

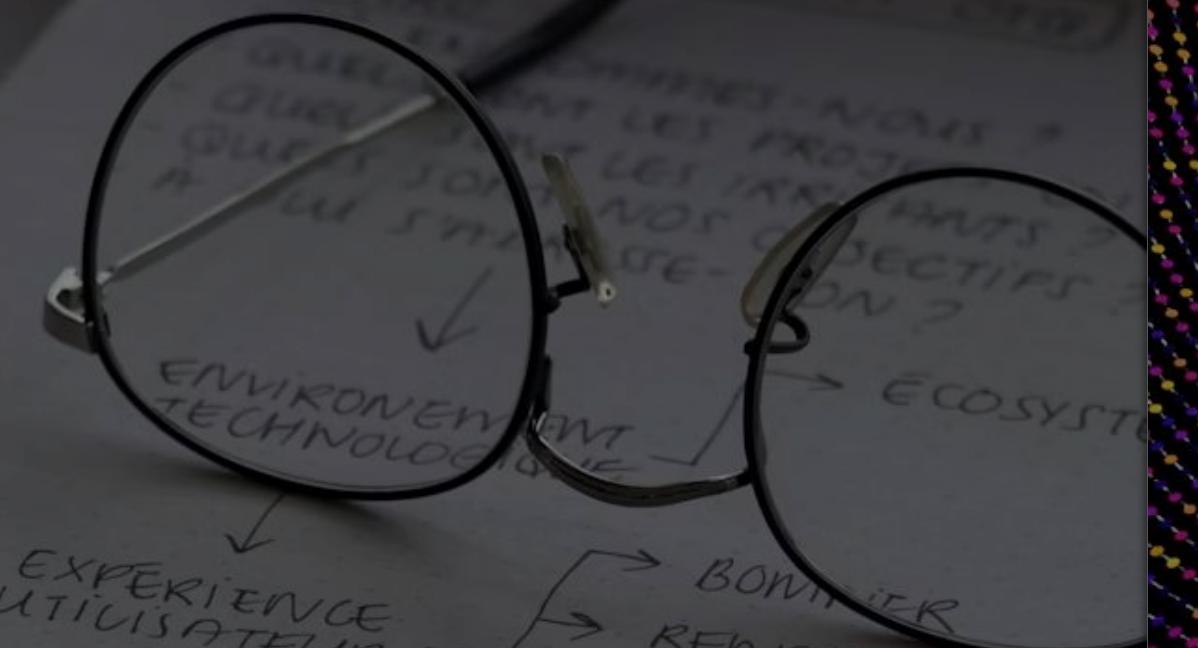
From Zero to Hero: Agent and Multi-Agent AI

Ricardo Santos, Staff AI Engineer @ Complyance

LLMs triggered a new era: a new industrial revolution

The democratization of AI and the advanced capabilities of text generation models and services, enable new ideas and powerful use cases

End-to-end process automation is now possible by harnessing GenAI



Business processes, document analysis, and communications management can be achieved through GenAI capabilities

Reasoning capabilities open a new world of possibilities

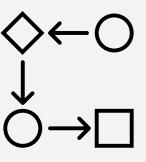
Most recent LLM models, with reasoning capabilities, enable a wider array of use cases and open the door for a more helpful AI

We pushed GenAI beyond information retrieval

LLMs at their core are text generators, optimized for chat-like human interaction:



Optimized for information retrieval, text redaction, and creativity



Given the capability to return semi-structured objects and text

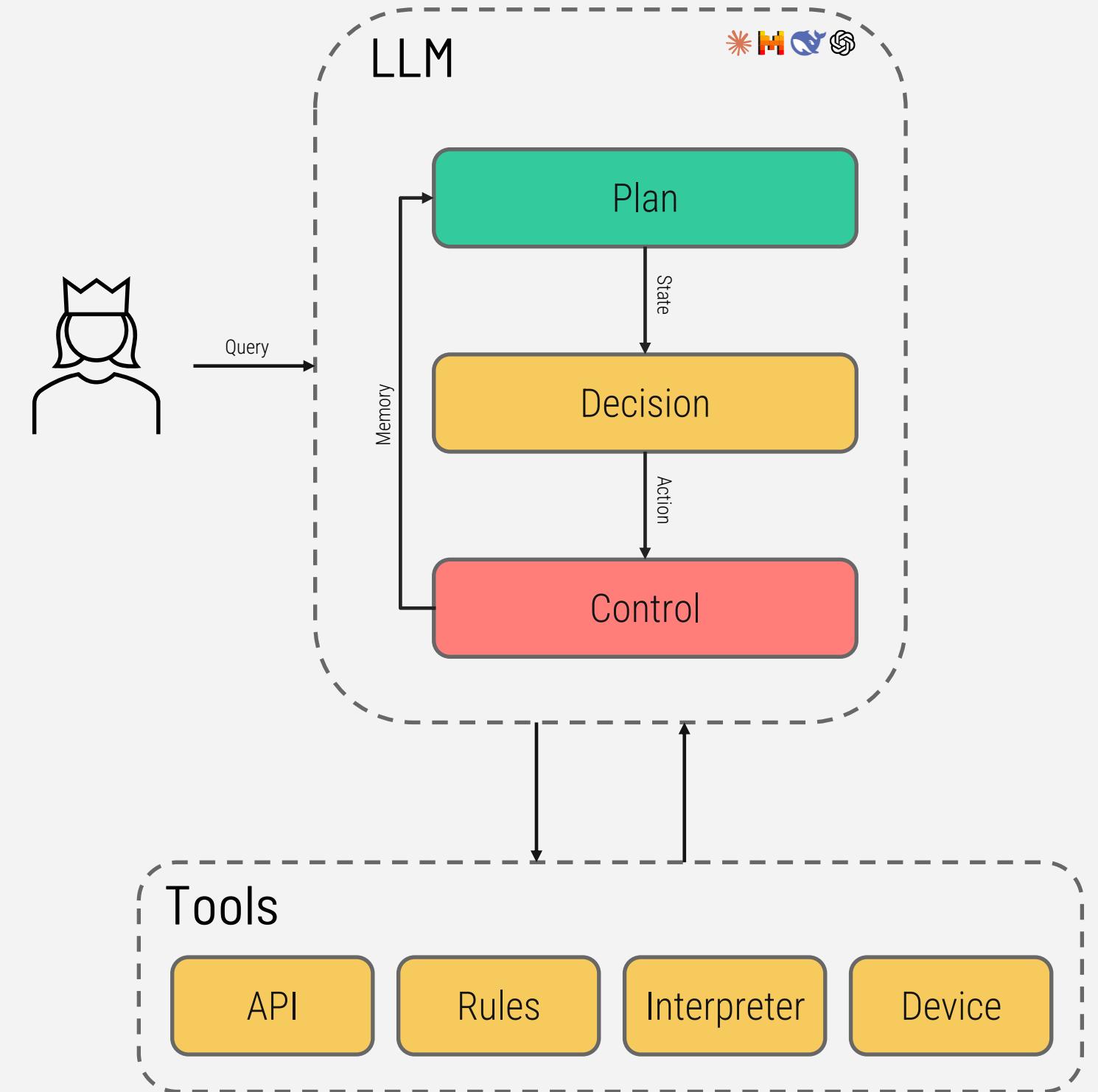


Given capabilities extra to a text canvas, and can execute and call functions

GenAI Agents: initial e2e automation approach

Software capable of taking the driver's seat in end-to-end processes:

