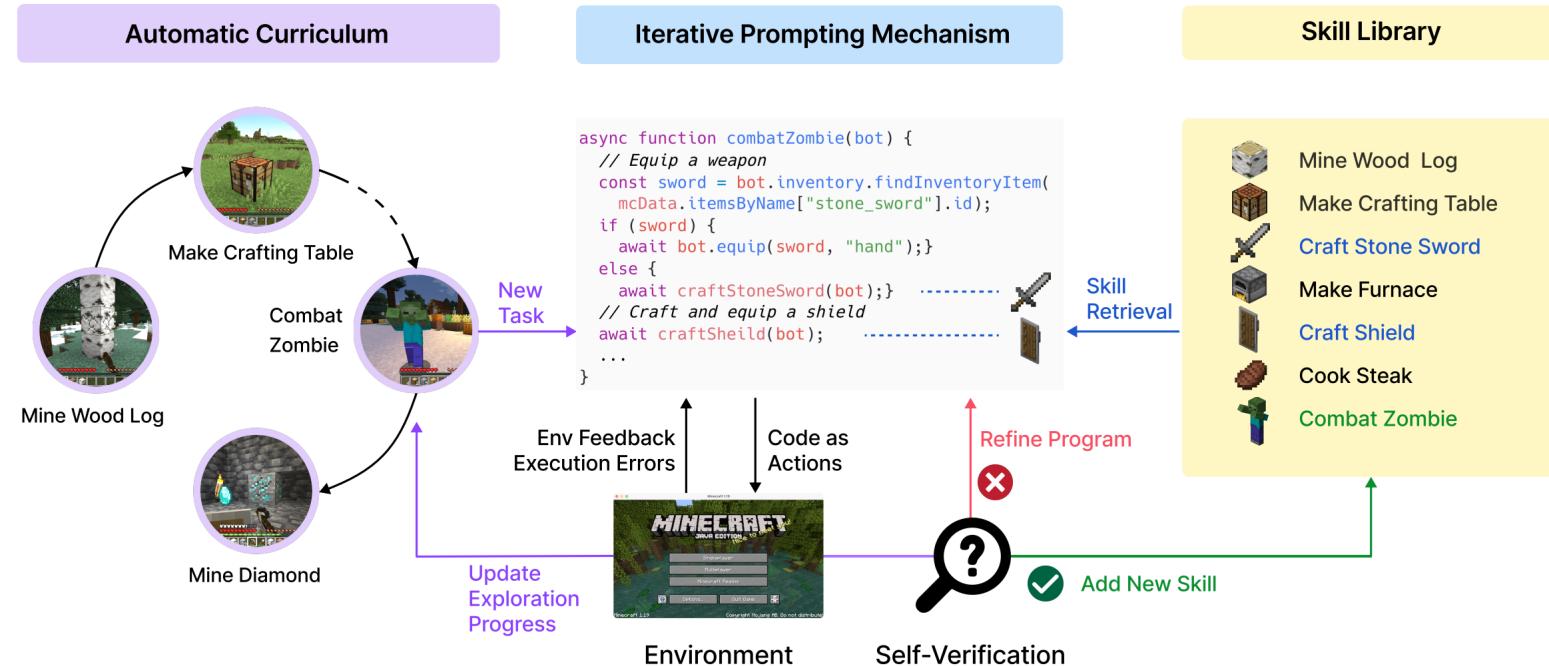


Figure 2: VOYAGER consists of three key components: an automatic curriculum for open-ended exploration, a skill library for increasingly complex behaviors, and an iterative prompting mechanism that uses code as action space.

# „Voyager: An Open-Ended Embodied Agent with Large Language Models“

[\[Link\]](#)

