

Разработка в NLP

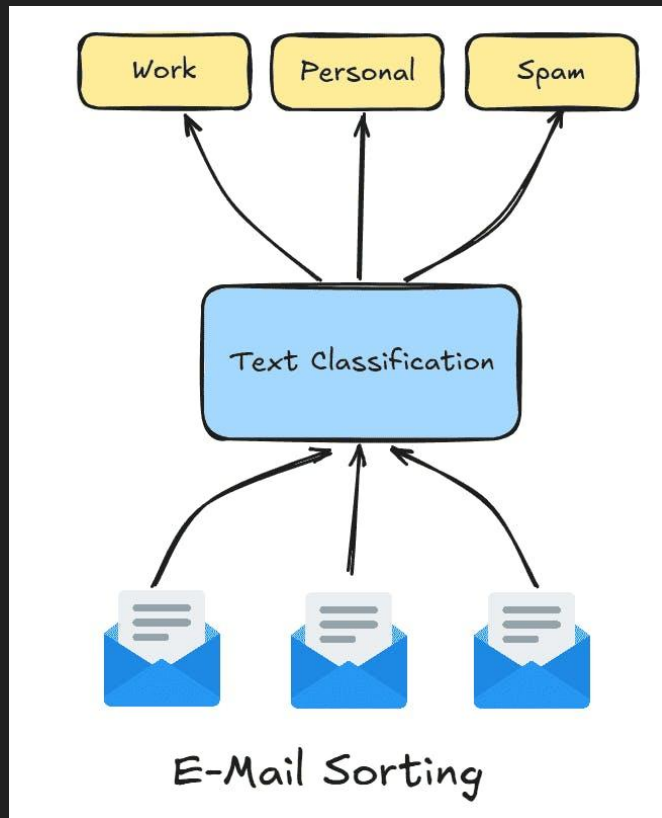
langchain/langgraph

Задачи: NER

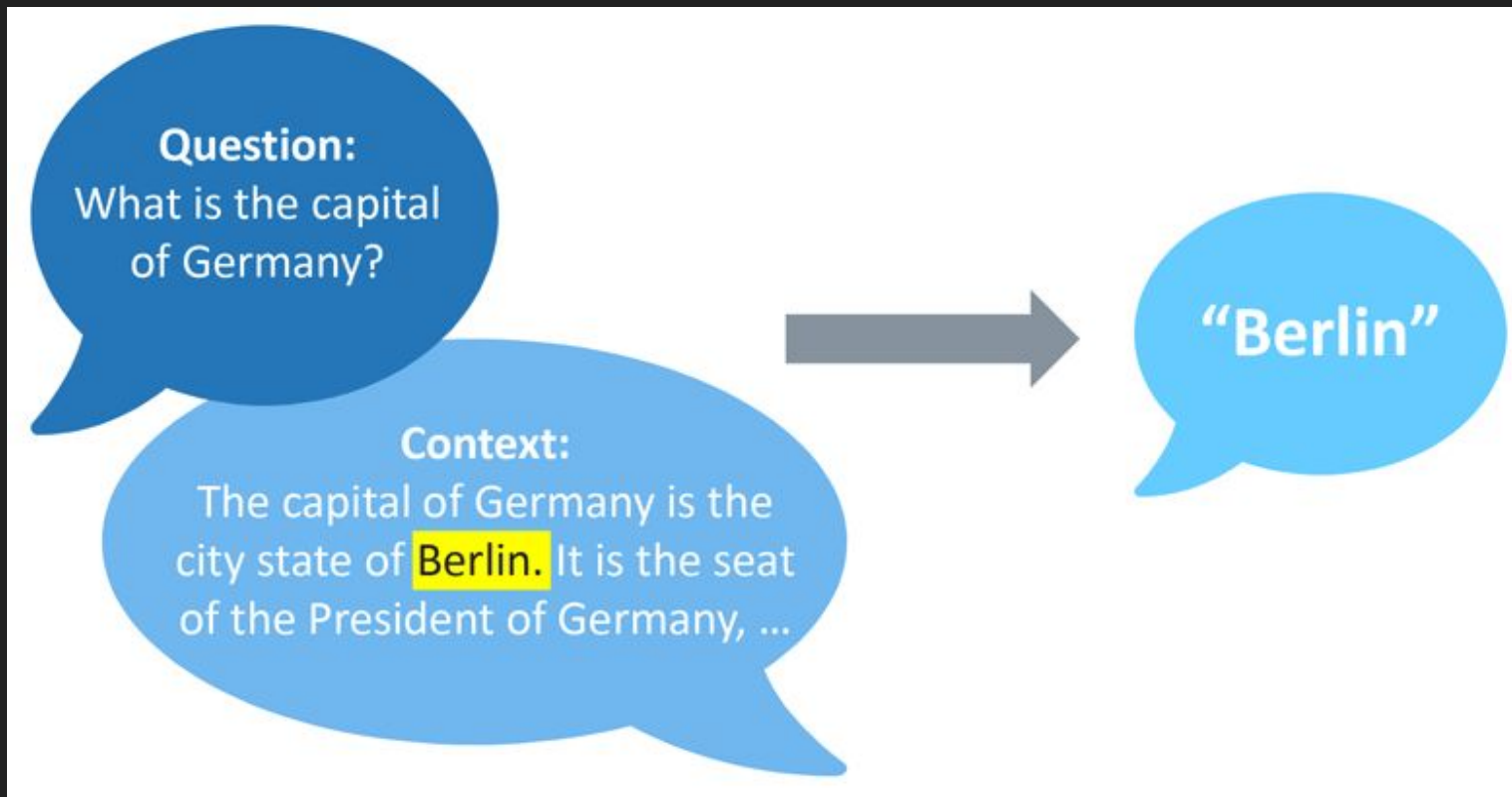
In fact, the Chinese NORP market has the three CARDINAL most influential names of the retail and tech space – Alibaba GPE , Baidu ORG , and Tencent PERSON (collectively touted as BAT ORG), and is betting big in the global AI GPE in retail industry space . The three CARDINAL giants which are claimed to have a cut-throat competition with the U.S. GPE (in terms of resources and capital) are positioning themselves to become the ‘future AI PERSON platforms’. The trio is also expanding in other Asian NORP countries and investing heavily in the U.S. GPE based AI GPE startups to leverage the power of AI GPE . Backed by such powerful initiatives and presence of these conglomerates, the market in APAC AI is forecast to be the fastest-growing one CARDINAL , with an anticipated CAGR PERSON of 45% PERCENT over 2018 - 2024 DATE .

