



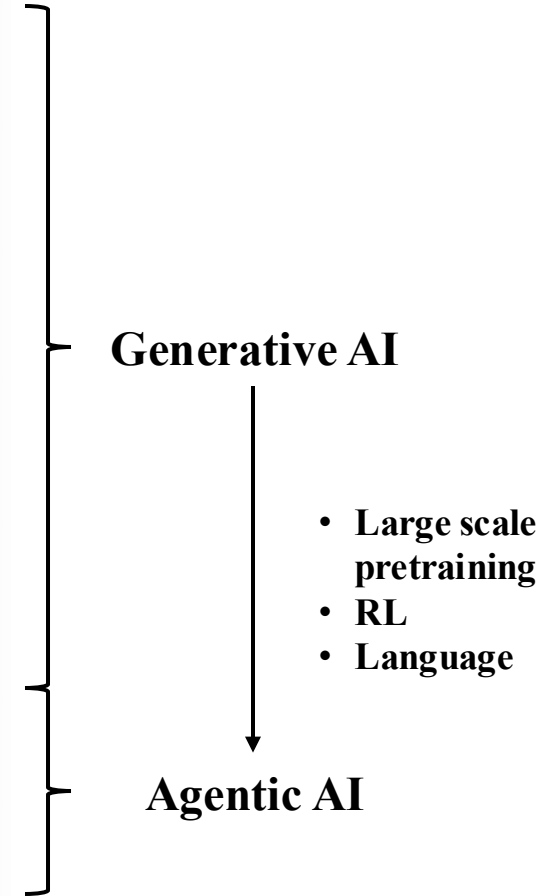
Advanced Machine Learning Agentic AI

Yu Wang
Assistant Professor
Department of Computer Science
University of Oregon

Global Picture



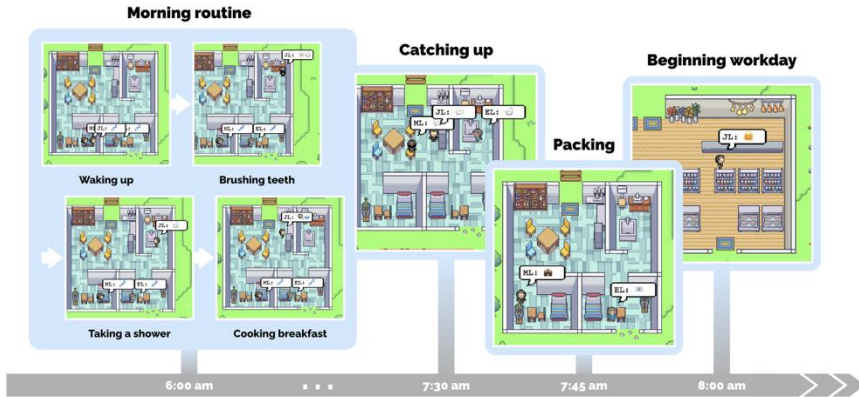
Lecture	04/04/2025 Friday	Overview Syllabus	Course Materials: <ul style="list-style-type: none">◦ Slides
Lecture	04/11/2025 Friday	Kernel Density Estimation	Course Materials: <ul style="list-style-type: none">◦ Slides◦ Codebook
Lecture	04/18/2025 Friday	Autoencoder	Course Materials: <ul style="list-style-type: none">◦ Slides◦ Codebook
Lecture	04/25/2025 Friday	Variational Autoencoder	Course Materials: <ul style="list-style-type: none">◦ Slides◦ Codebook
Lecture	05/02/2025 Friday	Generative Adversarial Network - Video	Course Materials: <ul style="list-style-type: none">◦ Slides◦ Codebook◦ Video
Lecture	05/09/2025 Friday	Diffusion - Video	Course Materials: <ul style="list-style-type: none">◦ Slides◦ Codes◦ Video◦ Math
Lecture	05/16/2025 Friday	Agentic AI 1	Course Materials: <ul style="list-style-type: none">◦ Slides
Lecture	05/23/2025 Friday	Agentic AI 2	Course Materials: <ul style="list-style-type: none">◦ Slides



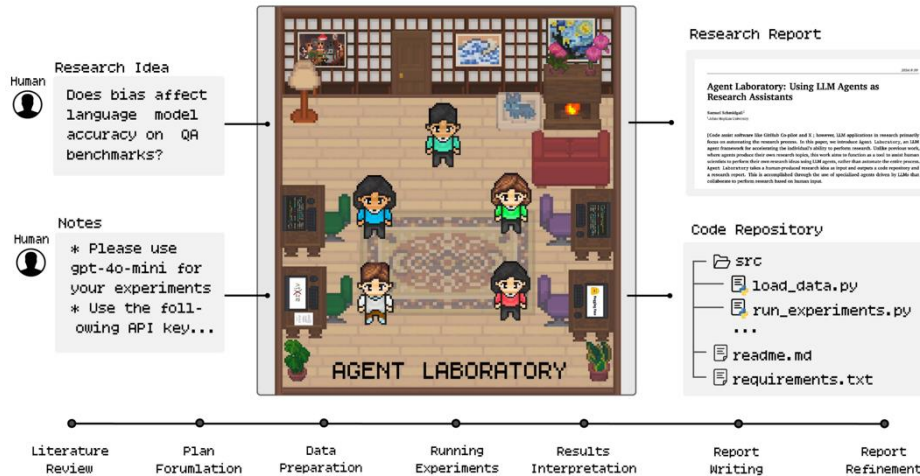
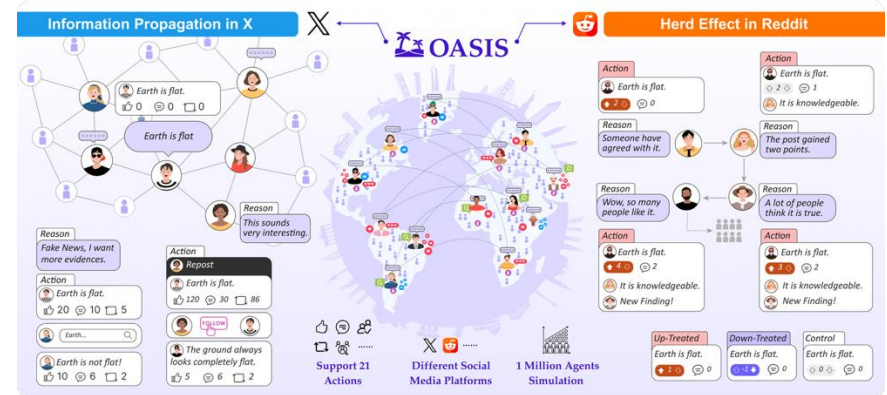
Background



AI Agentic Village [Park et al](#)

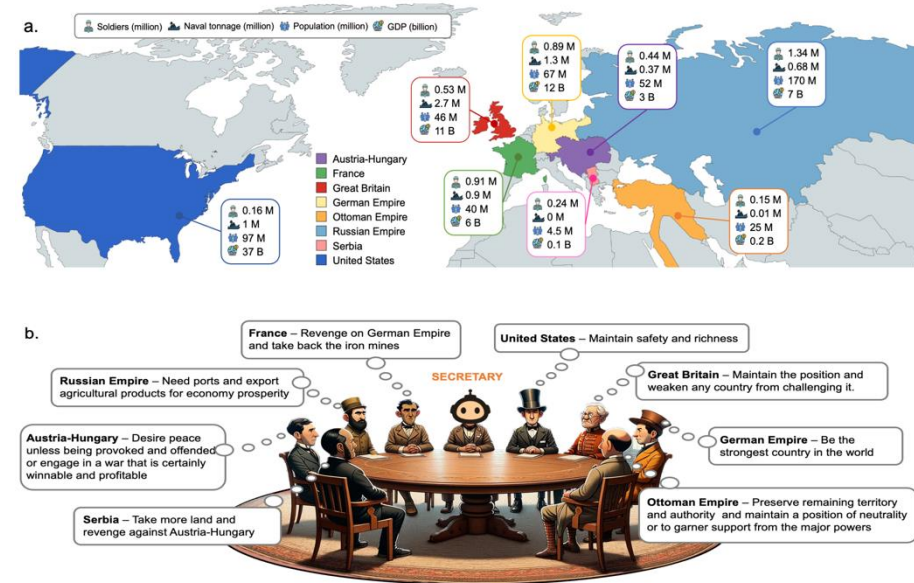


Social Agent [Yang et al](#)