- 1.NVIDIA
- 2.Caltech
- 3.UT Austin
- 4.Stanford
- 5.ASU

## Agent Architecture

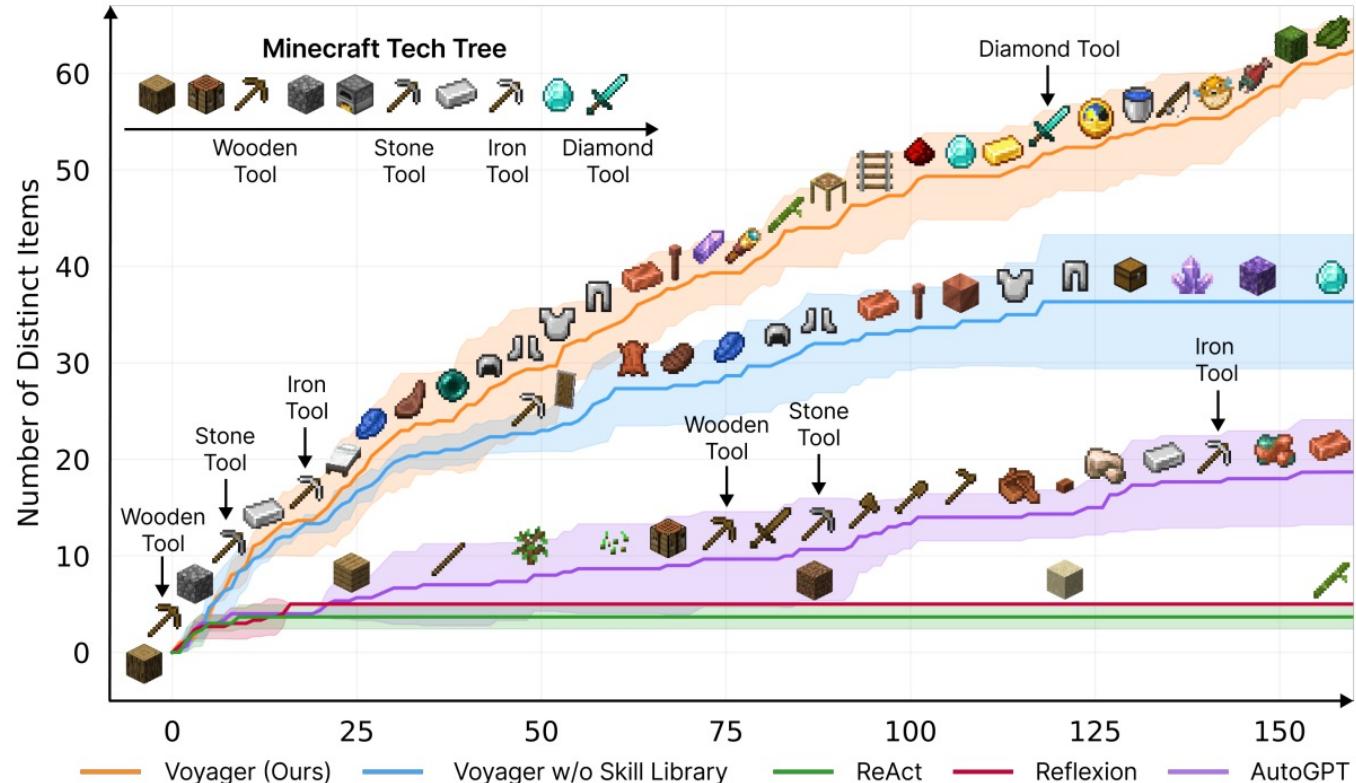


Figure 1: VOYAGER discovers new Minecraft items and skills continually by self-driven exploration, significantly outperforming the baselines. X-axis denotes the number of prompting iterations.

# „Voyager: An Open-Ended Embodied Agent with Large Language Models“

[\[Link\]](#)