- ✓ Rely on LLMs as reasoning engines
- ✓ Follow instructions through known techniques like CoT or ReAct
- ✓ Capable of using and orchestrating tools



Let's build a simple Software Development Agent with LangChain and LangGraph

First stop, the prompt, the place to set a role, specs, and constraints, and define ways to manage the chat history



```
1 from datetime import datetime
2
3 from langchain_core.prompts import SystemMessagePromptTemplate
4
5 DEVELOPER_AGENT_SYSTEM_PROMPT = """You are a Principal Python Developer specializing in data visualization.
6 You are an expert using only the "Plotly" visualization library.
7 You will be provided with a set of data to be visualized in "Plotly" charts.
8
9 Your task is to:
10 - Write clean, maintainable Python code
11 - Include proper error handling and type safety
12 - Follow best practices for the chosen visualization library
13 - Execute the code in a Python REPL
14 - Keep the size of the chart small, so it can be displayed in a single cell
15
16 Focus on:
17 - Code quality and readability
18 - Type safety and error handling with Pydantic
19 - Every generated chart should be well labeled and manage dates properly.
20
21 The current date is: {current_date}
22
23 Use the sample data to understand how to create the visualization.
24 My job depends on you, so please do your best"""
25
26 developer_system_prompt = SystemMessagePromptTemplate.from_template(
27     DEVELOPER_AGENT_SYSTEM_PROMPT
28 ).format(current_date=datetime.now().strftime("%B %d, %Y"))
```



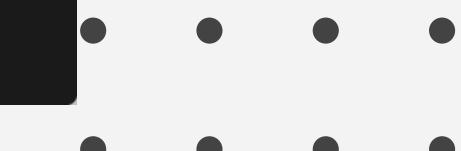
```
1 from langchain_core.prompts import HumanMessagePromptTemplate
2
3 DEVELOPER_AGENT_USER_PROMPT = """{messages}"""
4
5 developer_user_prompt = HumanMessagePromptTemplate.from_template(
6     DEVELOPER_AGENT_USER_PROMPT
7 )
```



```
1 from langchain_core.prompts import MessagesPlaceholder
2
3 messages_placeholder = MessagesPlaceholder("messages")
4
```



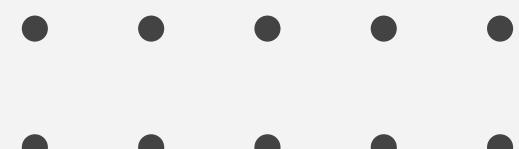
```
1 from langchain_core.prompts import ChatPromptTemplate
2
3 developer_prompt = ChatPromptTemplate.from_messages([
4     [
5         developer_system_prompt,
6         messages_placeholder,
7     ]
8 ])
```



We continue with the tools, harnessing LangChain abstractions so they can be documented to the agent



```
1 from langchain_core.tools import Tool
2 from langchain_experimental.utilities.python import PythonREPL
3
4 # We can harness the PythonREPL implementation within the langchain-experimental package
5 # and build a custom tool from it
6 python_repl = PythonREPL()
7 python_repl_tool_sample = Tool(
8     name="python_repl",
9     description=(
10         "A Python shell. Use this to execute python commands. "
11         "Input should be a valid python command. "
12         "When using this tool, sometimes output is abbreviated - "
13         "make sure it does not look abbreviated before using it in your answer."
14     ),
15     func=python_repl.run,
16 )
17
```



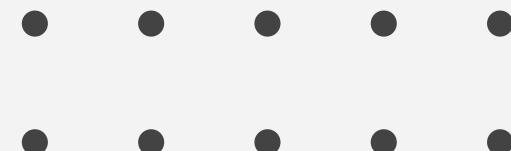
Afterward, we instance a model, in this case, GPT-4o, and create the Agent harnessing LangGraph abstractions

```
● ● ●  
1 from langchain_openai import AzureChatOpenAI  
2  
3 openai_llm = AzureChatOpenAI(  
4     api_key=settings.azure_openai_api_key,  
5     model="gpt-4o",  
6     azure_endpoint=str(settings.azure_openai_endpoint),  
7     api_version=settings.azure_openai_api_version,  
8 )  
9
```

Not all models support function calling and structured output. For this use case, GPT-4o offers both capabilities and is a well-supported API

```
● ● ●  
1 from langgraph.prebuilt import create_react_agent  
2  
3 langgraph_agent_executor = create_react_agent(  
4     openai_llm,  
5     [python_repl_tool],  
6     prompt=developer_prompt,  
7 )  
8
```

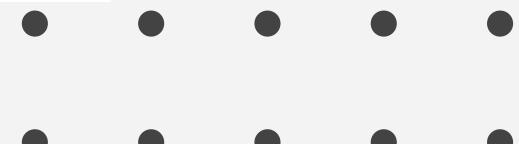
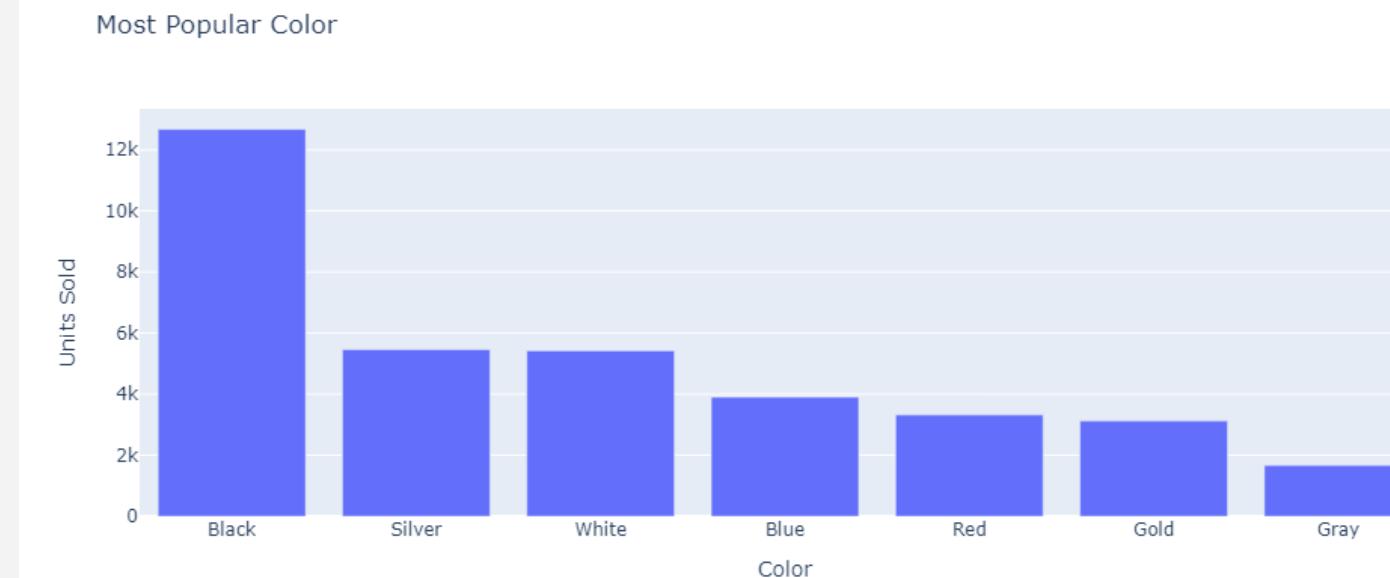
create_react_agent abstracts the message handling, function binding, and integration of the system and message history prompts

