

CALLING “HOUSE” ON YOUR AI BUZZWORD BINGO CARD

CARLY RICHMOND

ABOUT ME

- Developer Advocate Lead @
 elastic
- Frontend Engineer, Speaker & Blogger



SCAN ME





BINGO CARD

AI

LLM

AGENTS

PROMPT
ENGINEERING

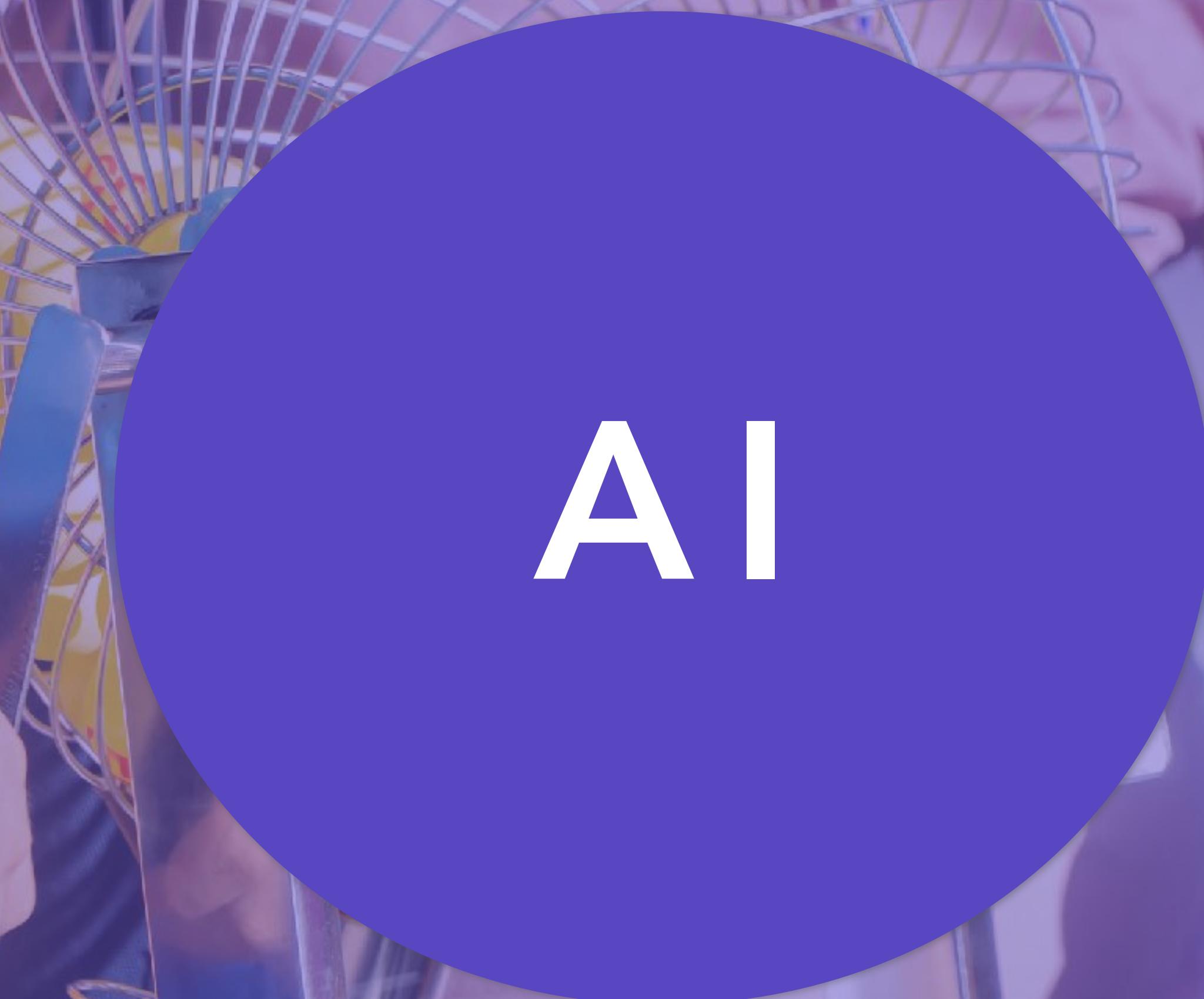
CONTEXT
ENGINEERING

VECTOR