## Agent Architecture

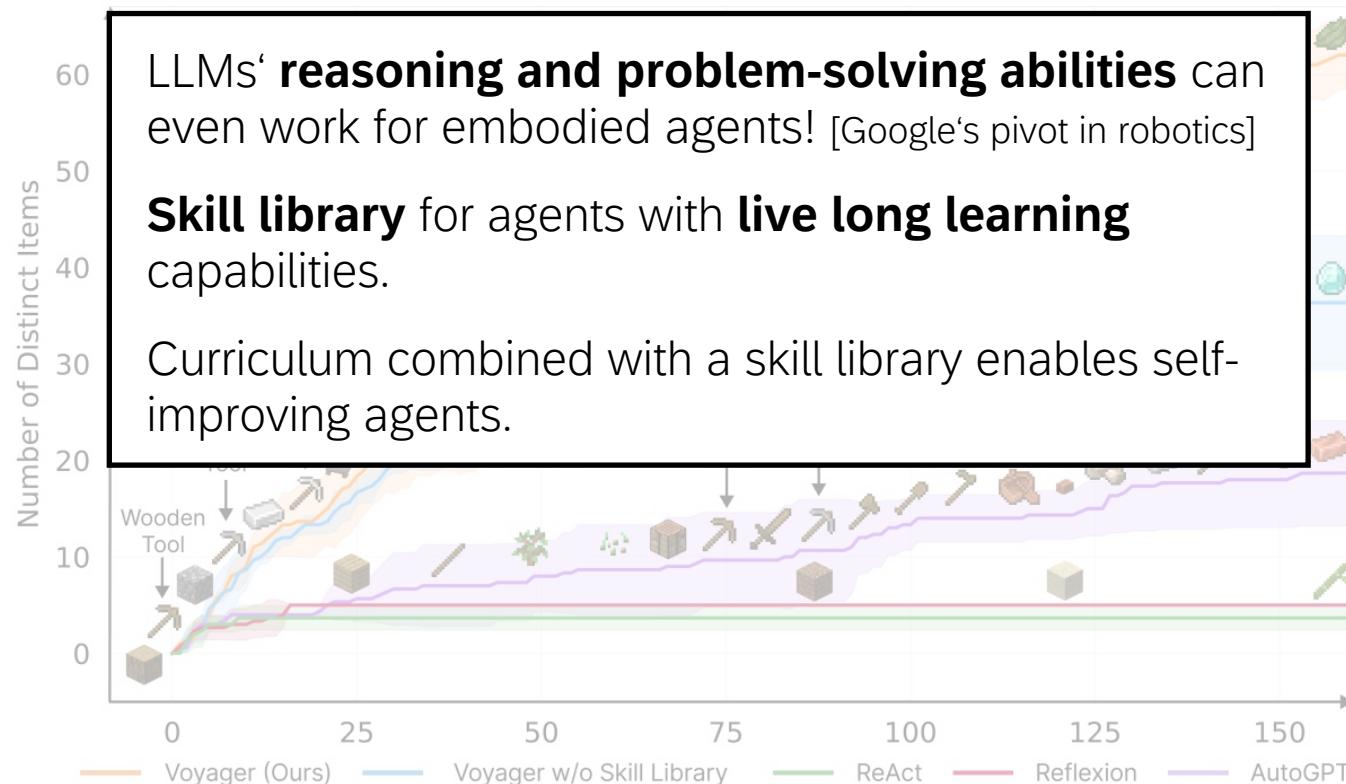
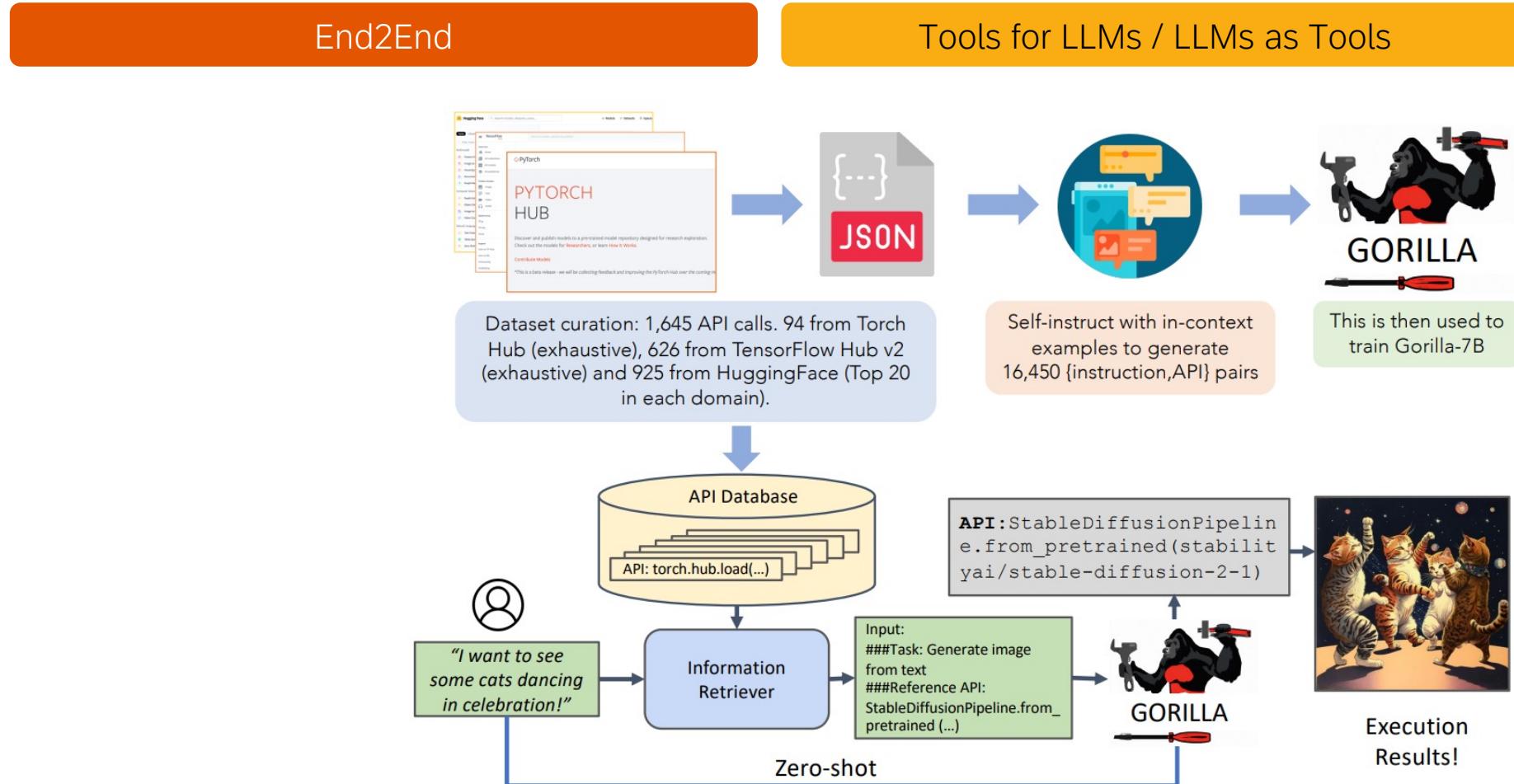


Figure 1: VOYAGER discovers new Minecraft items and skills continually by self-driven exploration, significantly outperforming the baselines. X-axis denotes the number of prompting iterations.