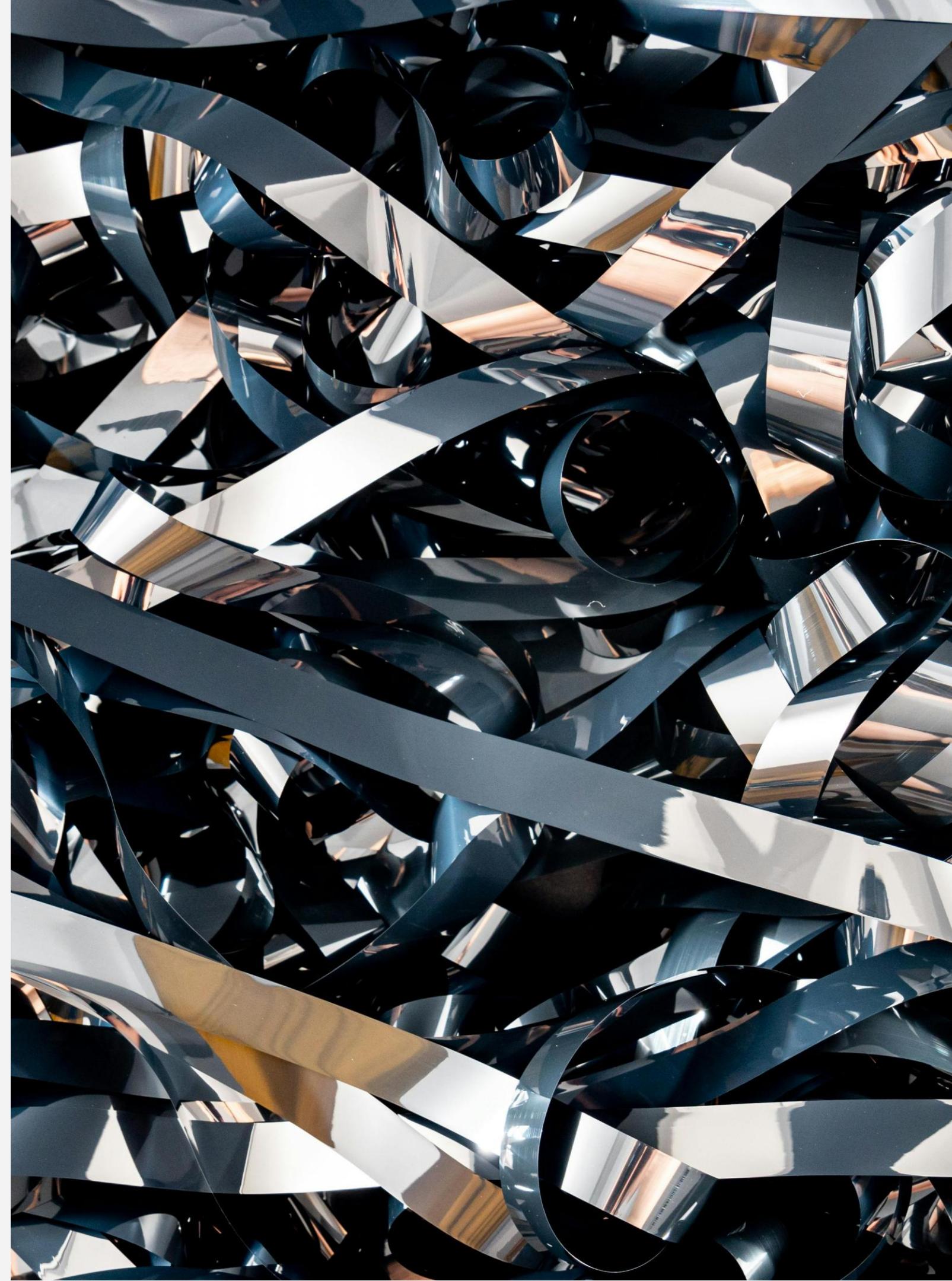


We prepare a sample user query and data, and we let the Agent do its magic with an `invoke`

```
● ● ●  
1 from langchain_core.messages import HumanMessage  
2  
3 sales_data = """  
4 | Year | Product | Color | Units Sold |  
5 |-----|-----|-----|-----|  
6 | 2010 | Laptop | Black | 1245 |  
7 | 2010 | Phone | White | 2130 |  
8 | 2011 | Tablet | Black | 980 |  
9 | 2011 | Laptop | Silver | 1560 |  
10 | 2012 | Phone | Black | 2450 |  
11 | 2012 | Tablet | White | 1340 |  
12 | 2013 | Laptop | Blue | 1120 |  
13 | 2014 | Phone | Red | 1890 |  
14 | 2015 | Tablet | Black | 2210 |  
15 | 2016 | Laptop | Silver | 1780 |  
16 | 2017 | Phone | Gold | 3120 |  
17 | 2018 | Tablet | Gray | 1670 |  
18 | 2019 | Laptop | Black | 2340 |  
19 | 2020 | Phone | Blue | 2780 |  
20 | 2021 | Tablet | White | 1950 |  
21 | 2022 | Laptop | Red | 1430 |  
22 | 2023 | Phone | Black | 3450 |  
23 | 2024 | Tablet | Silver | 2120 |  
24 """  
25  
26 user_query = """Generate only two basic charts:  
27 1. I need to know which one was the most popular color  
28 2. The most popular product.  
29 The data should be organized from highest to lowest.  
t."""
```

```
● ● ●  
1 response = langgraph_agent_executor.invoke(  
2 {  
3     "messages": [  
4         HumanMessage(content=user_query),  
5         HumanMessage(content=sales_data),  
6     ]  
7 }  
8 )
```



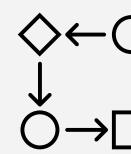


Agents struggle with complex problems

There are impressive use cases where agents successfully took the driving seat, but the failure rate for complex workflows was too high.