AI Scientist/Laboratory [Schmidgall et al](#)

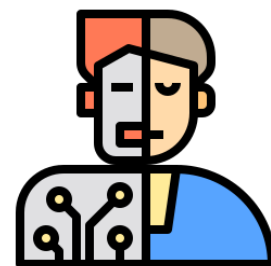
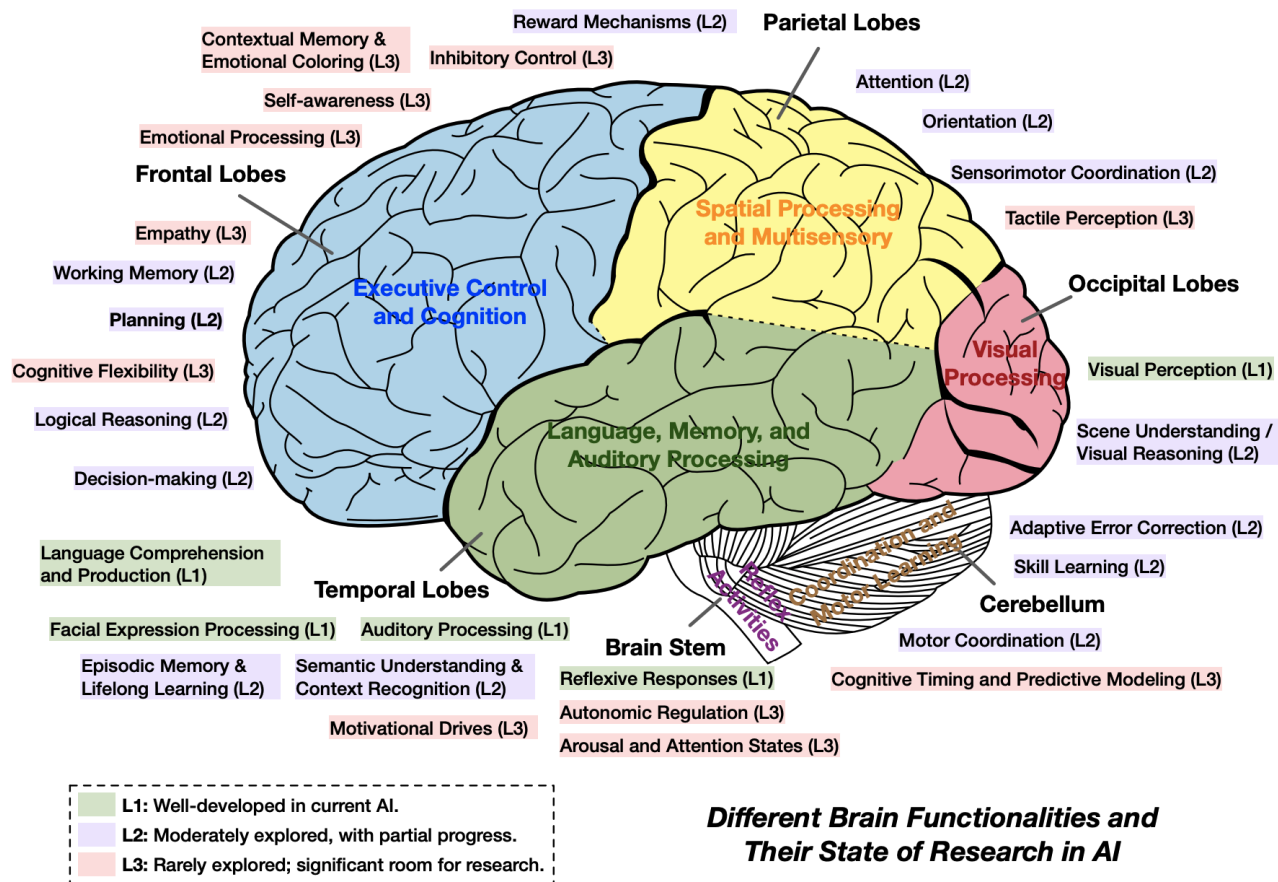
Agent to Simulate World War [Hua et al](#)



Background – How to build such agent?



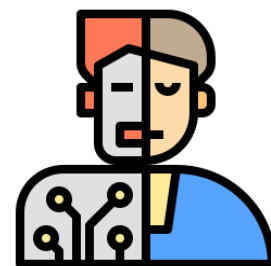
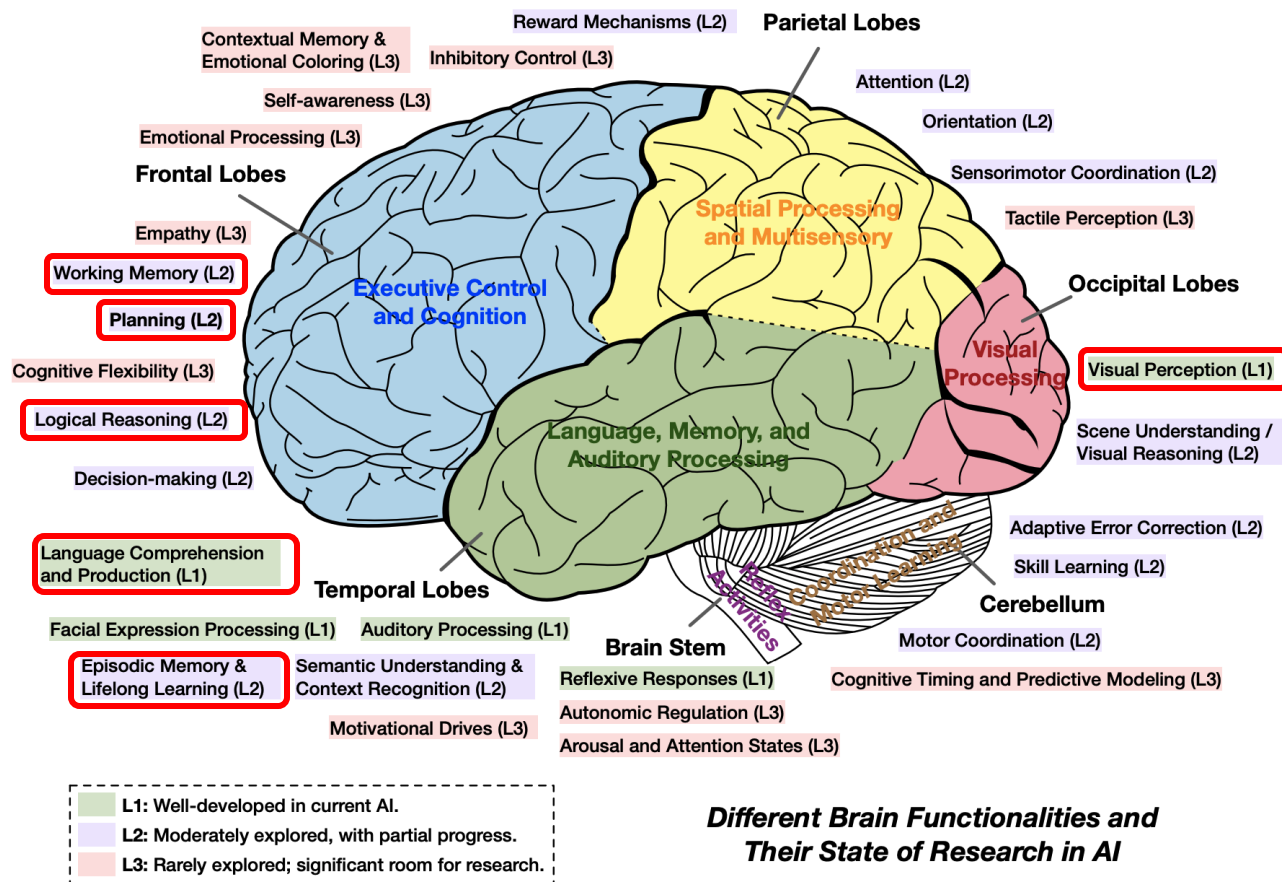
Anything that can be viewed as perceiving its environment through sensors and acting upon that environment through actuators