# „Gorilla: Large Language Model Connected with Massive APIs“ [\[Link\]](#)

1.UC Berkeley  
2.Microsoft Research



**Fine-tuning LLaMA (7B) for specific APIs  
with API retriever used during training**

# „Gorilla: Large Language Model Connected with Massive APIs“ [\[Link\]](#)

1.UC Berkeley  
2.Microsoft Research

End2End

Tools for LLMs / LLMs as Tools

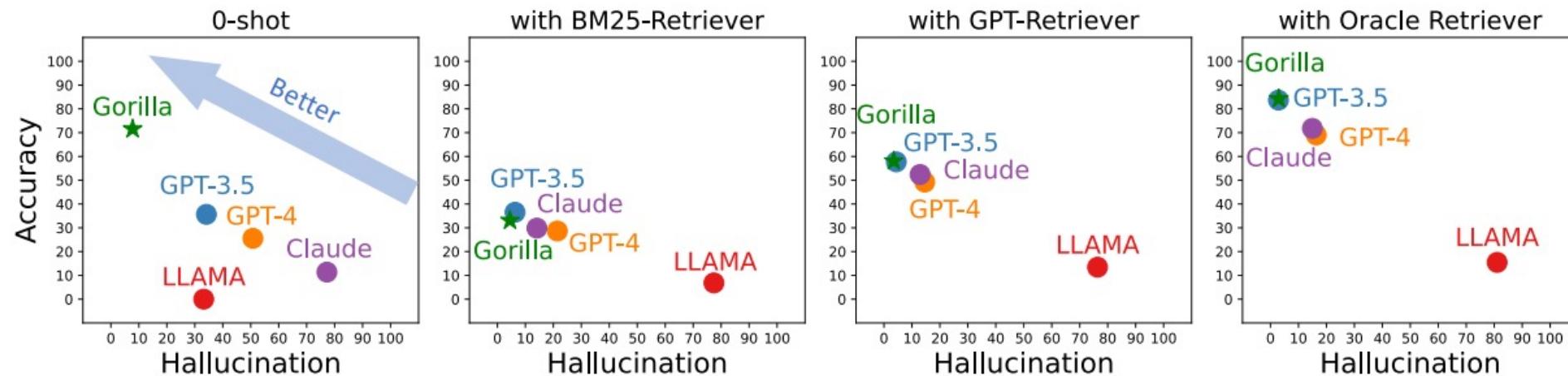


Figure 2: **Accuracy (vs) hallucination** in four settings, that is, *zero-shot* (i.e., without any retriever), and *with retrievers*. BM25 and GPT are commonly used retrievers and the oracle retriever returns relevant documents at 100%, indicating an upper bound. Higher in the graph (higher accuracy) and to the left is better (lower hallucination). Across the entire dataset, our model, Gorilla, improves accuracy while reducing hallucination.