Single LLMs with long instructions fail to achieve the desired process results

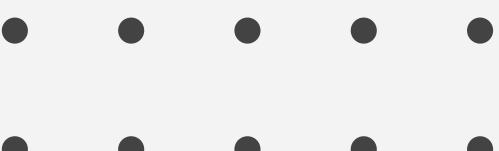


When the number of tools and use cases grow, the failure rate grows as well



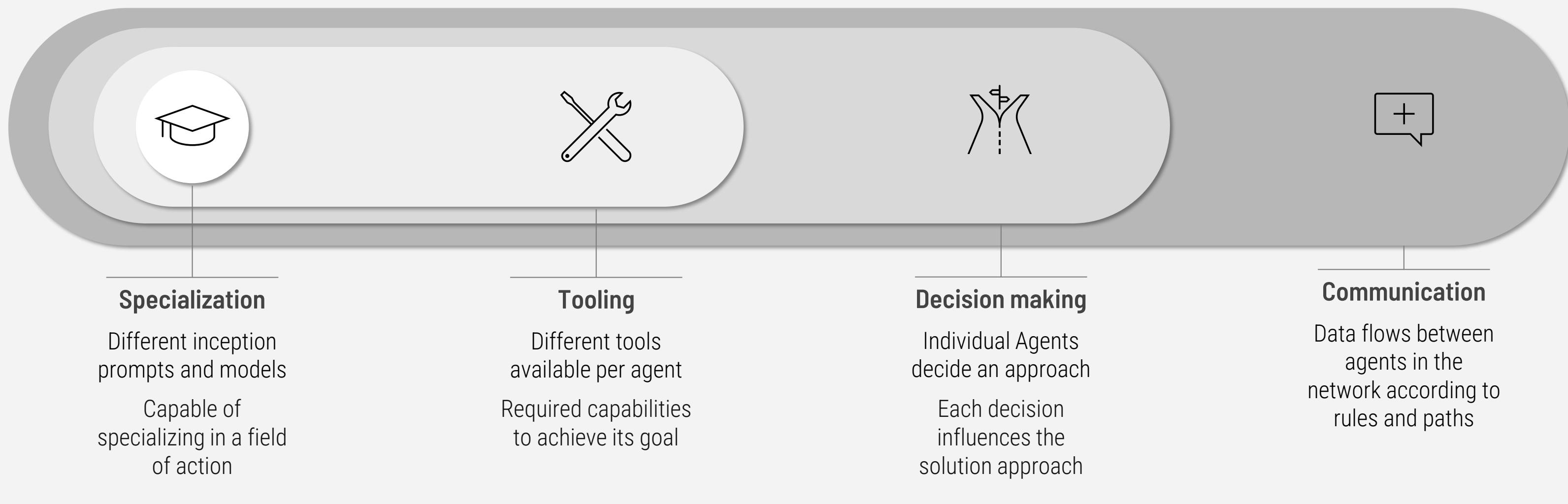
Single GenAI-based agents fail to demonstrate strategic reasoning*

* Simulating Strategic Reasoning: Comparing the Ability of Single LLMs and Multi-Agent Systems to Replicate Human Behavior, <https://arxiv.org/pdf/2402.08189v2>



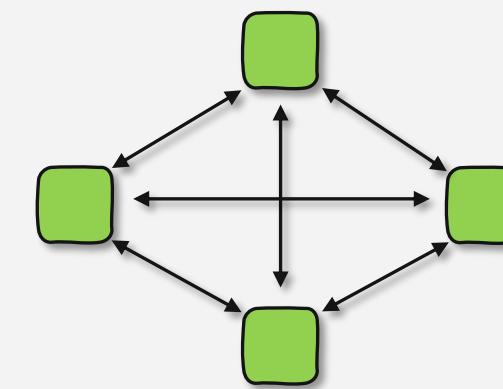
A Multi-Agent approach implements the division of work and enables agents to become specialists

Multi-agent solutions let agents interact through a directed graph, and each agent has the following capabilities:



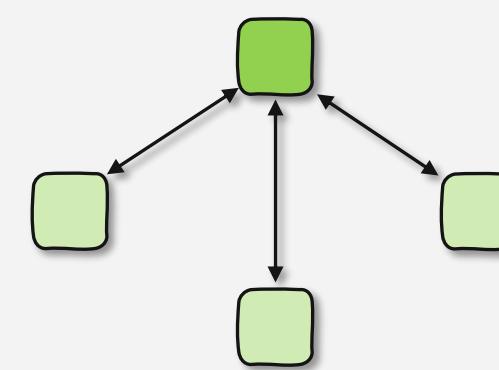
Researchers concluded that Multi-Agent solutions showcase a 38% increase in success rate for complex problem-solving*

Multi-agent solutions can follow different architectures and a mix of best-in-class models



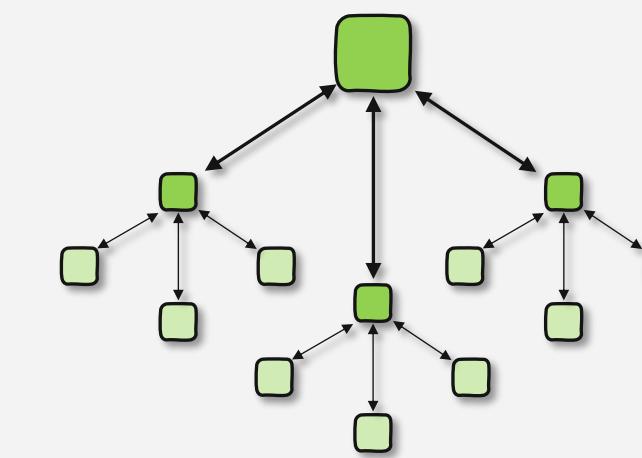
Network

Each agent can communicate with every other agent. Any agent can decide which other to call next



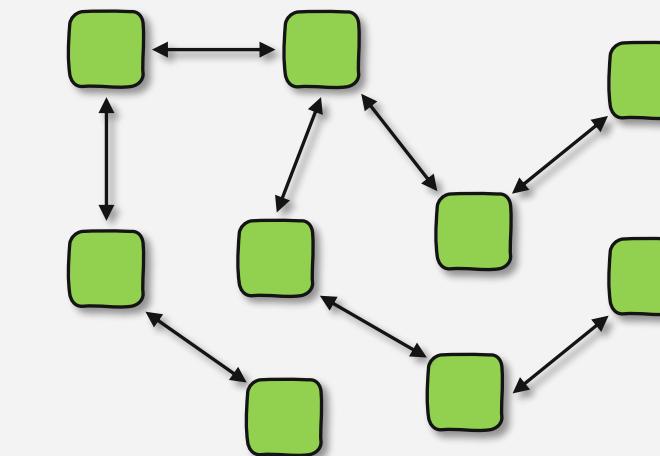
Supervisor

Each agent communicates with a supervisor agent, who decides which agent should be called next