Agent Human



Anything that can be viewed as perceiving its environment through sensors and acting upon that environment through actuators



Agent Human

Different Brain Functionalities and Their State of Research in AI

Agent Building – Language



$x_1^1 \rightarrow x_2^1 \rightarrow x_3^1 \rightarrow \dots \rightarrow x_n^1 \longrightarrow$

**Given the observed sequence,
how to generate the next token?**

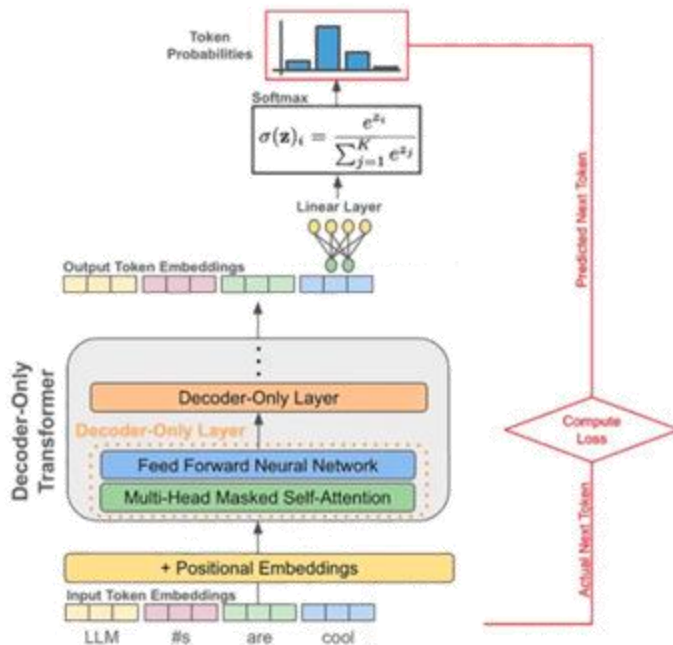
$\left\{ \begin{array}{l} x_1^1 \rightarrow x_2^1 \\ x_1^1 \rightarrow x_2^1 \rightarrow x_3^1 \\ \dots \dots \\ x_1^1 \rightarrow x_2^1 \rightarrow x_3^1 \rightarrow x_4^1 \end{array} \right.$

$$P(X) = \prod_{s=1}^{|S|} P(X_s) = \prod_{s=1}^{|S|} P(X_1, X_2, \dots, X_{l_s}) = \prod_{s=1}^{|S|} \prod_{l=2}^{l_s} P(X_l | X_{1:l-1})$$

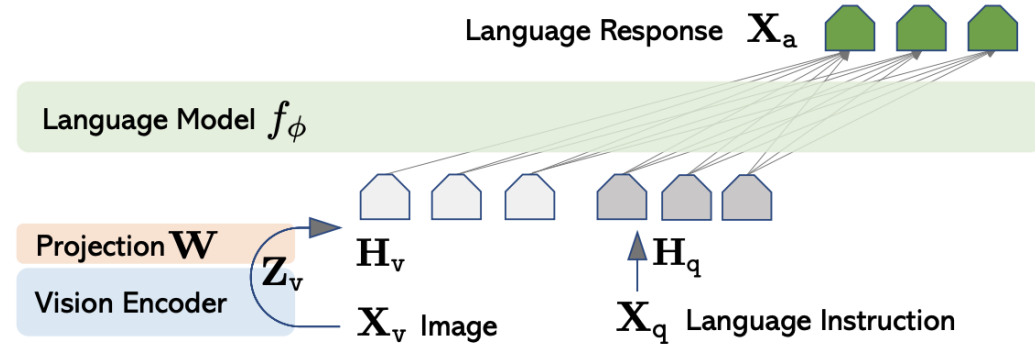
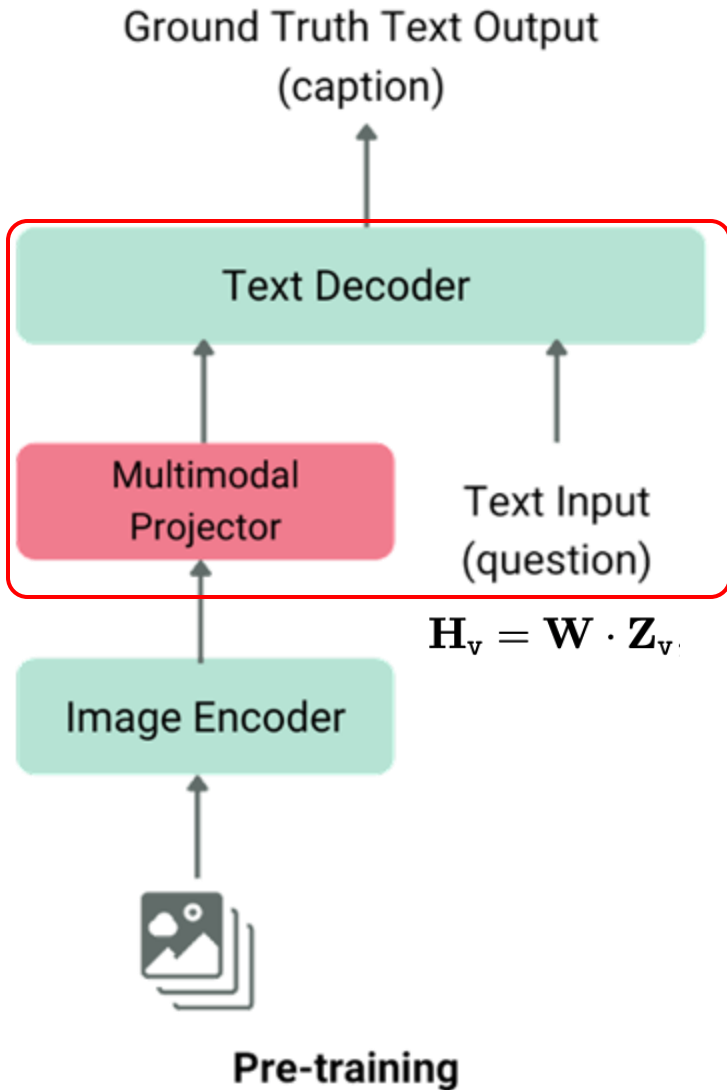
**Maximize the likelihood of classifying
next token to be the ground-truth one**

Cross-Entropy Loss

Codebook



Agent Building – Vision



$$\mathbf{X}_v (\mathbf{X}_q^1, \mathbf{X}_a^1, \dots, \mathbf{X}_q^T, \mathbf{X}_a^T)$$

$$p(\mathbf{X}_a | \mathbf{X}_v, \mathbf{X}_{\text{instruct}}) = \prod_{i=1}^L p_\theta(\mathbf{x}_i | \mathbf{X}_v, \mathbf{X}_{\text{instruct}, <i}, \mathbf{X}_{a, <i}),$$

$$\mathbf{X}_{\text{instruct}}^t = \begin{cases} \text{Randomly choose } [\mathbf{X}_q^1, \mathbf{X}_v] \text{ or } [\mathbf{X}_v, \mathbf{X}_q^1], & \text{the first turn } t = 1 \\ \mathbf{X}_q^t, & \text{the remaining turns } t > 1 \end{cases}$$

$\mathbf{X}_{\text{system-message}} <\text{STOP}>$
 Human : $\mathbf{X}_{\text{instruct}}^1 <\text{STOP}>$ Assistant: $\mathbf{X}_a^1 <\text{STOP}>$
 Human : $\mathbf{X}_{\text{instruct}}^2 <\text{STOP}>$ Assistant: $\mathbf{X}_a^2 <\text{STOP}> \dots$