# „Gorilla: Large Language Model Connected with Massive APIs“ [Link]

End2End

Tools for LLMs / LLMs as Tools

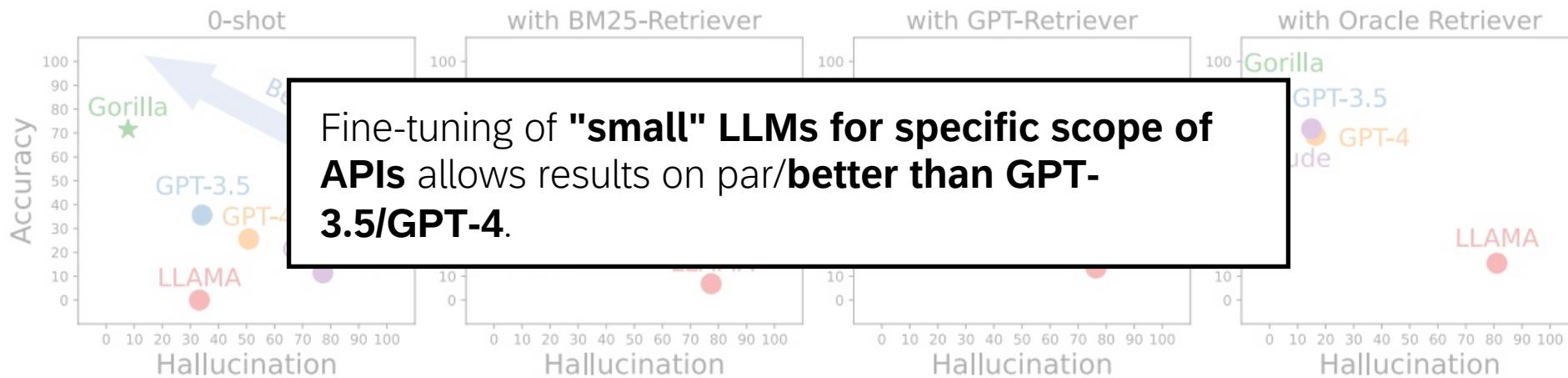


Figure 2: **Accuracy (vs) hallucination** in four settings, that is, *zero-shot* (i.e., without any retriever), and *with retrievers*. BM25 and GPT are commonly used retrievers and the oracle retriever returns relevant documents at 100%, indicating an upper bound. Higher in the graph (higher accuracy) and to the left is better (lower hallucination). Across the entire dataset, our model, Gorilla, improves accuracy while reducing hallucination.

# „ToolLLM: Facilitating Large Language Models to Master 16000+ Real-World APIs“ [\[Link\]](#)

End2End

Tools for LLMs / LLMs as Tools

- 1.Tsinghua University ModelBest Inc.
- 2.Renmin University of China
- 3.Yale University
- 4.WeChat AI, Tencent Inc.
- 5.Zhihu Inc.

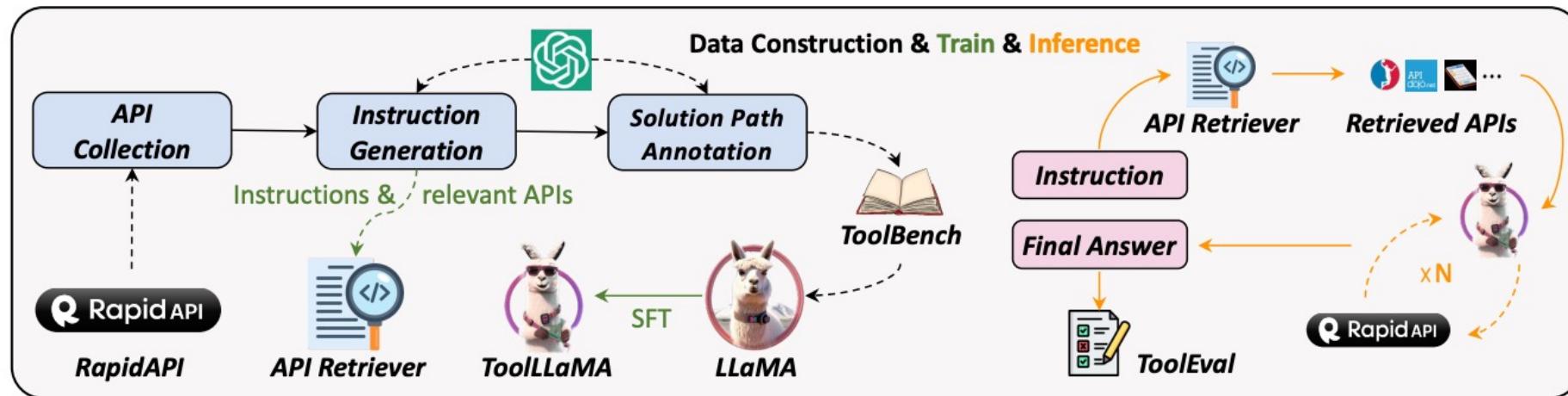


Figure 1: Three phases of constructing ToolBench and how we train our API retriever and ToolLLaMA. During inference of an instruction, the API retriever recommends relevant APIs to ToolLLaMA, which performs multiple rounds of API calls to derive the final answer. The whole reasoning process is evaluated by ToolEval.

# „ToolLLM: Facilitating Large Language Models to Master 16000+ Real-World APIs“ [\[Link\]](#)

End2End

Tools for LLMs / LLMs as Tools

- 1.Tsinghua University ModelBest Inc.
- 2.Renmin University of China
- 3.Yale University
- 4.WeChat AI, Tencent Inc.
- 5.Zhihu Inc.