To further elaborate on the geographical trends, North America LOC has procured more than 50% PERCENT of the global share in 2017 DATE and has been leading the regional landscape of AI GPE in the retail market. The U.S. GPE has a significant credit in the regional trends with over 65% PERCENT of investments (including M&As, private equity, and venture capital) in artificial intelligence technology. Additionally, the region is a huge hub for startups in tandem with the presence of tech titans, such as Google ORG , IBM ORG , and Microsoft ORG .

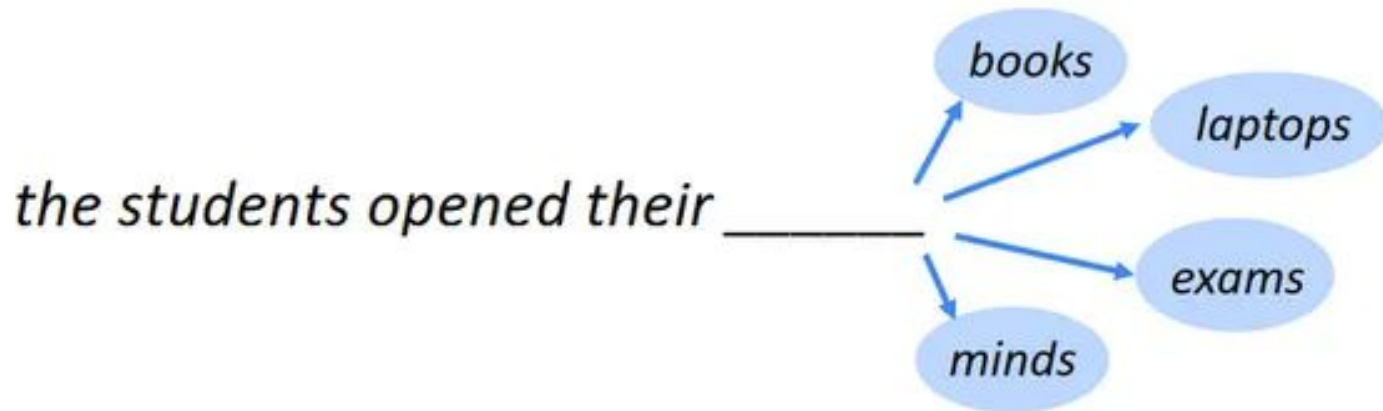
Задачи: классификация текста



Задачи: QA



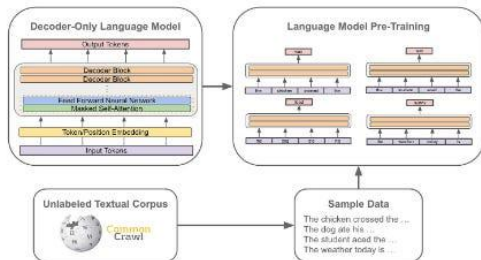
Задачи: авторегрессия



LLMs

Alignment

Pre-Training



SFT

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior.

This data is used to fine-tune GPT-3 with supervised learning.



RLHF

A prompt and several model outputs are sampled.

A labeler ranks the outputs from best to worst.

This data is used to train our reward model.



A new prompt is sampled from the dataset.