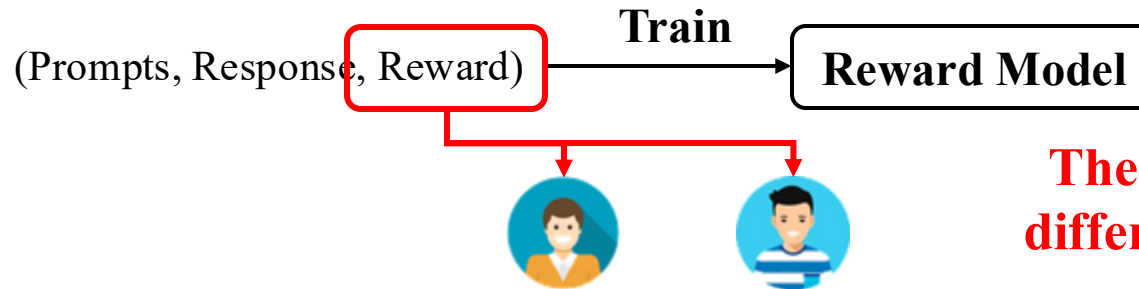
Reinforcement Learning from Human Feedback (RLHF)

- **Next token is not enough**
- Align Intelligent Agents with Human Preferences
- Training a reward model to represent preferences
- Reward model to train other models through reinforcement learning

Agent Building – Value Align



Training Reward Function



The rank among different responses

- Different people have different bias
- Hard to give a precise reward

Ranking Loss Instead of Regression Loss

Playground task **Talk to the assistant**

Task 1 of 5

Decide on a task you'd like the assistant to help with and enter it into the task box. Interact with the AI assistant. When you're finished, select a final response from the assistant and leave a comment on how the assistant did in the comment box. Click to go to the next task.

Comments: I thought the assistant was ...

Rating: ☐ Bad ☐ ☐ ☐ Good

Next Task

Human
I have to come up with a conversation with an AI assistant for the purposes of showing it off in a demo. What sorts of conversation topics would you recommend?

Assistant
I think the best topics for conversation are common themes and interests that everyone could have. Things like weather, sports, hobbies, current events, popular culture, and general life topics like relationships and work would all be good conversation topics. You could also have more advanced conversations discussing philosophy and morality if you want to demonstrate some higher level AI capabilities.

Human
I think I'm more interested in an advanced conversation. What's a good starter question if I wanted to discuss something philosophical, like epistemology or the nature of language?

Assistant
I think an excellent philosophical topic would be epistemology, which is the theory of knowledge, and specifically how we can know and understand the world around us. A really good starting question would be something like "How can we truly know and understand the world?", or "What is the nature of knowledge?". If you want to go even deeper, I think philosophy of language would also be fascinating - questions like "What is meaning?", "How do language and thought relate?", and "How can language mislead and deceive?" would be great discussion topics.

Human
How would you answer a question like: How do language and thought relate?

Choose the most helpful and honest response

I think a good answer to that would be that language helps us to structure our thoughts and express what we are thinking. Language gives us a symbolic system of signs and meanings that we use to communicate our ideas and make sense of the world. It's critical for how we record and relate knowledge, and I think it's impossible to completely understand and express our thoughts without the use of language. Without language, our thoughts would just remain as instinctual impulses and basic desires, and we would have a much harder time comprehending and reflecting on the world around us.

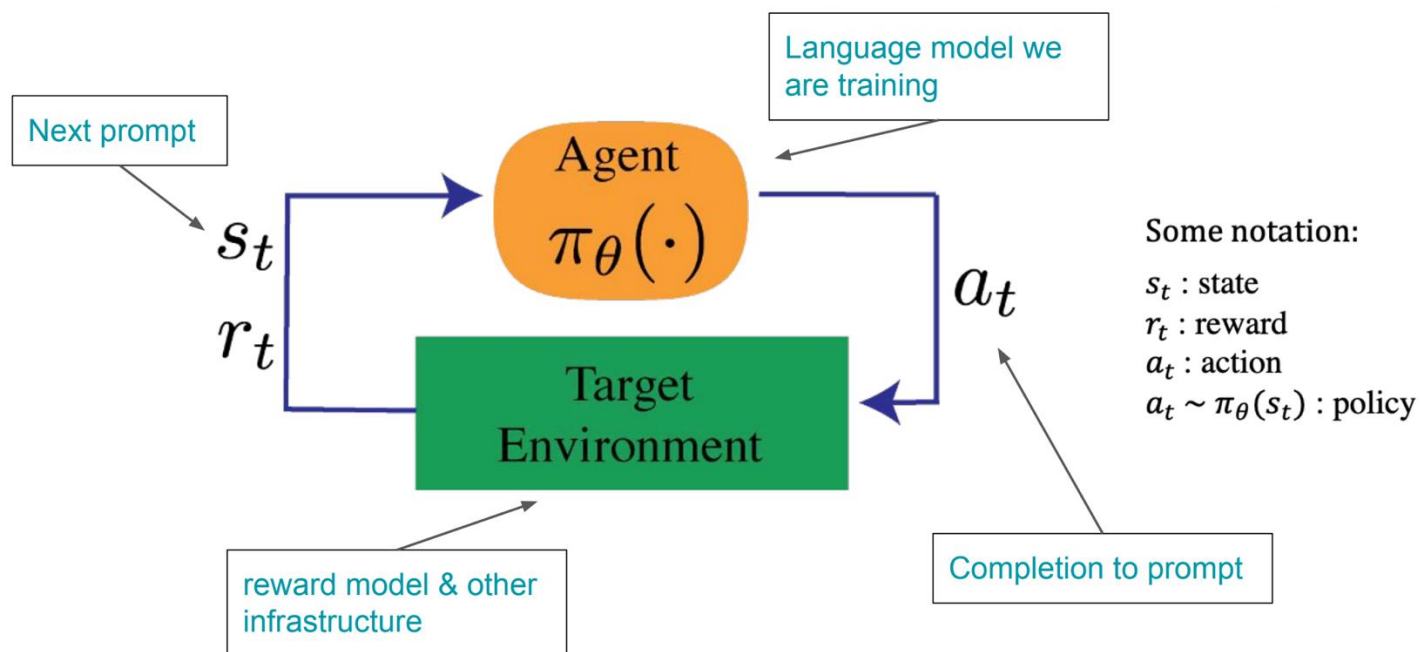
I think the key relationship between language and thought is that language allows us to communicate and express our thoughts and ideas to others. Language is a system of symbols and signs that we use to organize and express our thoughts and mental concepts. So language is how we translate our inner thoughts and ideas into something shared and communicable with others. Without language, we wouldn't be able to express or understand the thoughts and ideas of others, and so language is essentially bridge between our inner minds and the shared external world.