SEARCH

RAG

A2A

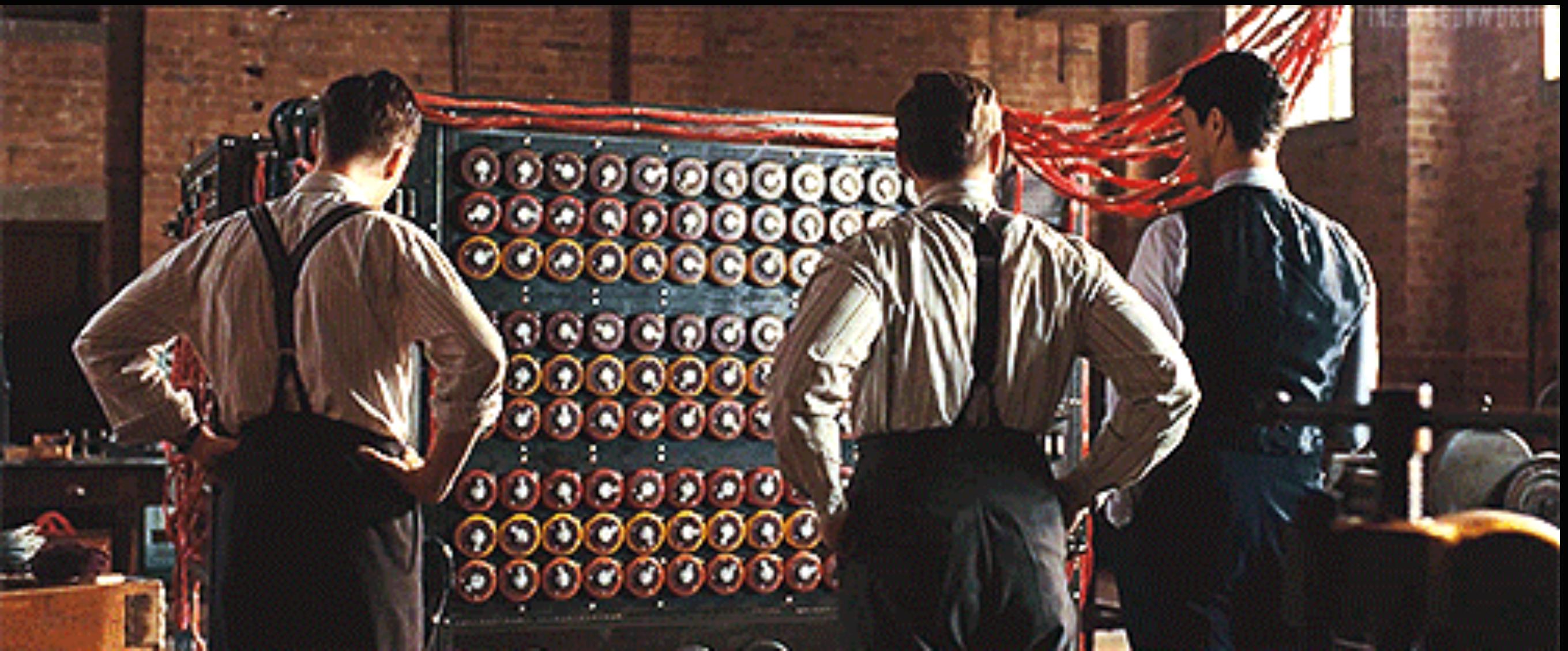
MCP



AI

“Artificial intelligence (AI) is the capability of computational systems to perform tasks typically associated with human intelligence, such as learning, reasoning, problem-solving, perception, and decision-making.”