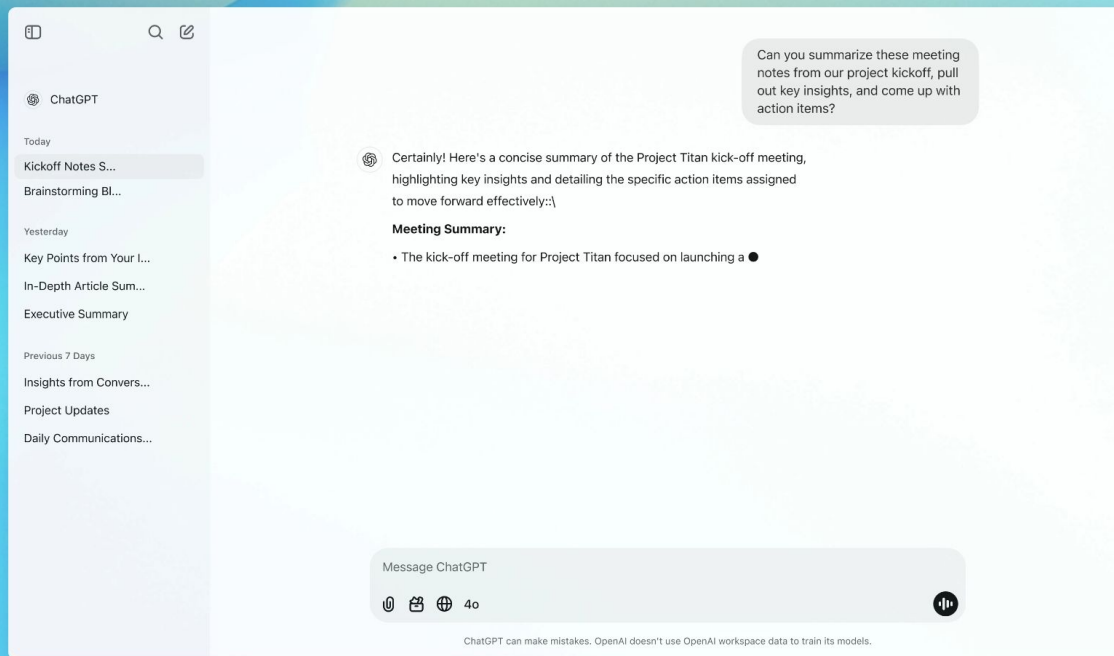
The policy generates an output.

The reward model calculates a reward for the output.

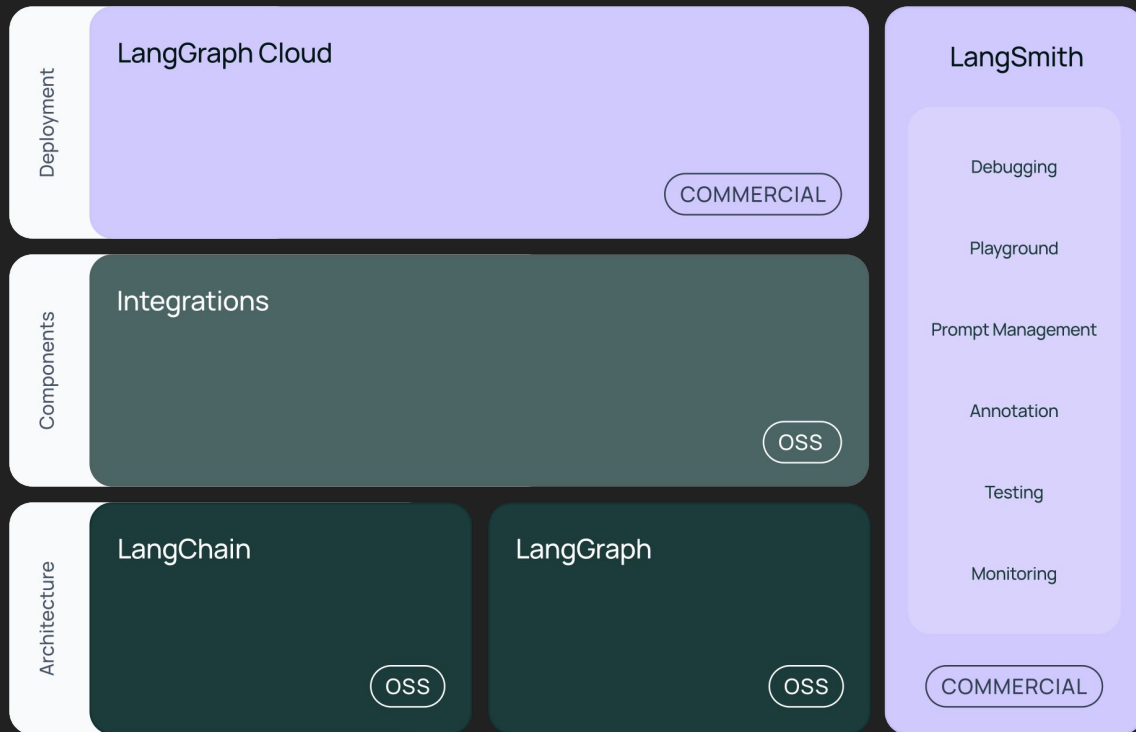
The reward is used to update the policy using PPO.