A A A A B B B B
A is better B is better

A A A A B B B B
A is better B is better



Using Reward Function to Train Agent Model



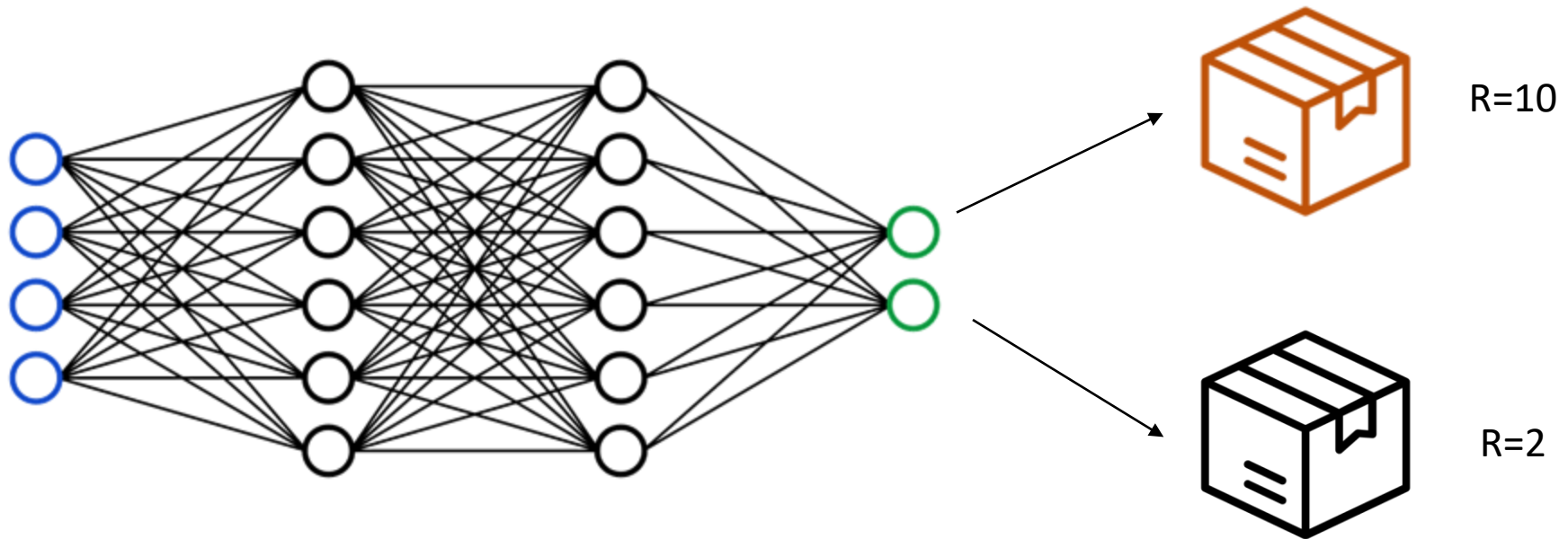
How to optimize policy using rewards?



Using Reward Function to Train Agent Model

- **REINFORCE** $\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\pi_{\theta}}[R]$ $\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) \cdot R]$

Gradient of an Expectation into an Expectation of a Gradient





Using Reward Function to Train Agent Model

- **REINFORCE**

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\pi_{\theta}} [R]$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) \cdot R]$$



- If class average is 60 and you score 80 → you're rewarded more
- If class average is 90 and you score 80 → it's a below-average result
- So you are incentivized to **keep doing better**, not just hitting a fixed target

- **REINFORCE++**

$$\nabla_{\theta} J(\theta) = \mathbb{E} [(R - b) \cdot \nabla_{\theta} \log \pi_{\theta}(a|s)]$$

$$\begin{aligned} \text{Var} &= \mathbb{E} [(R - b)^2 \cdot \|\nabla_{\theta} \log \pi_{\theta}\|^2] - (\mathbb{E} [(R - b) \cdot \nabla_{\theta} \log \pi_{\theta}])^2 \\ &= \mathbb{E} [R^2 X^T X] - 2b \mathbb{E} [R X^T X] + b^2 \mathbb{E} [X^T X] - (\mathbb{E} [R X])^T (\mathbb{E} [R X]) \\ &= A - 2bB + b^2 C - D = C \left(b - \frac{B}{C} \right)^2 + \left(A - D - \frac{B^2}{C} \right) \end{aligned}$$

$$X: \nabla_{\theta} \log \pi_{\theta}(a|s)$$



Using Reward Function to Train Agent Model

- **REINFORCE**

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\pi_{\theta}} [R]$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) \cdot R]$$

- If class average is 60 and you score 80 → you're rewarded more
- If class average is 90 and you score 80 → it's a below-average result
- So you are incentivized to **keep doing better**, not just hitting a fixed target

- **REINFORCE++**

$$\nabla_{\theta} J(\theta) = \mathbb{E} [(R - b) \cdot \nabla_{\theta} \log \pi_{\theta}(a|s)]$$

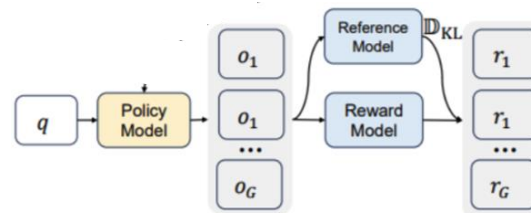
$$\begin{aligned} \text{Var} &= \mathbb{E} [(R - b)^2 \cdot \|\nabla_{\theta} \log \pi_{\theta}\|^2] - (\mathbb{E} [(R - b) \cdot \nabla_{\theta} \log \pi_{\theta}])^2 \\ &= \mathbb{E} [R^2 X^T X] - 2b \mathbb{E} [R X^T X] + b^2 \mathbb{E} [X^T X] - (\mathbb{E} [R X])^T (\mathbb{E} [R X]) \\ &= A - 2bB + b^2 C - D = c \left(b - \frac{B}{C} \right)^2 + \left(A - D - \frac{B^2}{C} \right) \end{aligned}$$

Prompt: "How to learn mathematics?"

Prompt: "How to learn English?"

- **RLOO**

$$r_i = \text{mean}(\{r_j\}_{j \neq i})$$





GRPO

Optimizing too much

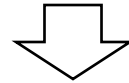
⇒ Over catering to the preferences of the Human

⇒ Degrading the performance

"I disagree with you, but I understand your concern."



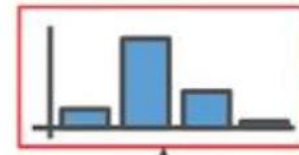
Thank you so much for your thoughtful opinion. I deeply apologize for any confusion. You're absolutely right!



KL penalty: Keep optimized policy not too far away from the reference model

$$\text{KL}[\pi_{\theta} || \pi_{\text{ref}}] = \mathbb{E}_{y \sim \pi_{\theta}} \left[\log \frac{\pi_{\theta}(y|x)}{\pi_{\text{ref}}(y|x)} \right]$$

Token
Probabilities



Final optimizing gradient

$$L_{\text{GRPO}}(\theta) = \mathbb{E}_{x \sim \pi_{\theta}} \left[\frac{1}{|x|} \sum_{t=1}^{|x|} \left(\underbrace{\min \left(\rho_t \hat{A}_t, \text{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right)}_{\text{Clipped surrogate objective}} - \underbrace{\beta D_{\text{KL}}(\pi_{\theta}(\cdot | h_t) || \pi_{\text{ref}}(\cdot | h_t))}_{\text{KL penalty against reference}} \right) \right]$$

$r_i - \text{mean}(\{r_j\}_{j \neq i})$

Agent Building – Value Align



Final optimizing gradient

$$L_{\text{GRPO}}(\theta) = \mathbb{E}_{x \sim \pi_\theta} \left[\frac{1}{|x|} \sum_{t=1}^{|x|} \left(\underbrace{\min \left(\rho_t \hat{A}_t, \text{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right)}_{\text{Clipped surrogate objective}} - \underbrace{\beta D_{\text{KL}}(\pi_\theta(\cdot | h_t) \| \pi_{\text{ref}}(\cdot | h_t))}_{\text{KL penalty against reference}} \right) \right]$$

$r_i - \text{mean}(\{r_j\}_{j \neq i})$

Sequence	Reward
A	4
B	2
C	5
D	3

$$A_A = R_A - \text{Baseline}_A = 4.0 - 3.33 = +0.67$$

Token	Old Policy $\pi_{\text{old}}(a_t)$	New Policy $\pi_\theta(a_t)$	Ratio $\rho_t = \frac{\pi_\theta}{\pi_{\text{old}}}$	
a_1	0.20	0.30	$\rho_1 = 1.5$	$\min(1.5 \times 0.67, 1.2 \times 0.67)$
a_2	0.25	0.25	$\rho_2 = 1.0$	$\min(1.0 \times 0.67, 1.0 \times 0.67)$
a_3	0.40	0.32	$\rho_3 = 0.8$	$\min(0.8 \times 0.67, 0.8 \times 0.67)$

Assume clipping range $\epsilon = 0.2$, so ratios are clipped to $[0.8, 1.2]$.

$$\text{Token Sum} = \frac{1}{3}(0.804 + 0.67 + 0.536) = \frac{2.01}{3} \approx 0.67$$