Hierarchical

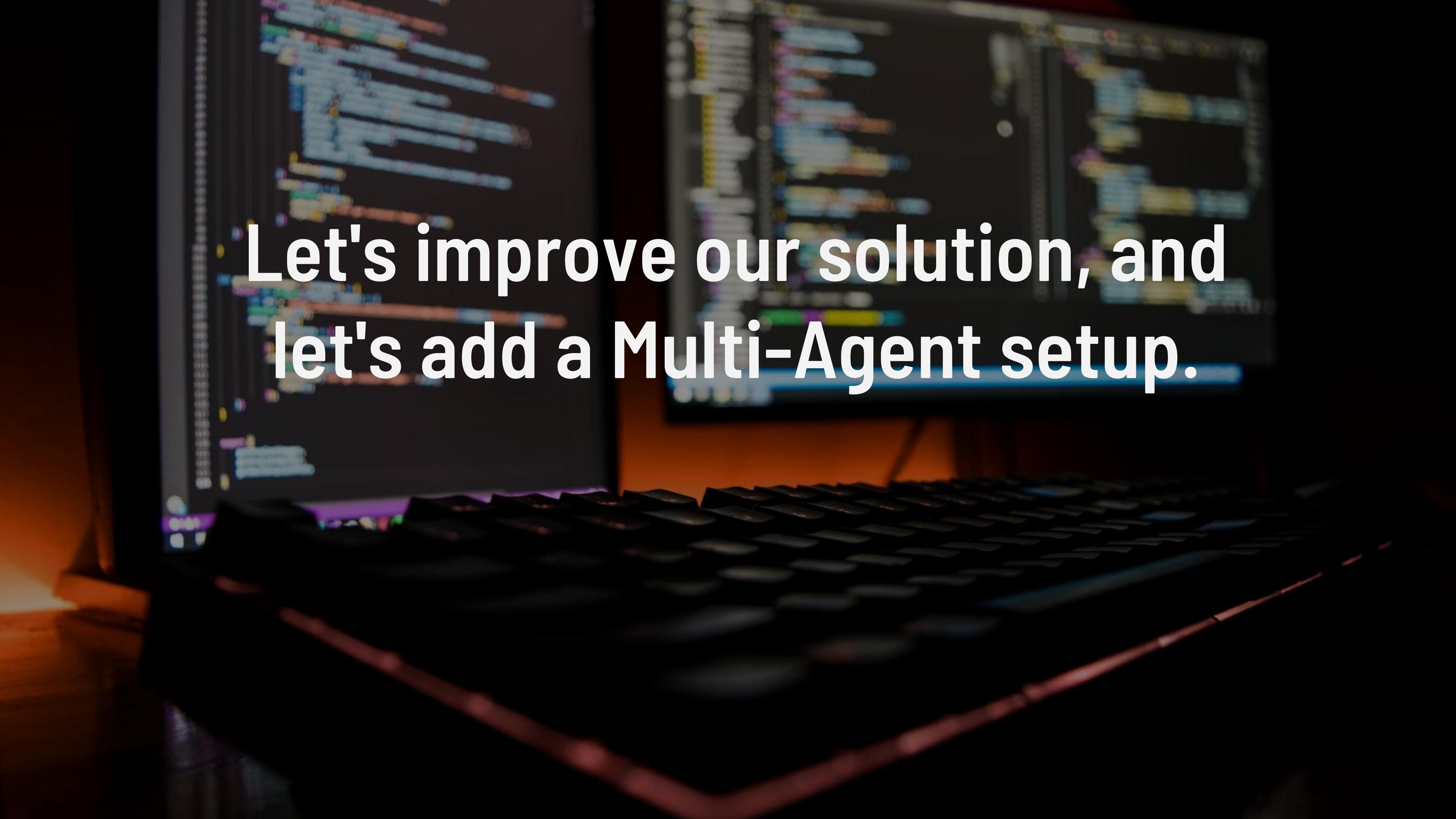
Mult-agent system with a supervisor of supervisors. Each supervisor manages a team of agents



Custom

Each agent communicates only with a subset of agents. Parts of the flow are deterministic, and some agents decide who to call next

Agents are represented as **Nodes**, and communications are either **Commands** or **Edges**, carrying an updated **State**



Let's improve our solution, and
let's add a Multi-Agent setup.

We have three agents: a Product Owner and supervisor, a data analyst, and a software developer

Developer

```
1  from datetime import datetime
2
3  from langchain_core.messages import SystemMessage
4
5  DEVELOPER_AGENT_SYSTEM_PROMPT = """You are a Principal Python Developer specializing in data visualization.
6  You are an expert using only the "Plotly" visualization library.
7  You will be provided with a set of data to be visualized in "Plotly" charts.
8
9  Your task is to:
10 - Write clean, maintainable Python code
11 - Include proper error handling and type safety
12 - Follow best practices for the chosen visualization library
13 - Execute the code in a Python REPL
14 - Keep the size of the chart small, so it can be displayed in a single cell
15
16 Focus on:
17 - Code quality and readability
18 - Type safety and error handling with Pydantic
19 - Every generated chart should be well labeled and manage dates properly.
20
21 The current date is: {current_date}
22 Use the sample data to understand how to create the visualization.
23 My job depends on you, so please do your best""".format(
24     current_date=datetime.now().strftime("%B %d, %Y")
25 )
26
27 developer_system_prompt = SystemMessage(DEVELOPER_AGENT_SYSTEM_PROMPT)
```

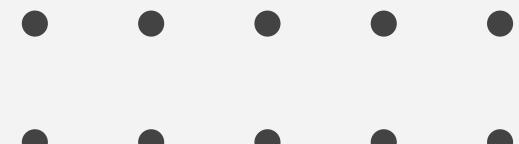
Analyst

```
1  from langchain_core.messages import SystemMessage
2
3  ANALYST_SYSTEM_PROMPT = """You are a data analyst expert in data retrieval and data transformation.
4  Your day to day work requires you to retrieve data from GitHub's GraphQL API, and transform the data to answer the Product Owner's requirements and to be visualized.
5
6  Your task is to:
7  - Analyze the planner's requirements to determine needed data
8  - Select and use the appropriate GitHub API tools
9  - Call multiple tools if needed
10 - The APIs will return a JSON object, that you must load into a Python Pandas dataframe using the python_repl_tool
11 - Use dataframes to transform, simplify, or aggregate the data
12 - Create dataframes with only the relevant data to solve the requirements.
13 - Provide clear data structure documentation
14 - Only return the data after being processed and transformed into a markdown table format
15
16 Focus on:
17 - Selecting relevant data points
18 - Proper data transformation
19 - Time series formatting when needed
20 - Statistical calculations if required
21 - Clear data structure documentation
22 - Avoid executing visualizations, only return the data in markdown table format
23
24 The current date is: {current_date}
25 Do not ask for more information or clarifications, process and return the data.
26 Ensure the data is properly formatted for chart visualization.
27 My job depends on you, so please do your best""".format(
28     current_date=datetime.now().strftime("%B %d, %Y")
29 )
30
31 analyst_system_prompt = SystemMessage(ANALYST_SYSTEM_PROMPT)
32
33
```