-ARTIFICIAL INTELLIGENCE (AI), WIKIPEDIA



A. M. Turing (1950) Computing Machinery and Intelligence. *Mind* 49: 433-460.

COMPUTING MACHINERY AND INTELLIGENCE

By A. M. Turing

1. The Imitation Game

I propose to consider the question, "Can machines think?" This should begin with definitions of the meaning of the terms "machine" and "think." The definitions might be framed so as to reflect so far as possible the normal use of the words, but this attitude is dangerous. If the meaning of the words "machine" and "think" are to be found by examining how they are commonly used it is difficult to escape the conclusion that the meaning and the answer to the question, "Can machines think?" is to be sought in a statistical survey such as a Gallup poll. But this is absurd. Instead of attempting such a definition I shall replace the question by another, which is closely related to it and is expressed in relatively unambiguous words.

“**Artificial General Intelligence (AGI)**—sometimes called **human-level intelligence AI**—is a type of artificial intelligence that would match or surpass human capabilities across virtually all cognitive tasks.”

—ARTIFICIAL GENERAL INTELLIGENCE (AGI), WIKIPEDIA

BINGO CARD

AI

LLM

AGENTS

PROMPT
ENGINEERING

CONTEXT
ENGINEERING

VECTOR
SEARCH

RAG

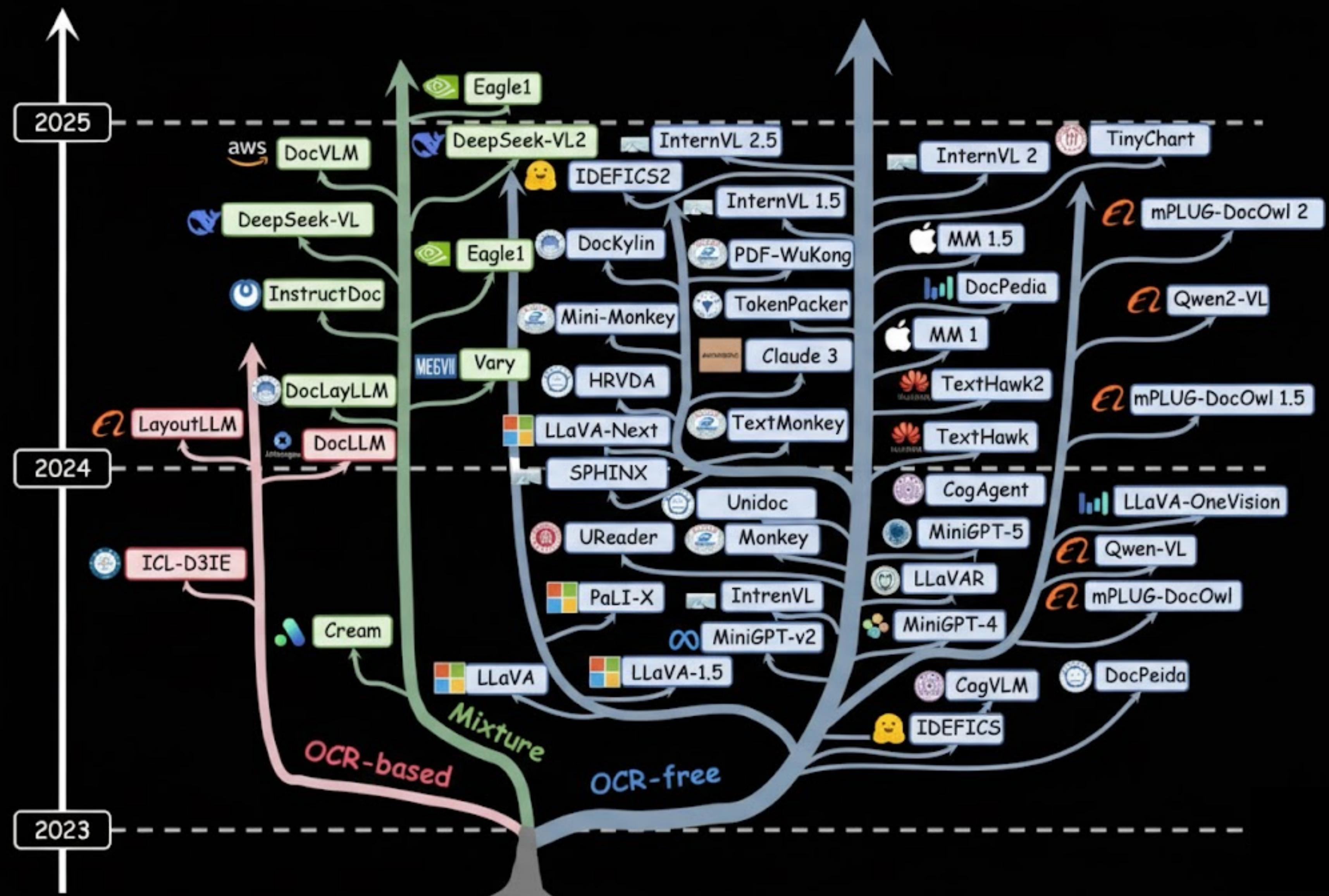
A2A

MCP



A photograph of a person's hand reaching into a metal lottery cage. The cage is filled with numerous small, colorful balls, each marked with a different number. The background is slightly blurred, showing other people and what might be a festive or event setting. A large, solid blue circle is overlaid on the lower half of the image, containing the letters "LLM" in white.