Chat models



Langchain

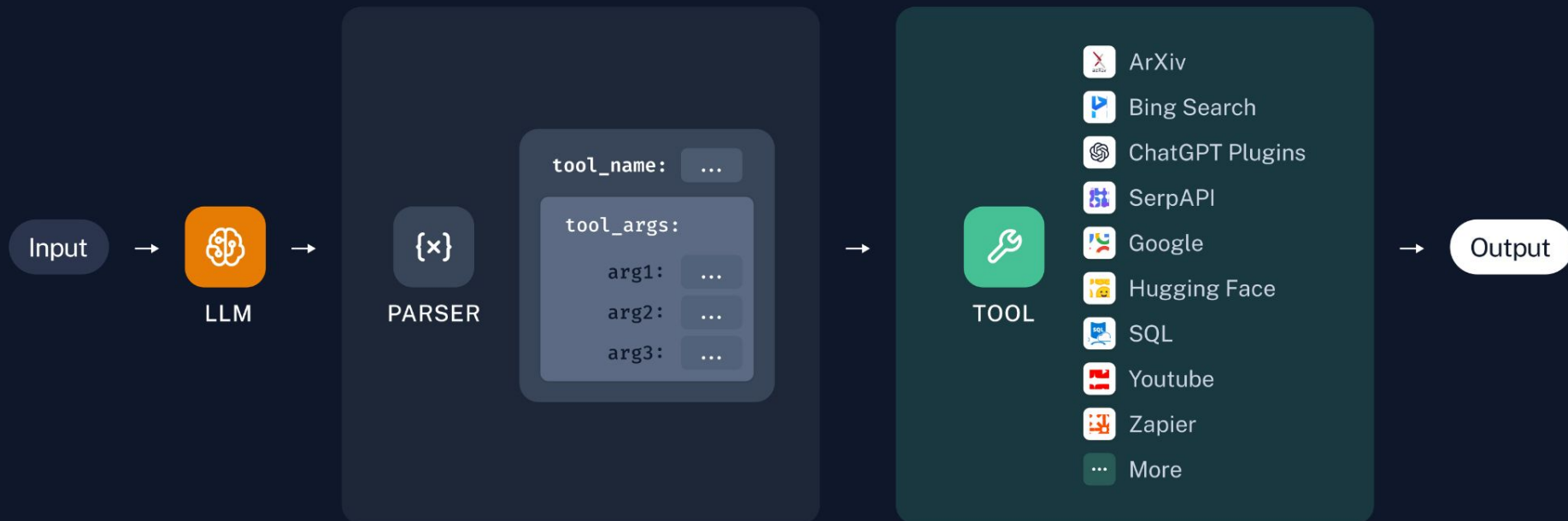


Langsmith

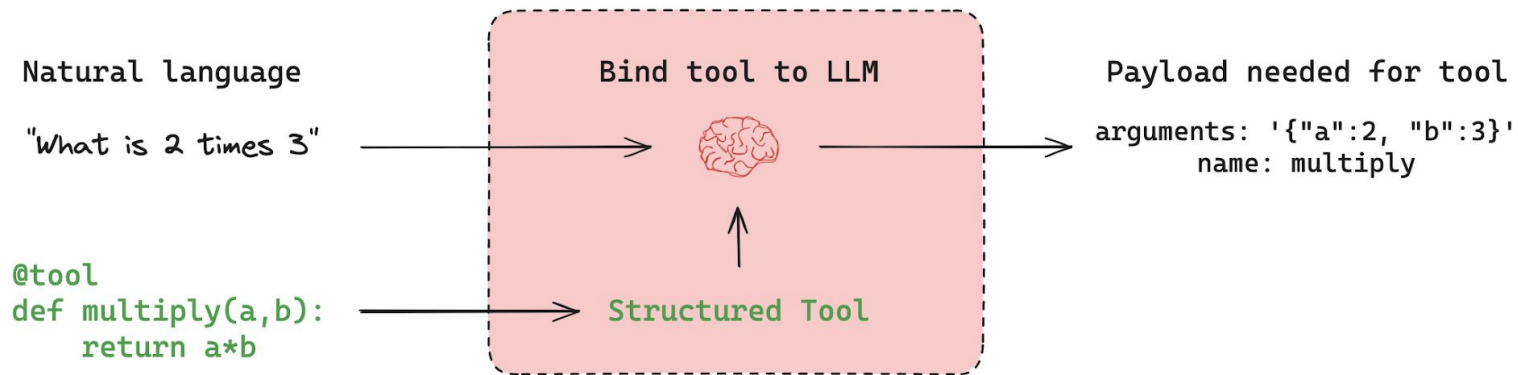
The screenshot displays the LangSmith web interface for a project named "chat-langchain". The interface is divided into several sections:

- Left Sidebar:** Contains navigation links for Projects (360), Annotation Queues (17), Deployments, Datasets & Testing (468), and Hub (36). At the bottom, there are links for API Keys and Documentation, and contact information for LangChain Inc.
- Header:** Shows the project name "chat-langchain" and a link to the project page.
- Traces Section:** A table of LLM runs with columns for Name, Input, Start Time, Latency, Dataset, Annotation Queue, Tokens, and Cost. The table is filtered for the "Last 7 days". A modal is open for the selected run, showing details for "Chat LangChain Review".
- Details Panel:** Located on the right, it provides summary statistics for the project:
 - Run Count:** 26,180
 - Total Tokens:** 199,037,036 / \$93,143,7495
 - Median Tokens:** 5,832
 - Error Rate:** 1%
 - % Streaming:** 99%
 - Latency:** P50: 5.30s, P99: 68.73s
 - First Token:** P50: 1.50s, P99: 10.63s
 - Feedback:** A list of feedback items including ALWAYS_PASS, CORRECTNESS, FAITHFULNESS_SCORE, GRAMMATIC_CORRECTNESS-OPTIM, JRP, WOW WHAT DUDE, USER_SCORE, and USER_CLICK.