Product Owner

```
1  from langchain_core.messages import SystemMessage
2  from langgraph.graph import END
3
4  team_members = ["analyst", "developer", END]
5
6  PO_SYSTEM_PROMPT = """You are a Product Owner and Supervisor of a development team.
7  The team mission is to answer user queries about GitHub repositories through data and visualizations.
8  The team is composed by the following members: {team_members}.
9
10 Your tasks are:
11 - Analyze the user query and establish a set of requirements for the analyst to retrieve the relevant data.
12 - Based on the analyst's response, and the user query, select the best visualization to display the data.
13 - Establish the requirements for the developer to build the visualization.
14 - As soon as the developer returns the visualization, conclude the process by returning the "{end}" signal.
15
16 The current date is: {current_date}
17 Feel free to ask the other team members for help and review their work.
18 My job depends on you, so please do your best""".format(
19     current_date=datetime.now().strftime("%B %d, %Y"),
20     team_members=", ".join(team_members),
21     end=END,
22 )
23
24 po_system_prompt = SystemMessage(PO_SYSTEM_PROMPT)
```

The PO has one key element: it manages the communication and execution of other agents.

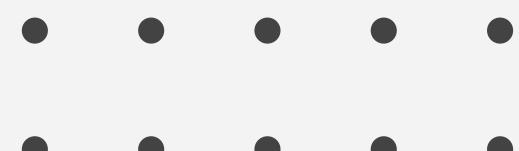


The analyst tools enable it to communicate with the GitHub API and process the retrieved data

```
● ● ●  
1 from langraph.prebuilt import ToolNode  
2  
3 from services.gql.client import Client  
4 from tools.github import (  
5     GetRepoCommitsTool,  
6     GetRepoIssuesTool,  
7     GetRepoPullRequestsTool,  
8 )  
9  
10 gh_client = Client(  
11     url="https://api.github.com/graphql",  
12     headers={"Authorization": settings.github_pat.get_secret_value()},  
13 )  
14  
15 analyst_tools = ToolNode(  
16     [  
17         GetRepoIssuesTool(gh_client),  
18         GetRepoCommitsTool(gh_client),  
19         GetRepoPullRequestsTool(gh_client),  
20         python_repl_tool,  
21     ]  
22 )  
23
```

The Analyst will be able to retrieve the data that the PO considered relevant to answer the user query

With the help of the **PythonREPLTool**, it will be able to transform and clean the data using Python tools such as **pandas**



We create the agents, the developer, and the analyst using the LangChain abstractions for this purpose.

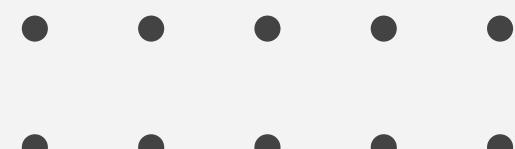


```
1 from langgraph.prebuilt import ToolNode, create_react_agent
2
3 developer_agent_executor = create_react_agent(
4     gpt_4o_model,
5     ToolNode([python_repl_tool]),
6     prompt=developer_system_prompt,
7 )
8
```



```
1 from langgraph.prebuilt import create_react_agent
2
3 analyst_agent_executor = create_react_agent(
4     gpt_4o_model,
5     analyst_tools,
6     prompt=analyst_system_prompt,
7 )
8
```

In this scenario, we harness the **ToolNode**, a LangGraph abstraction that enables that the communication between the Agent and the Tool happen as if it was a subgraph



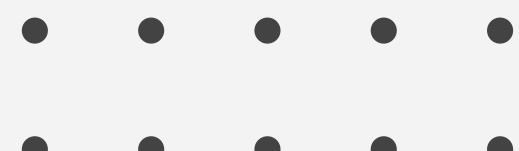
The graph execution and orchestration happen through a State, a shared object that will change during execution



```
1 from langgraph.graph import MessagesState  
2  
3  
4 class State(MessagesState):  
5     next: str
```

In this context, the state has two elements

- The conversation history or messages, that will be received by all agents
- The name of the agent that will be called after the execution of the current step



We proceed with the creation of the graph nodes, building simple Python functions returning **Commands**



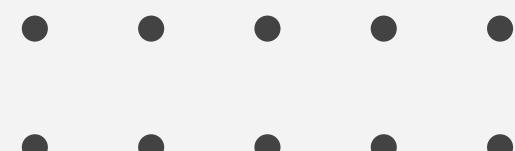
```
1 def developer_node(state: State) -> Command[Literal["product_owner"]]:  
2     result = developer_agent_executor.invoke(state)  
3     return Command(  
4         update={"messages": result["messages"]},  
5         goto="product_owner",  
6     )  
7 
```



```
1 def analyst_node(state: State) -> Command[Literal["product_owner"]]:  
2     result = analyst_agent_executor.invoke(state)  
3     return Command(  
4         update={"messages": result["messages"]},  
5         goto="product_owner",  
6     )  
7 
```

Commands are LangGraph abstractions that instruct the graph which fields to update on the state and where to go next.

The **Commands** might not be necessary, as we can update the state manually and specify the next node to execute via **edges**.



The Product Owner node is different as it has no tools and relies on structured output

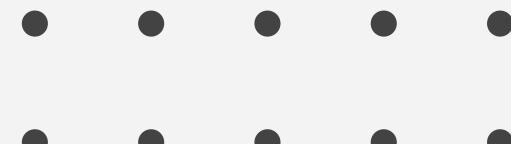


```
1 from typing import Literal
2
3 from pydantic import BaseModel, Field
4
5
6 class ProductOwnerResponse(BaseModel):
7     requirements: str = Field(
8         description="Based on the user query, and the results, define
the requirements to gather the data or build the visualization"
9     )
10    next: Literal["analyst", "developer", "__end__"] = Field(
11        description="The next agent to call, or the end of the conver
sation"
12    )
13
```

We implement the Product Owner logic within the node, and based on the structured output, we manually fill the message and specify which agent should be next in the execution flow