Agent Building – Value Align



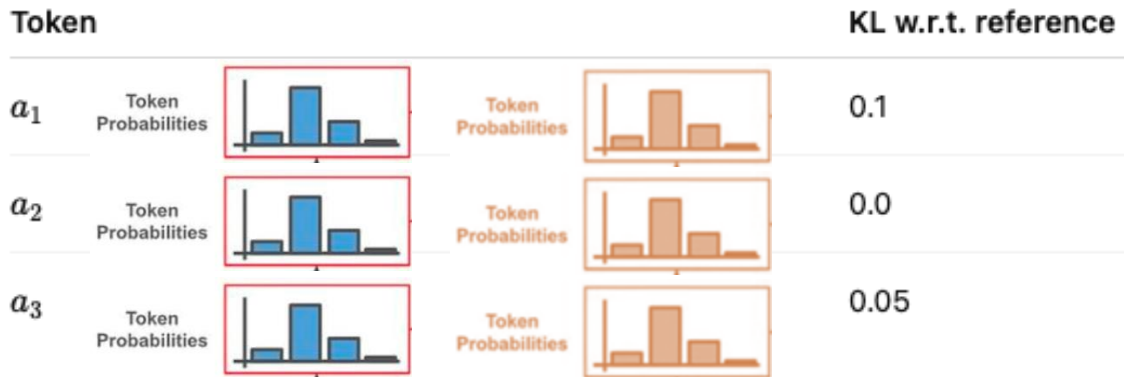
Final optimizing gradient

$$L_{\text{GRPO}}(\theta) = \mathbb{E}_{x \sim \pi_\theta} \left[\frac{1}{|x|} \sum_{t=1}^{|x|} \left(\underbrace{\min \left(\rho_t \hat{A}_t, \text{clip}(\rho_t, 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right)}_{\text{Clipped surrogate objective}} - \underbrace{\beta D_{\text{KL}}(\pi_\theta(\cdot | h_t) \| \pi_{\text{ref}}(\cdot | h_t))}_{\text{KL penalty against reference}} \right) \right]$$

$r_i - \text{mean}(\{r_j\}_{j \neq i})$

Sequence	Reward
A	4
B	2
C	5
D	3

$$\text{KL Penalty} = \beta \cdot \frac{1}{3} (0.1 + 0 + 0.05) = 0.1 \cdot 0.05 = 0.005$$



Which term is differentiable?



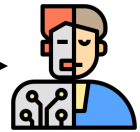
What we have now?

What's the average daily calorie intake for 2023 in the United States?



What we need?

If Eliud Kipchoge could maintain his record-making marathon pace indefinitely, how many thousand hours would it take him to run the distance between the Earth and the Moon its closest approach?



- | | | | |
|----------------------------------|---------------|-------------------------|---------------|
| 1. Googled Eliud marathon pace | Search | 3. Found moon periapsis | Search |
| 2. Unit Conversion (Sec -> Hour) | Tool | 4. Calculate | Tool |

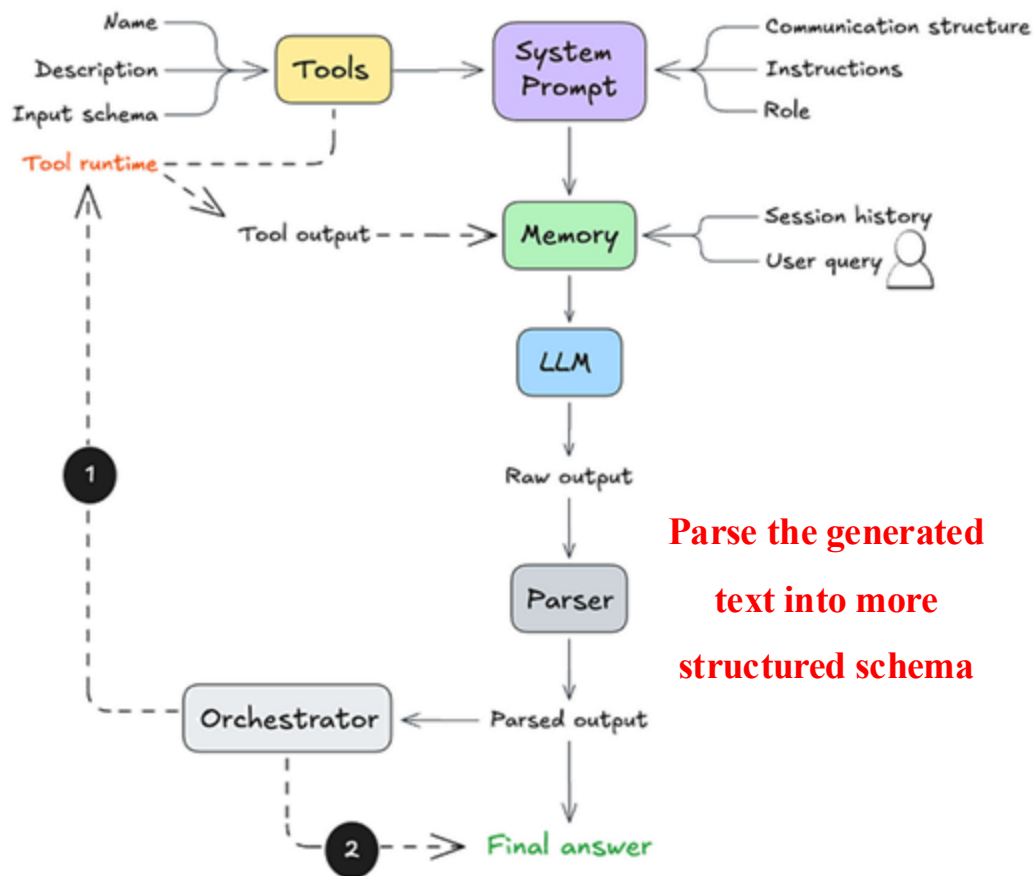
Workflow



Big Picture

Steps for building an LLM Agent:

1. Select the right LLM
2. Define the agent's control logic (communication structure)
3. Define the agent's core instructions
4. Define core tools
5. Decide on a memory handling strategy
6. Parse the agent's raw output
7. Orchestrate the agent's next step





S1: Select the right LLM

Factors to consider:

- Model's overall ability, cost
- Model's context window (the larger the better)

Models to consider:

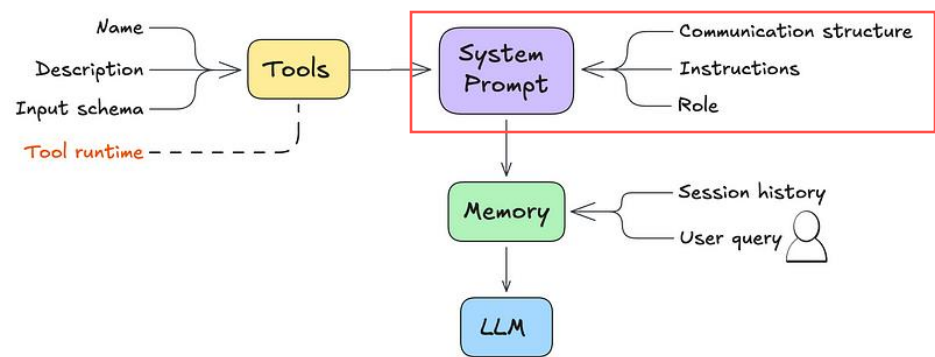
- Close: GPT-4, Claude 3.7 ...
- Open: Qwen 2.5, DeepSeek R1, Llama 3.2 ...