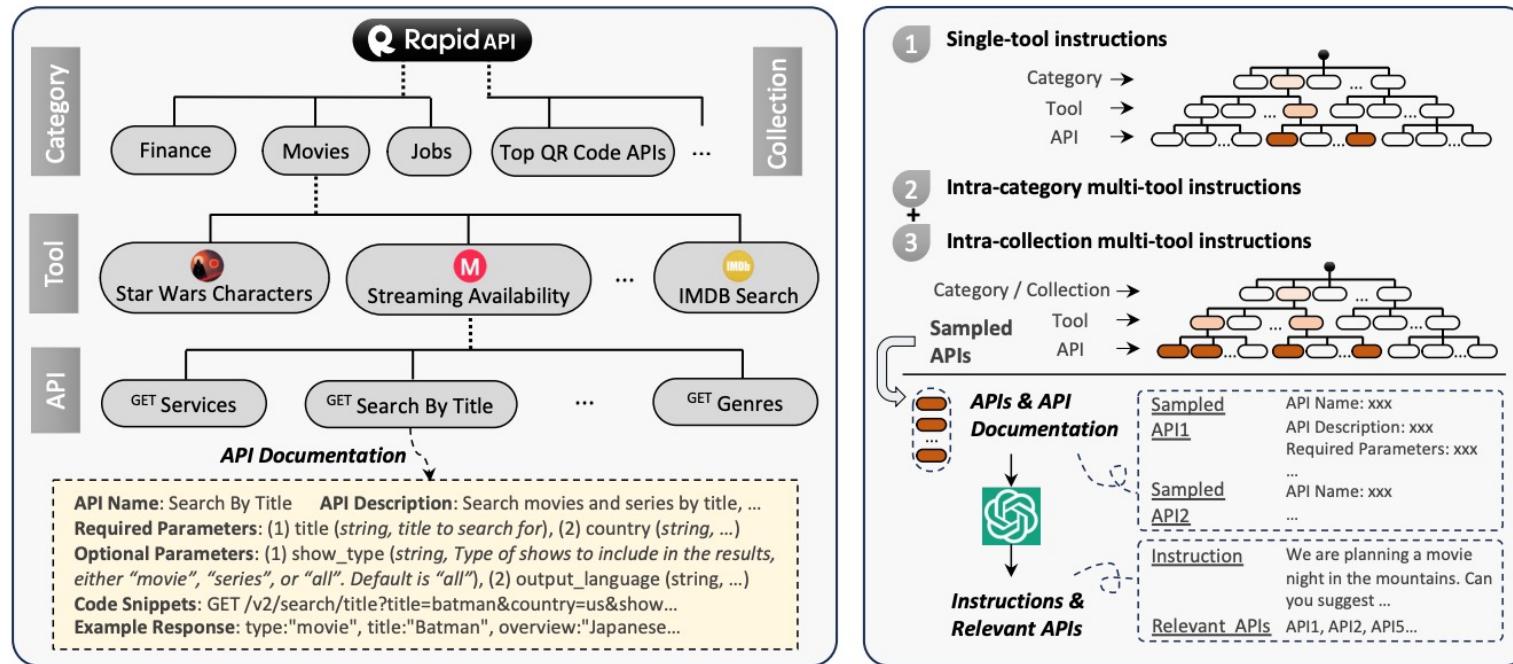


Figure 3: The hierarchy of RapidAPI (left) and the process of instruction generation (right).

# „ToolLLM: Facilitating Large Language Models to Master 16000+ Real-World APIs“ [\[Link\]](#)

End2End

Tools for LLMs / LLMs as Tools

- 1.Tsinghua University ModelBest Inc.
- 2.Renmin University of China
- 3.Yale University
- 4.WeChat AI, Tencent Inc.
- 5.Zhihu Inc.