LLM

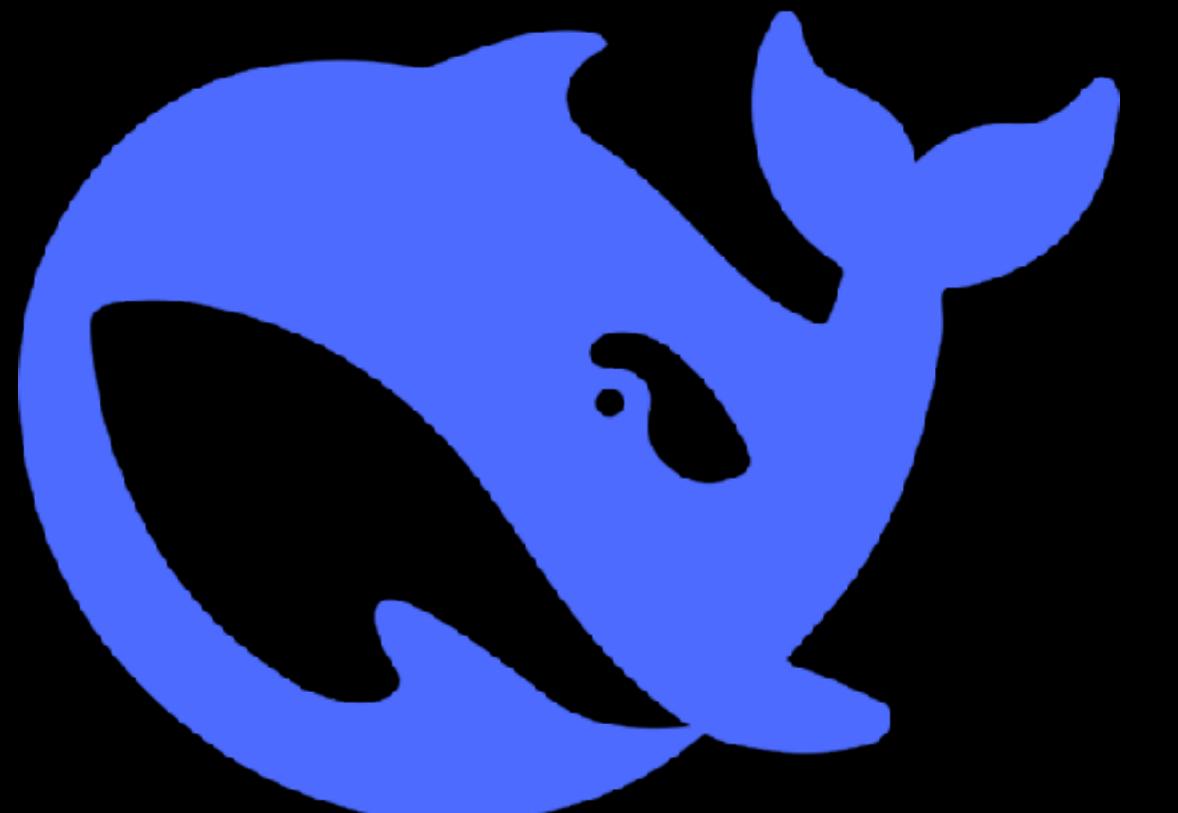