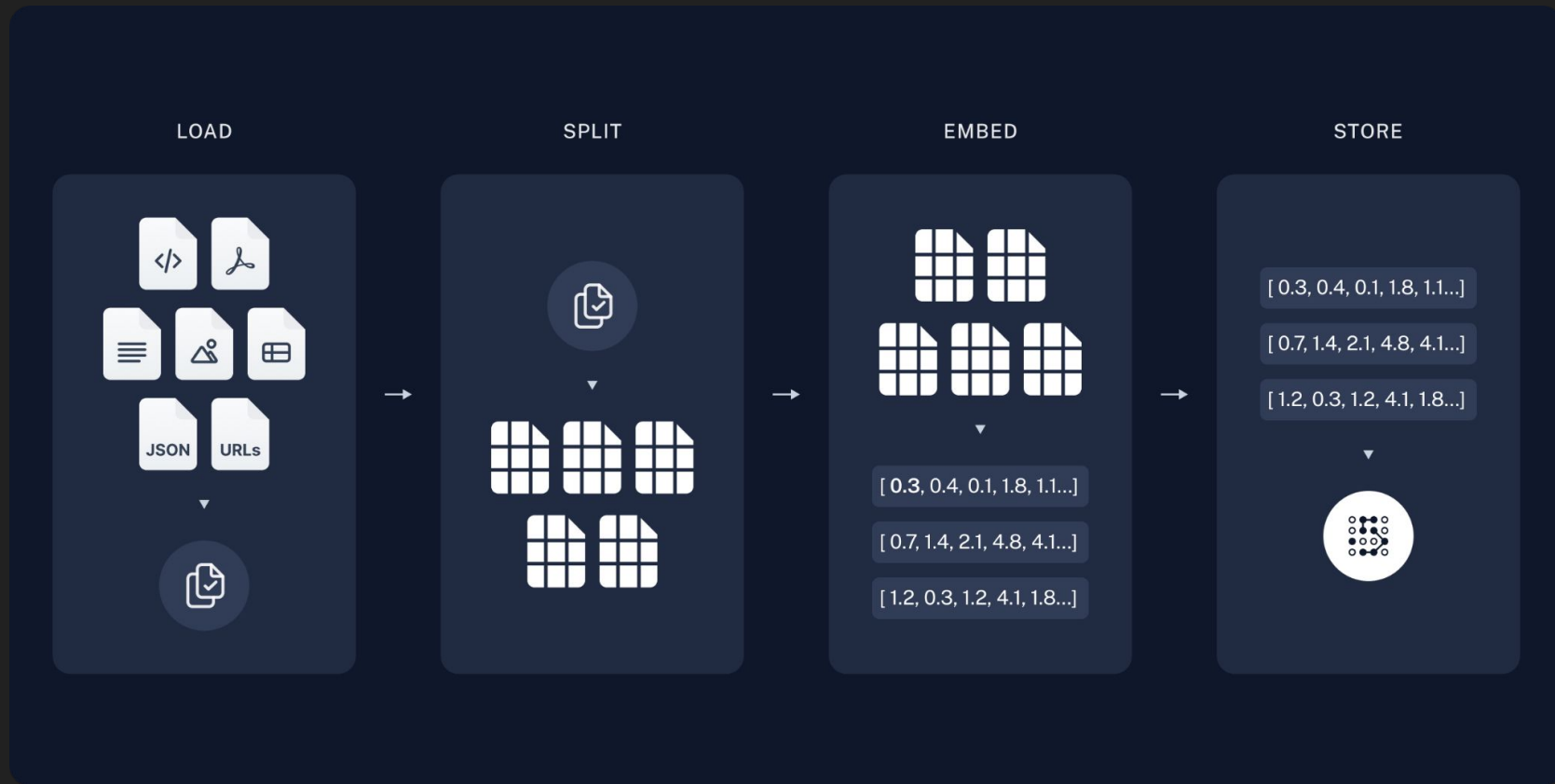
Tools



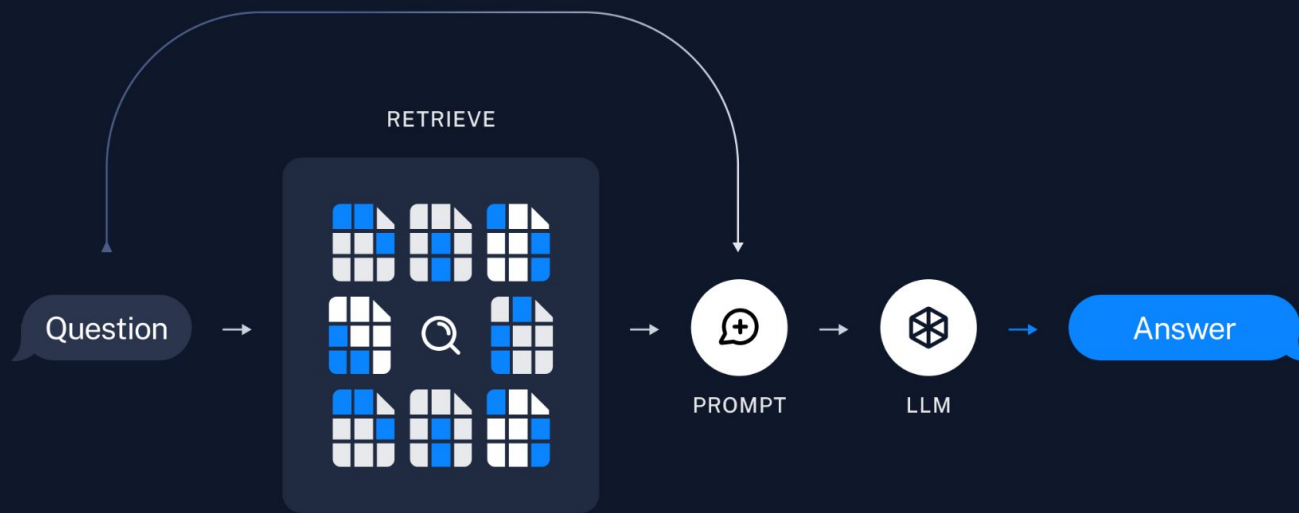
Tools



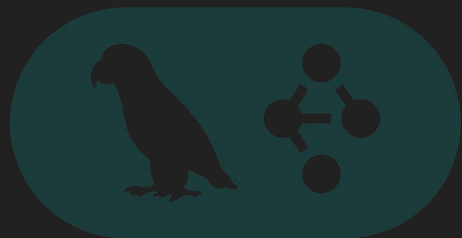
RAG: indexing