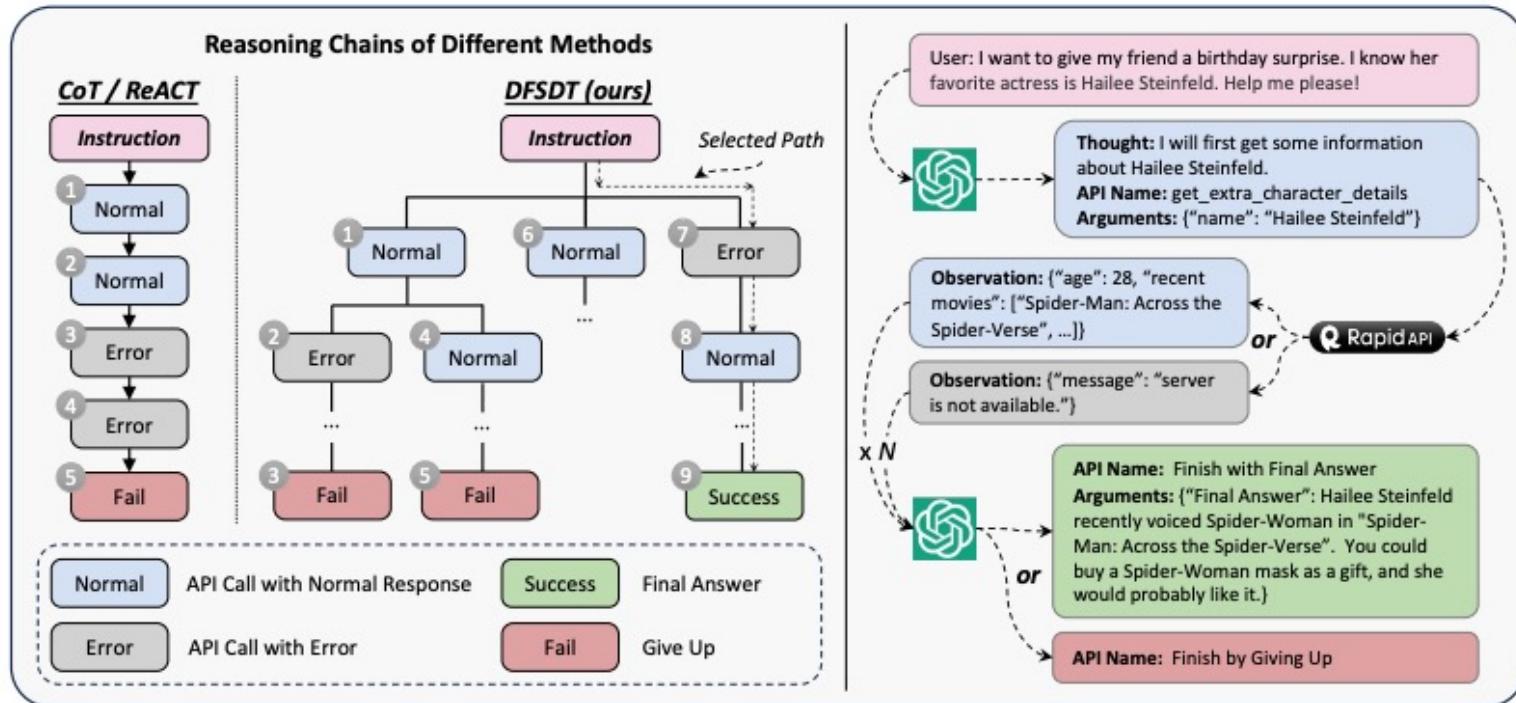


Figure 4: A comparison of our DFSDT and conventional CoT or ReACT during model reasoning (left). We show part of the solution path annotation process using ChatGPT (right).

# „ToolLLM: Facilitating Large Language Models to Master 16000+ Real-World APIs“ [\[Link\]](#)

End2End

Tools for LLMs / LLMs as Tools

- 1.Tsinghua University ModelBest Inc.
- 2.Renmin University of China
- 3.Yale University
- 4.WeChat AI, Tencent Inc.
- 5.Zhihu Inc.

Model	I1-Inst.		I1-Tool		I1-Cat.		I2-Inst.		I2-Cat.		I3-Inst.		Average	
	Pass	Win												
ChatGPT-ReACT	56.0	-	62.0	-	66.0	-	28.0	-	22.0	-	30.0	-	44.0	-
Vicuna (ReACT & DFSDT)	0.0	-	0.0	-	0.0	-	0.0	-	0.0	-	0.0	-	0.0	-
Alpaca (ReACT & DFSDT)	0.0	-	0.0	-	0.0	-	0.0	-	0.0	-	0.0	-	0.0	-
Text-Davinci-003-DFSDT	53.0	46.0	58.0	38.0	61.0	39.0	38.0	46.0	38.0	45.0	39.0	48.0	47.8	43.7
ChatGPT-DFSDT	<b>78.0</b>	<b>68.0</b>	<b>84.0</b>	<b>59.0</b>	<b>89.0</b>	<b>57.0</b>	<b>51.0</b>	<b>78.0</b>	<b>58.0</b>	<b>77.0</b>	<b>57.0</b>	<b>77.0</b>	<b>69.6</b>	<b>69.3</b>
ToolLLaMA-DFSDT	<u>68.0</u>	<u>68.0</u>	<u>80.0</u>	<u>59.0</u>	<u>75.0</u>	<u>56.0</u>	<u>47.0</u>	<u>75.0</u>	<u>56.0</u>	<u>80.0</u>	<u>40.0</u>	<u>72.0</u>	<u>61.0</u>	<u>68.3</u>

Table 4: Main experiments on the test set of ToolBench. Win rate is calculated by comparing each model with ChatGPT-ReACT. A win rate higher than 50% means the model performs better than ChatGPT-ReACT.