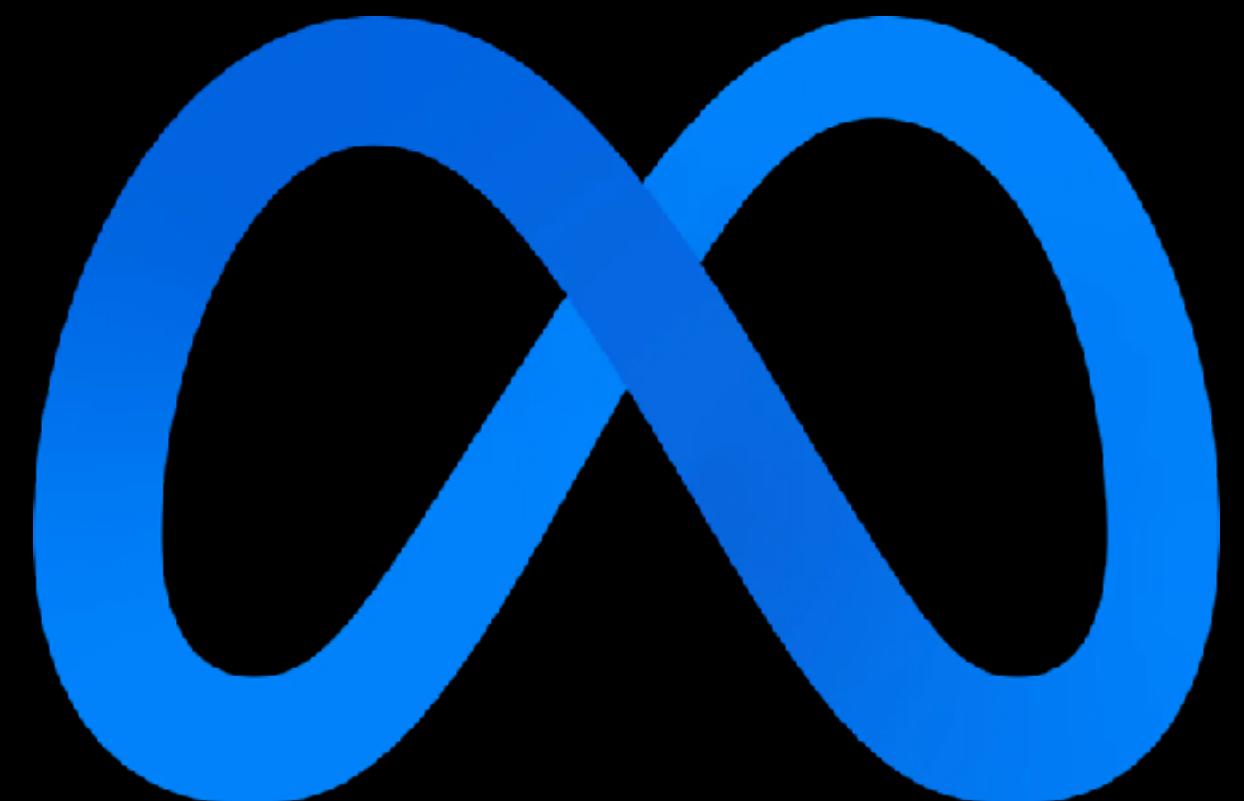