RAG: pipeline



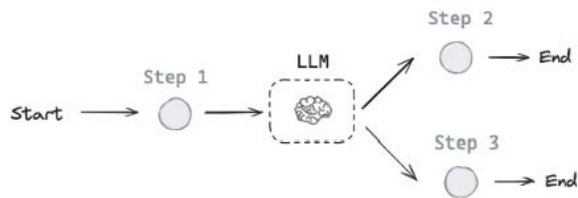
Langgraph



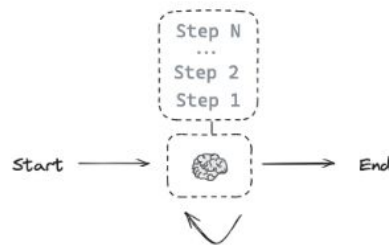
LangGraph

Agents

Router



Fully Autonomous

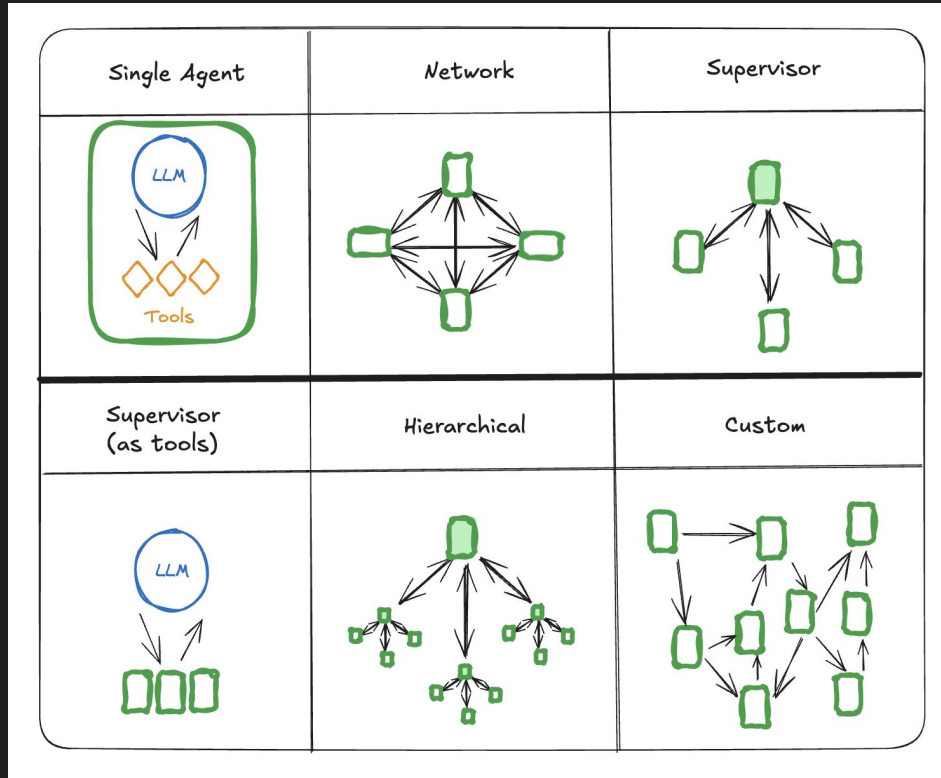


Less

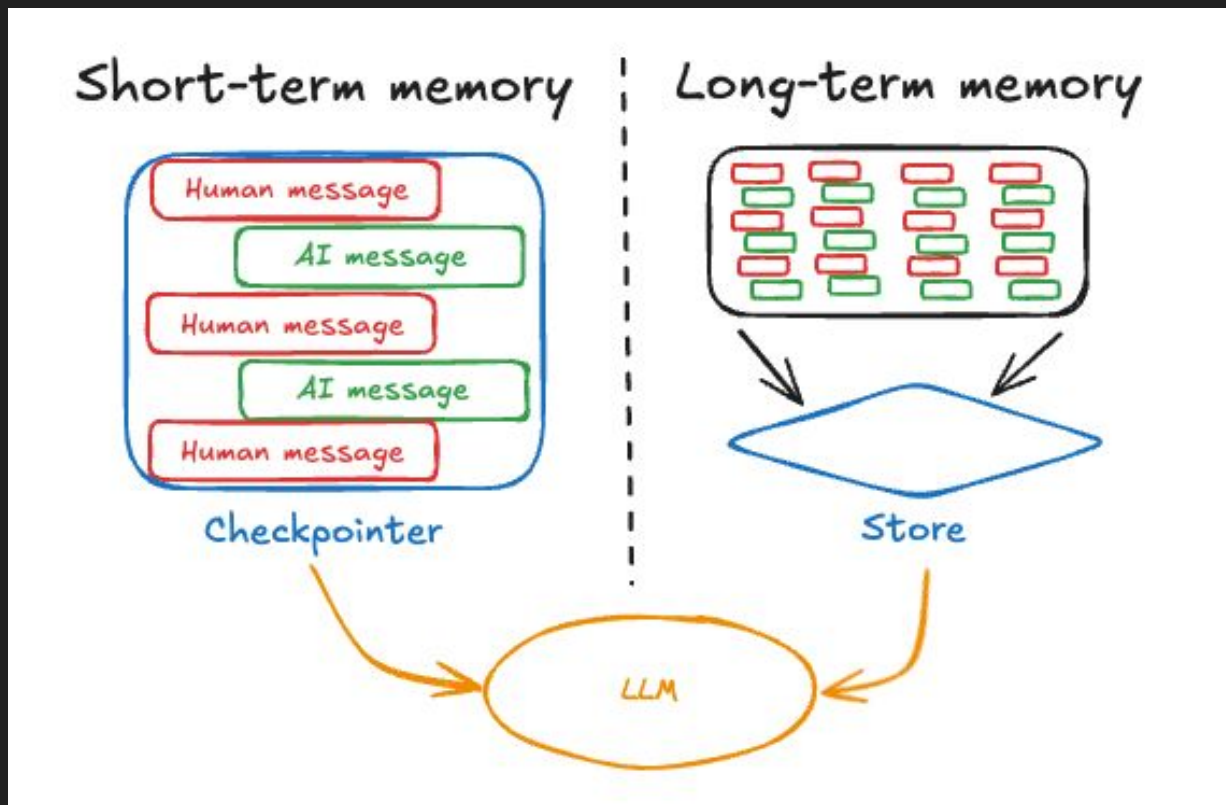
Control

More

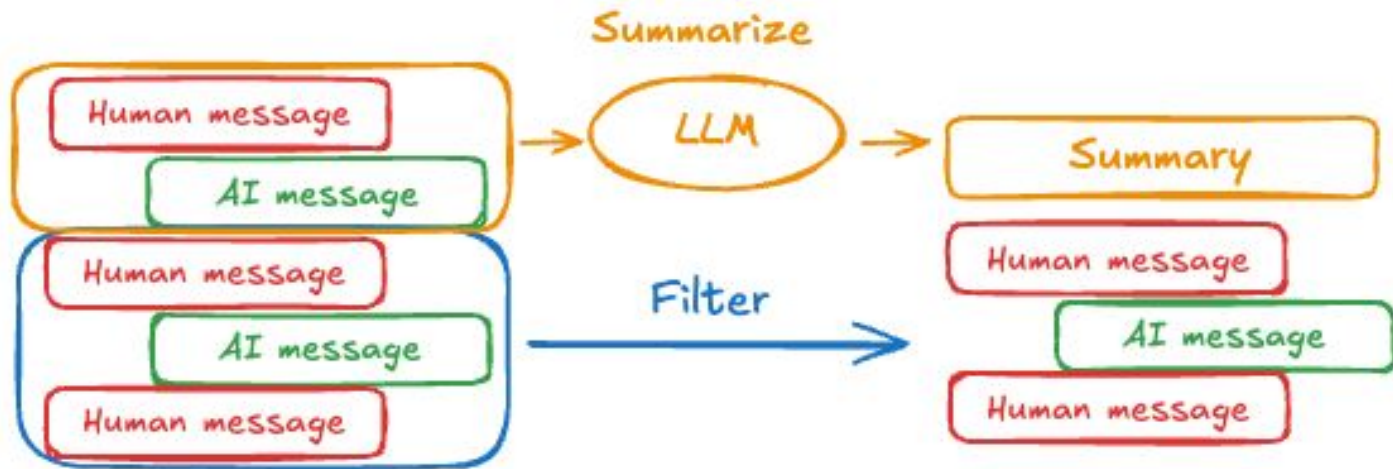
Multi-agents