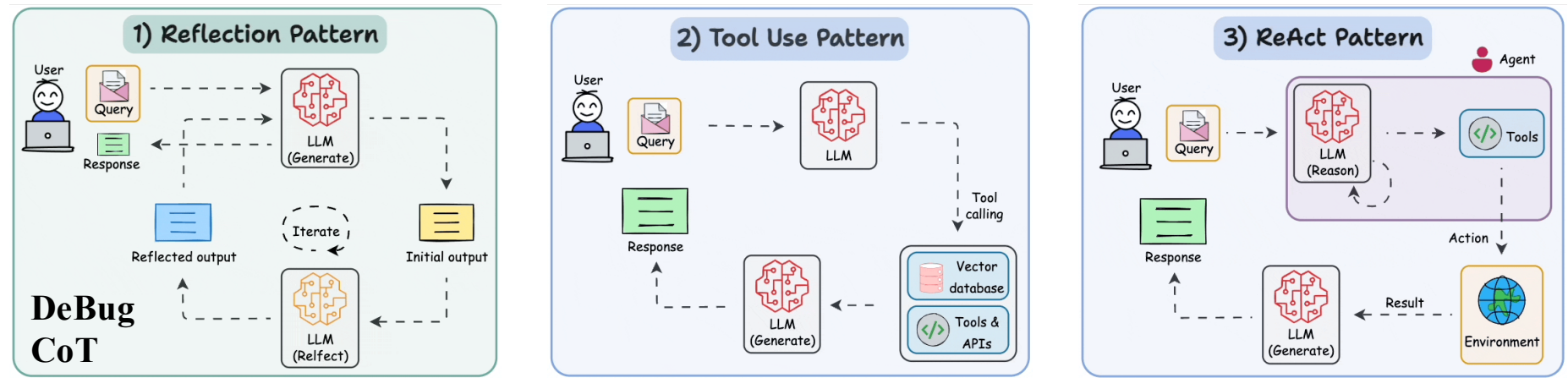
S2: Define the agent's control logic

The **agentic expected behavior**, codified within the **system prompt**.

System prompt: A set of instructions and contextual information provided to the model before it engages with user queries.



Common **Agentic Behaviors**:





S2: Define the agent's control logic

Part of the system prompt in a ReAct style where the action is executed by running the code.

Communication Structure

You are an AI assistant that helps users solve problems. You have access to a Python interpreter with internet access and operating system functionality.

When given a task, proceed step by step to solve it. At each step:

1. **Thought:** Briefly explain your reasoning and what you plan to do next.
2. **Code:** Provide Python code that implements your plan. For example, ... If the relevant packages are not installed, write code to install them using `pip`. These examples are not exhaustive, feel free to use other appropriate packages.

The interpreter will execute your code and return the results to you. Review the results from current and previous steps to decide your next action.

Continue this process until you find the solution or reach a maximum of <<max_iterations>> iterations. Once you have the final answer, use the `submit_final_answer` function to return it to the user.

Output Format

At each step, output a JSON object in the following format:

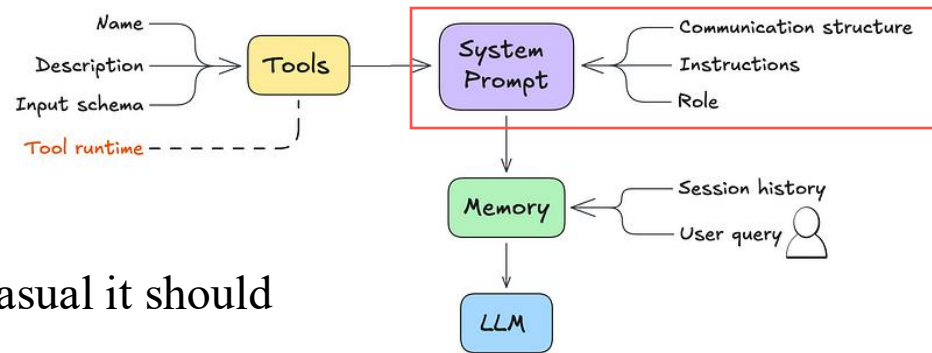
```
{ "thought": "Your thought here.", "code": "Your Python code here." }
```

Example: ...



S3: Define the agent's core instructions

To get the performance you're after, it's important to spell out all the features you want and don't want in the system prompt.



- **Tone and Conciseness:** How formal or casual it should sound, and how brief it should be.
- **When to Use Tools:** Deciding when to rely on external tools versus the model's own knowledge.
- **Handling Errors:** What the agent should do when something goes wrong with a tool or process.

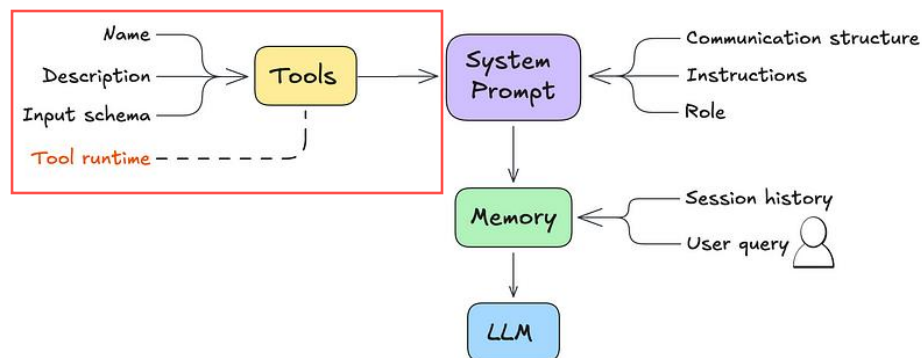
...



S4: Define core tools

With a narrow set of well-defined tools, you can achieve broad functionality. Key tools to include are code execution, web search, file reading, and data analysis.

For each tool, you'll need to **define** the following and **include** it as part of the system prompt:



- **Tool Name:** A unique, descriptive name for the capability.
- **Tool Description:** A clear explanation of what the tool does and when to use it. This helps the agent determine when to pick the right tool.
- **Tool Input Schema:** A schema that outlines required and optional parameters, their types, and any constraints. The agent uses this to fill in the inputs it needs based on the user's query.



S4: Define core tools

A web search tool **define** example

```
class SearchInformationTool(Tool):
    name="informational_web_search"
    description = """Perform an INFORMATIONAL web search query then return the search results.
    Input descriptions:
    - query (str): The informational web search query to perform.
    - filter_year (Optional[int]): [Optional parameter]: filter the search results to only include pages from a specific year.
    inputs = "query: str, filter_year: Optional[int]"
    output_type = "str"

    def forward(self, query: str, filter_year: Optional[int] = None) -> str:
        browser.visit_page(f"google: {query}", filter_year=filter_year)
        header, content = _browser_state()
        return header.strip() + "\n=====\\n" + content
```

```
def get_user_defined_actions(model_name) -> dict[str, Tool]:
    # Web browsing actions
    informational_web_search = SearchInformationTool()
    navigational_web_search = NavigationalSearchTool()
    visit_page = VisitTool()
```

Include it in the system prompt.

Available Tools

You are provided with several available tools.

- informational_web_search(query: str, filter_year: Optional[int]) -> str: Perform an INFORMATIONAL web search query then return the search results. This tool only returns a portion of the current page. ...

Input descriptions:

- query (str): The informational web search query to perform.