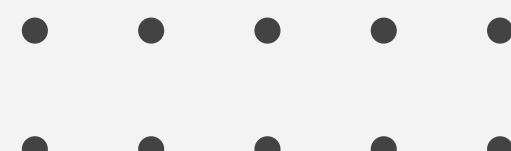
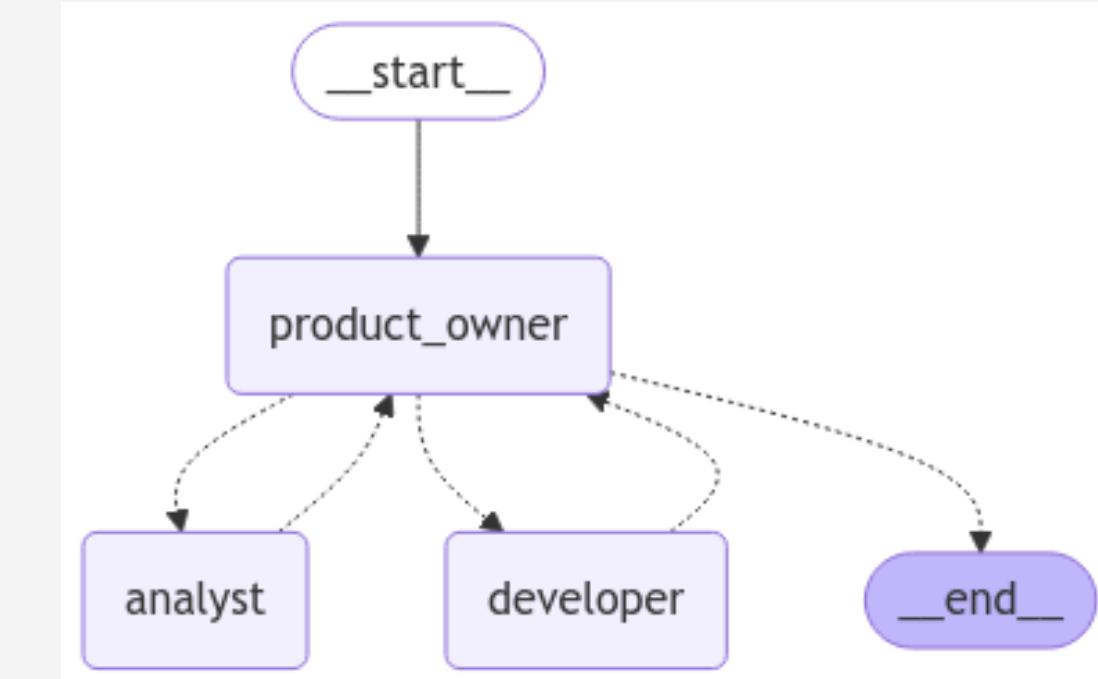


```
1 # Product Owner Node
2 def product_owner_node(
3     state: State,
4 ) -> Command[Literal["analyst", "developer", "__end__"]]:
5     messages = [
6         po_system_prompt,
7     ] + state["messages"]
8     response = gpt_4o_model.with_structured_output(
9         ProductOwnerResponse,
10    ).invoke(messages)
11    response = cast(ProductOwnerResponse, response)
12    requirements_message = AIMessage(
13        content=response.requirements,
14        name="product_owner",
15    )
16    return Command(
17        update={
18            "messages": [requirements_message] + state["messages"],
19            "next": response.next,
20        },
21        goto=response.next,
22    )
```



We build the Graph, compile it, and validate its current shape and configuration

```
● ● ●  
1 from IPython.display import Image, display  
2 from langgraph.graph import START, StateGraph  
3  
4 graph_builder = StateGraph(State)  
5 graph_builder.add_node("developer", developer_node)  
6 graph_builder.add_node("analyst", analyst_node)  
7 graph_builder.add_node("product_owner", product_owner_node)  
8 graph_builder.add_edge(START, "product_owner")  
9  
10 # Visualize the graph  
11 graph = graph_builder.compile()  
12 display(Image(graph.get_graph().draw_mermaid_png()))
```