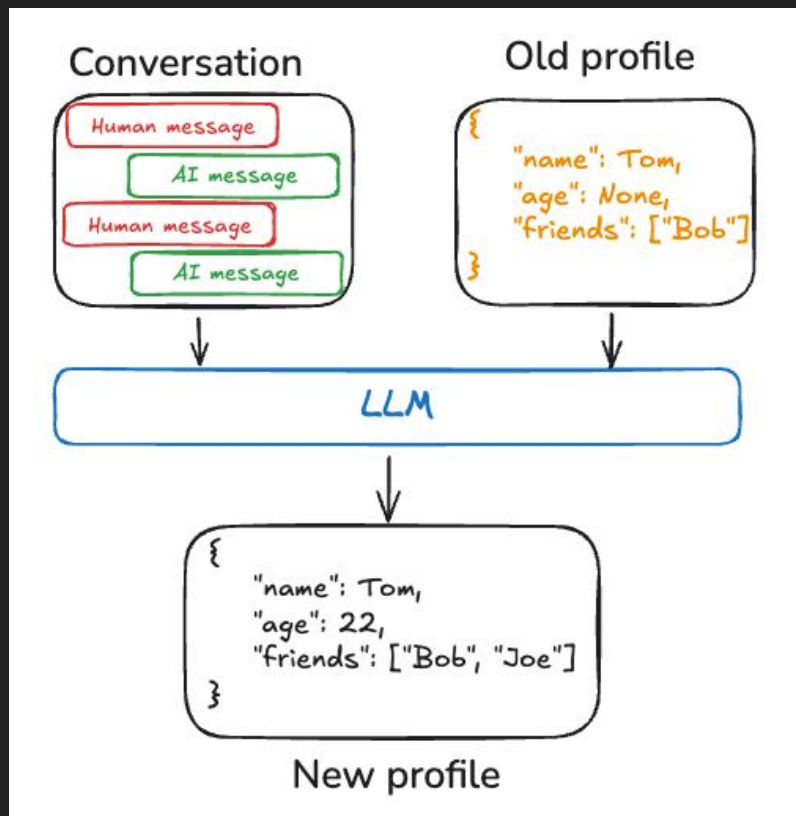
Memory



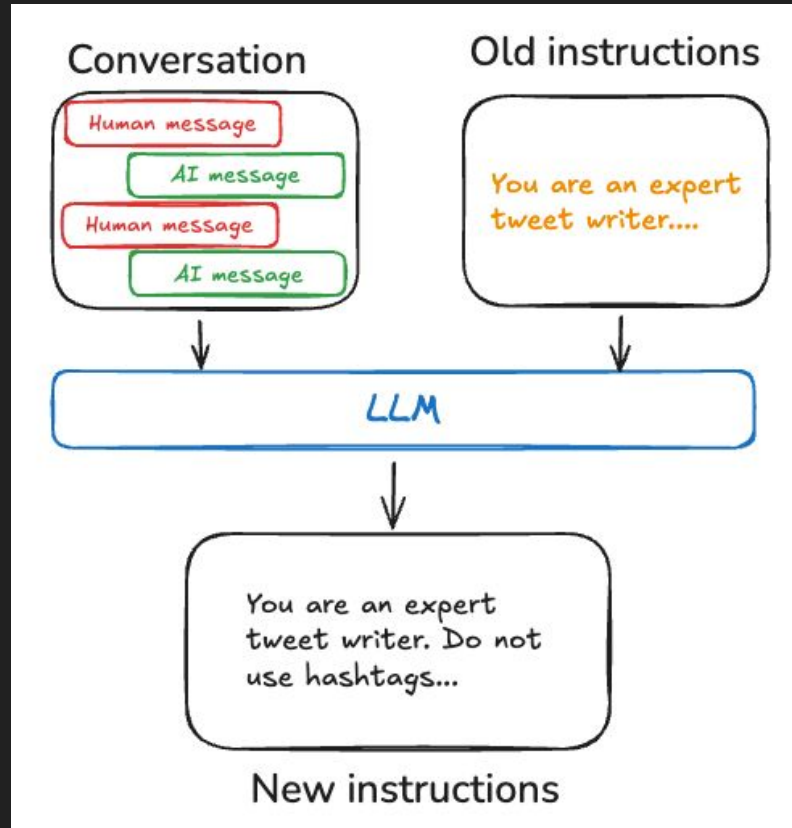
Short-term memory



Long-term memory



Meta-prompting



Meta-prompting