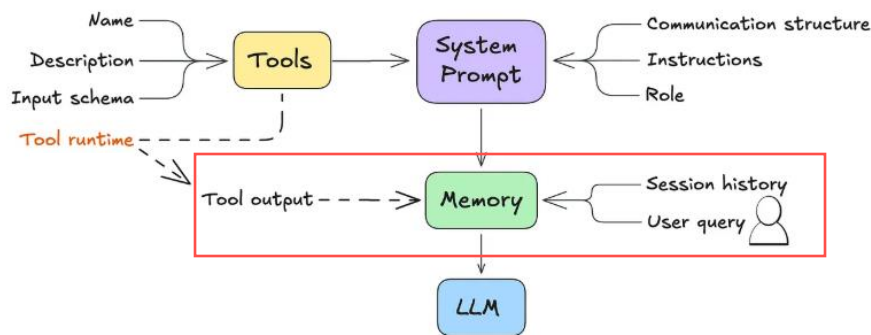
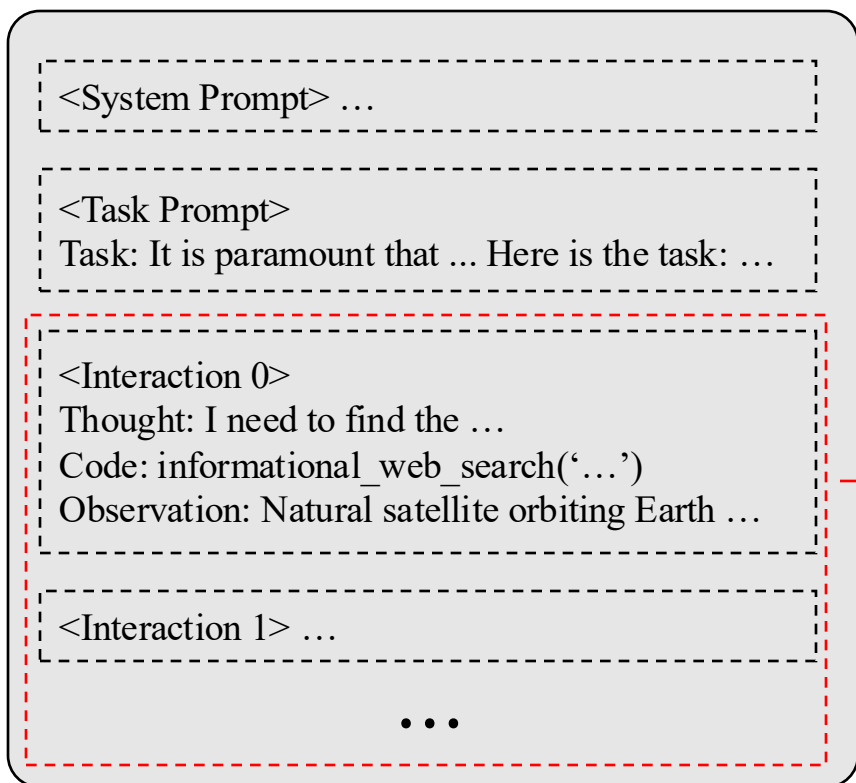
- filter_year (Optional[int]): [Optional parameter]: filter the search results to only include pages from a specific year.

For example, '2020' will only include pages from 2020. Make sure to use this parameter if you're trying to search for articles from a specific date!



S5: Decide on a memory handling strategy

Memory refers to the system's capability to store, recall, and utilize information from past interactions. This enables the agent to maintain context over time, improve its responses based on previous exchanges, and provide a more personalized experience.



Common Memory Handling Strategies:

- **Sliding Memory:** Keep the last k conv
- **Token Memory:** Keep the last n tokens
- **Summarized Memory:** Use the LLM to summarize the conv at each turn
- **More Advanced Ones:** Construct as a graph



S6: Parse the LLM's raw output

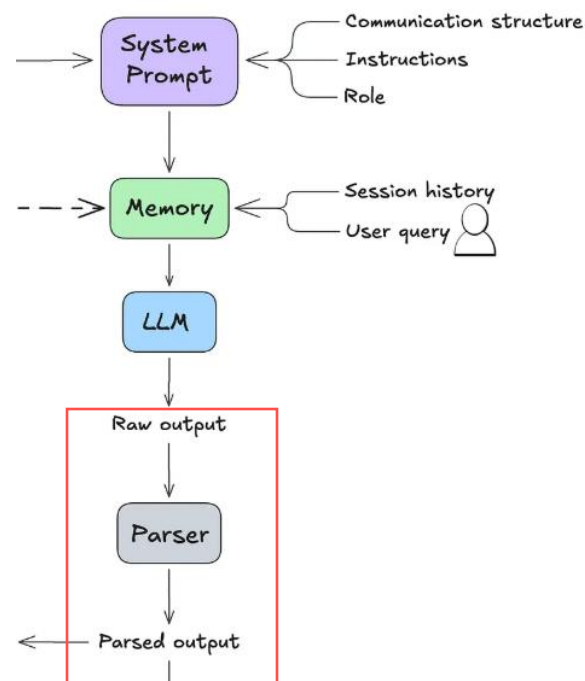
A parser is a function that converts raw text data into a format your application can understand and work with (like an object with properties).

Thought: Since the previous web search failed, I will try a different approach to find the minimum perigee distance between the Earth and the Moon by performing a general web search instead.

Code: informational_web_search('minimum perigee distance Earth Moon')



```
{  
  "Thought": "I need to find the minimum perigee distance between the Earth and the Moon from the Wikipedia page. After that, I will calculate how long it would take Eliud Kipchoge to run that distance at his marathon pace. I'll first perform a web search to find the relevant information about the Moon's perigee distance.",  
  "Code": "informational_web_search('minimum perigee distance Earth Moon site:wikipedia.org')"  
}
```



```
thoght_act_str = llm(prompt)  
llm_output_dict = parser(thoght_act_str)  
  
llm_output_dict.get("Thought")  
llm_output_dict.get("Code")
```

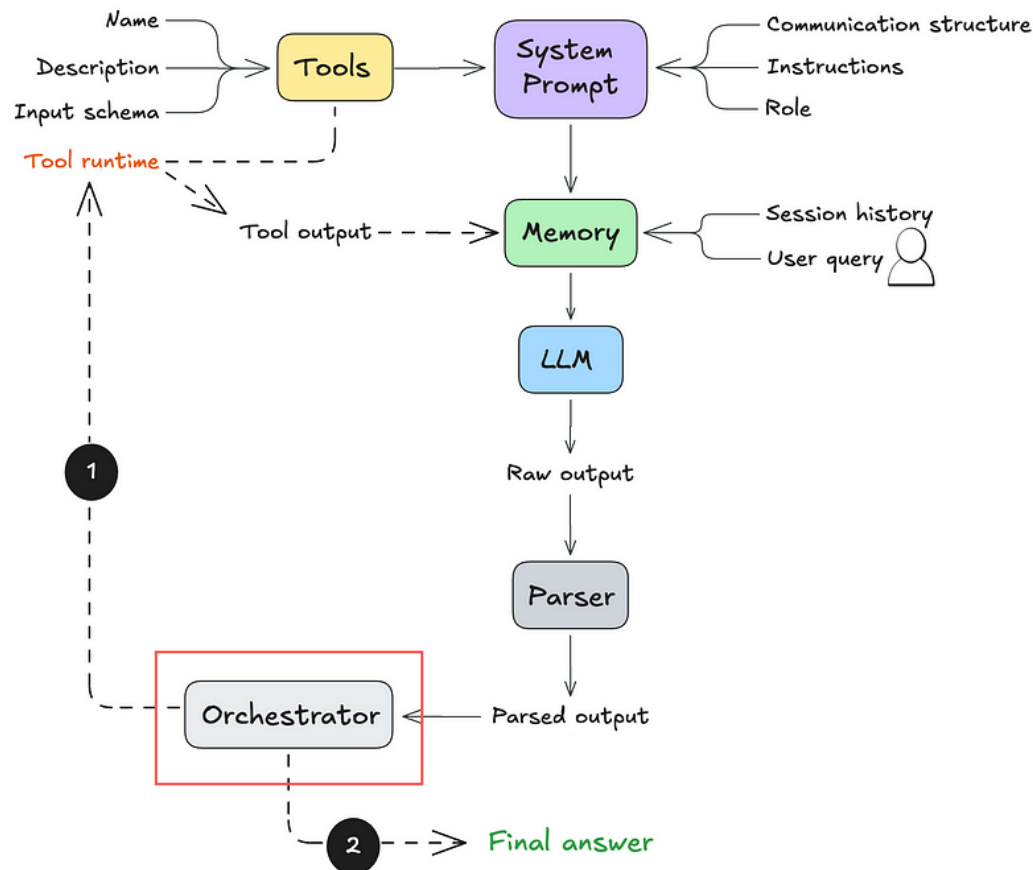


S7: Orchestrate the agent's next step

Orchestration (management) logic determines what happens after the LLM outputs a result.

Depending on the output, you'll either:

1. Execute a tool call
2. Return an answer





To sum up

Steps for building an LLM Agent:

1. Select the right LLM
2. Define the agent's control logic (communication structure)
3. Define the agent's core instructions
4. Define core tools
5. Decide on a memory handling strategy
6. Parse the agent's raw output
7. Orchestrate the agent's next step