# „ToolLLM: Facilitating Large Language Models to Master 16000+ Real-World APIs“ [\[Link\]](#)

End2End

Tools for LLMs / LLMs as Tools

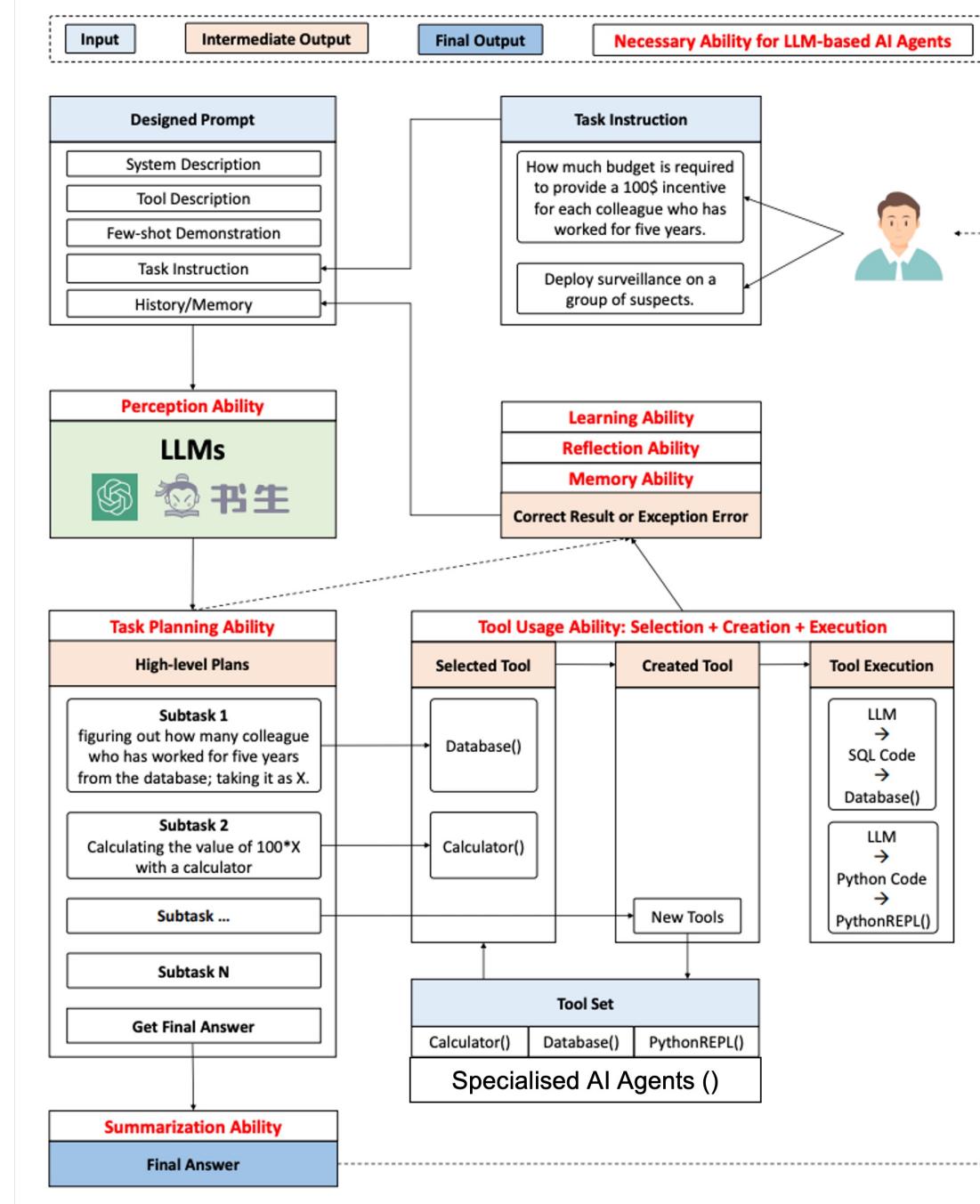
Model	Pass	Win
ChatGPT-ReAct	44.0	-
Vicuna (ReACT)	0.0	-
Alpaca (ReACT)	0.0	-
Text-Davinci-003	47.8	43.7
ChatGPT-DFS	<b>69.6</b>	<b>69.3</b>
ToolLLaMA-DT	61.0	68.3

Table 4: Main results for each model with ReACT.

**ToolLLM give a clear recipe on how to bootstrap a dataset to train a small LLM to generate API calls for a large number of tools (~3500 tools, 16500 endpoints).**

**Using a tree-based problem solving framework** (similar to „Tree of Thought“) also **improves tool usage.**

# Take Aways



# Take Aways

**Formulating a complete plan before solving** and giving the model a **chance to reflect** on it seems to **help the model reason significantly**.

**Tools** can be reasonably well controlled **by in-context learning (only providing documentation)** with LLMs >100B.

The **usefulness of few-shot examples** seems to be **limited** and sometimes even limits the model's ability to use the tools.

LLMs' **reasoning and problem-solving abilities** can even work for embodied agents! [Google's pivot in robotics]

**Skill library** for agents with **live long learning** capabilities.

Giving an LLM the ability to **go back and reconsider** old decisions further improves the ability to **solve complex problems**.

There is **no obvious best agent architecture**.

**Open source models can't compete with commercial models** when used for AI agents.

**ToolLLM give a clear recipe on how to bootstrap a dataset to train a small LLM to generate API calls for a large number of tools** (~3500 tools, 16500 endpoints).

Using a **tree-based problem solving framework** (similar to „Tree of Thought“) also **improves tool usage**.