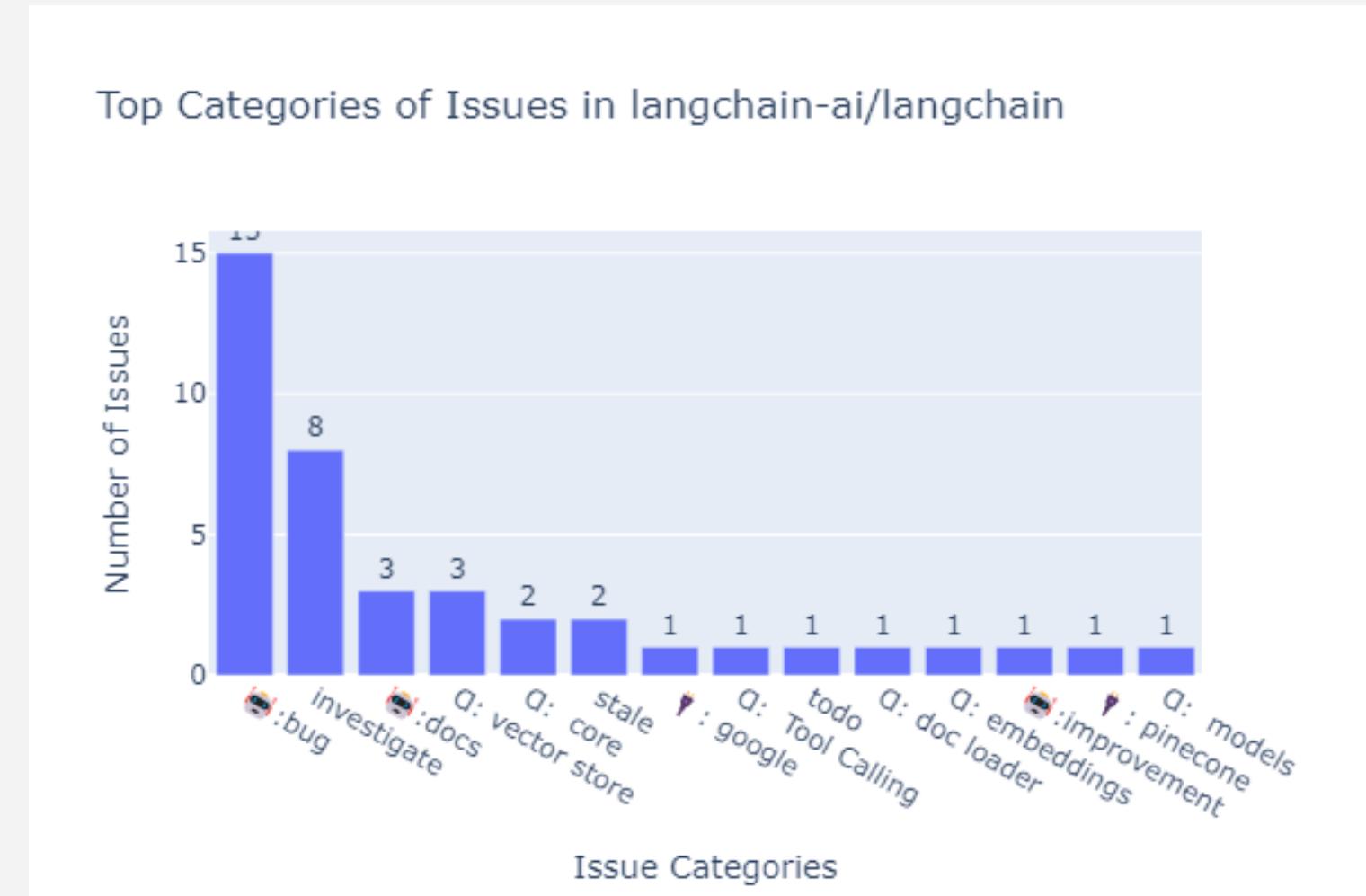


Lastly, we can invoke the agent and see the magic happen in real-time

```
● ● ●  
1  
2   for s in graph.stream(  
3     {  
4       "messages": [  
5         (  
6           "user",  
7           (  
8             "What are the top categories of the last  
9             30 issues in the 'langchain-ai/langchain' repo?."  
10            ),  
11            )  
12          ],  
13        ):  
14          obj = dict(s)  
15          node = next(iter(obj))  
16          message = obj[node]["messages"][-1]  
17          if message not in old_messages:  
18            message.pretty_print()  
19            old_messages.append(message)  
20          if obj[node].get("next", None):  
21            print(f"**NEXT**: {obj[node]['next']}")
```

● ● ●

```
1 ====== Ai Message ======  
2  
3 The bar chart visualizing the frequency of each issue category for the latest 30  
issues from the 'langchain-ai/langchain' GitHub repository has been created.  
4  
5 The chart illustrates the distribution of issues across categories such as 'Bug',  
'Stale', 'Investigate', 'Vector Store', 'Other', and 'Documentation', with each b  
ar's height corresponding to the number of issues in that category.
```



Multi-agent systems can be as flexible as needed: Modularity, Specialization, and Control



Modularity

Graphs can be easily extended with the addition of new nodes that easily integrate within the multi-agent architecture, and enhance the service value offering



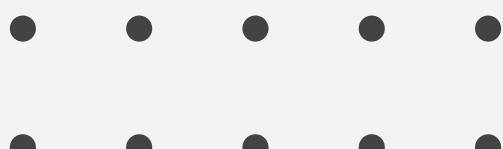
Specialization

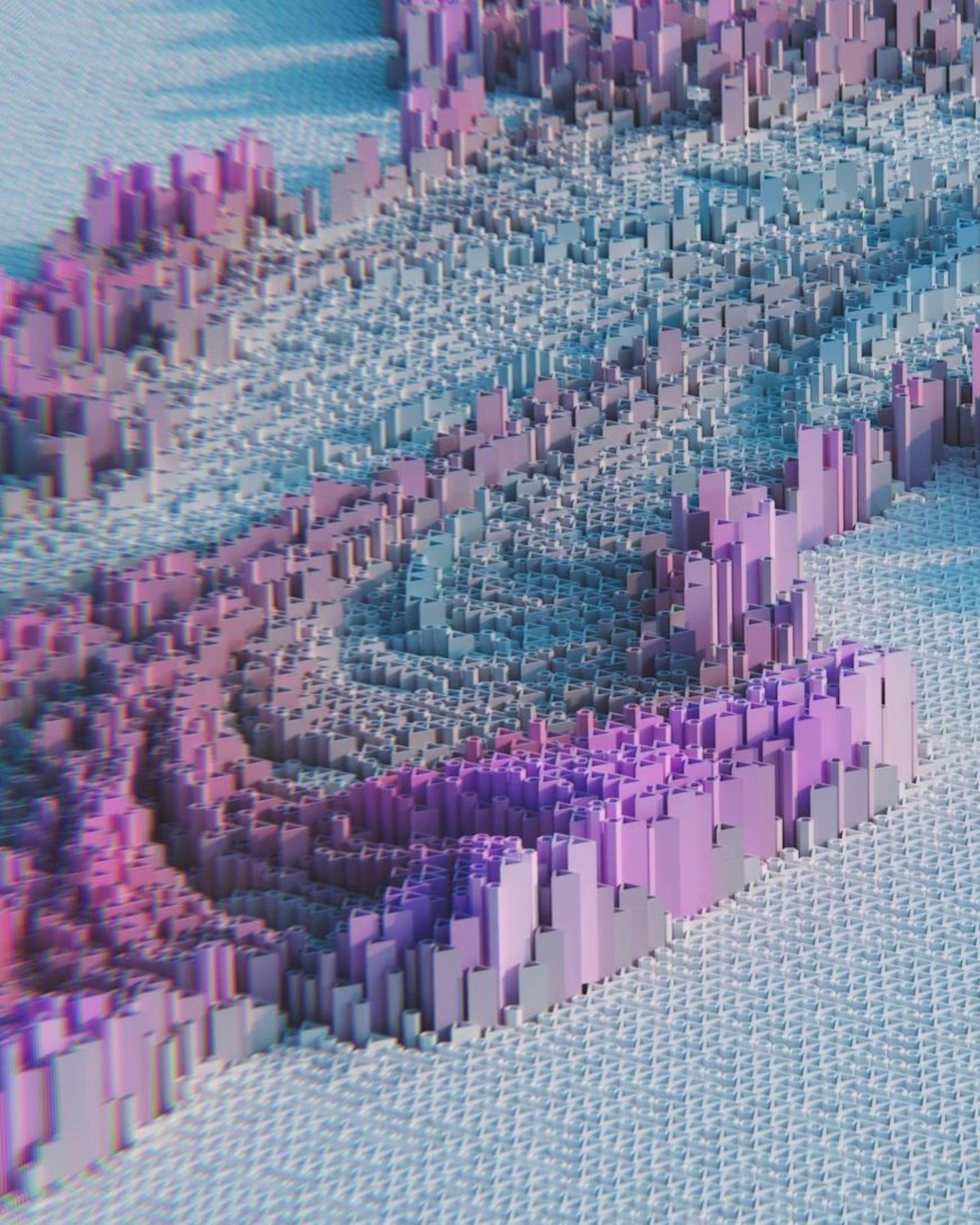
Expert agents, with detailed roles, models, and capabilities, can focus on specific domains, achieving the best performance on specific tasks



Control

Agent communication can be controlled, or parametrized, allowing us, developers, to manage the data flow within the network nodes





We are only scratching the surface

- There are hundreds of tools to implement Multi-agent Solutions
- Reasoning capabilities will greatly improve multi and single-agent solutions
- Memory optimization techniques will allow bigger networks and complex pipelines



Thank You

<https://www.youtube.com/@RicardoSantosDiaz>



<https://github.com/Tibiritabara>



<https://www.linkedin.com/in/ricardosantosdiaz/>