FOUNDATIONAL MODELS



Gemini

* Claude

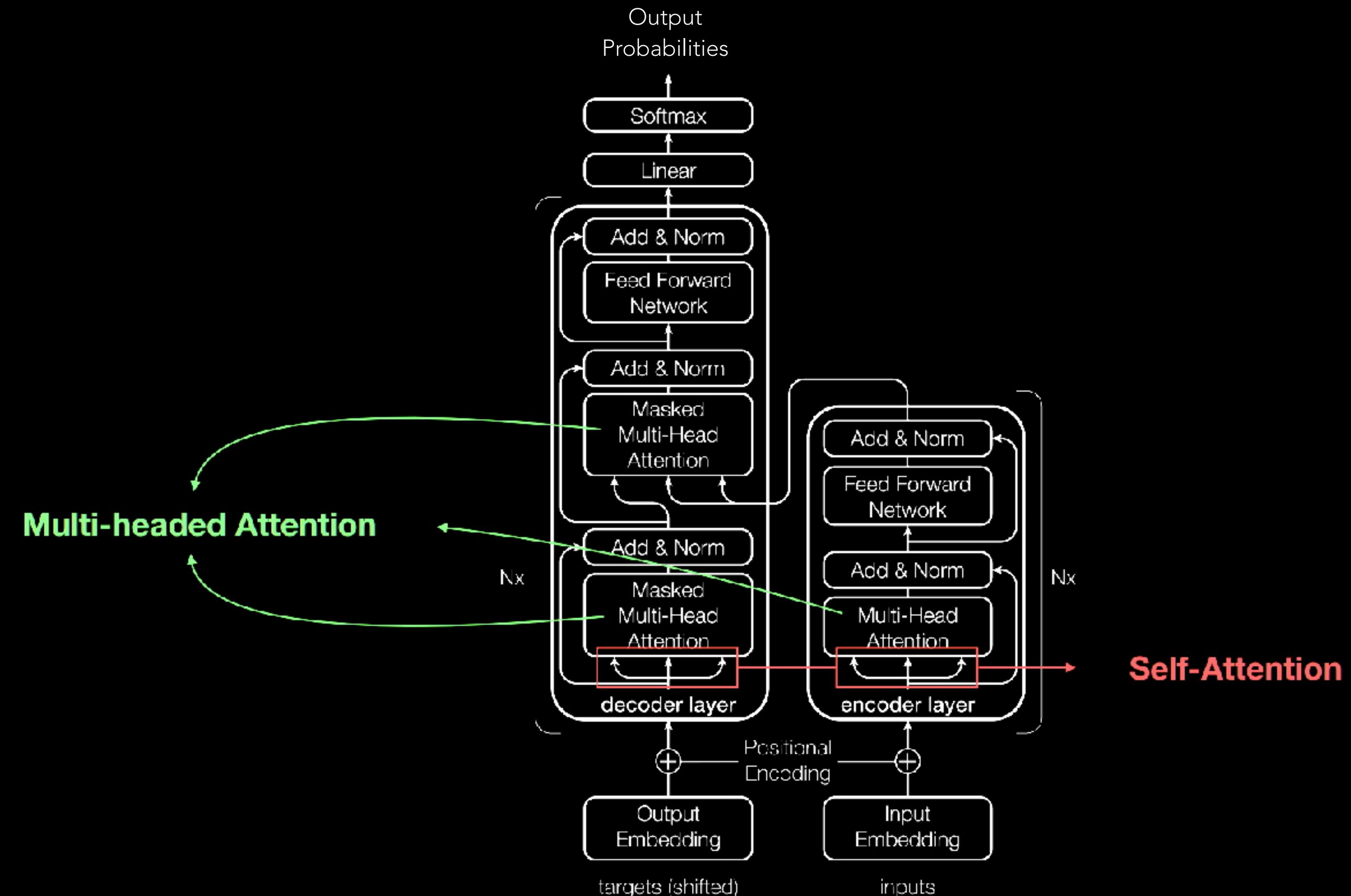


SMALL LANGUAGE MODELS

The screenshot shows a web browser window titled "small - Ollama Search" with the URL "ollama.com/search?q=small". The page displays a search bar with the query "small" and a navigation bar with links for "Models", "GitHub", "Discord", "Docs", "Cloud", "Sign in", and "Download". Below the navigation bar are filter buttons for "Cloud", "Embedding", "Vision", "Tools", "Thinking", and a dropdown menu for "Popular". The main content area lists four language models:

- smallthinker**
A new small reasoning model fine-tuned from the Qwen 2.5 3B Instruct model.
3b
81.5K Pulls, 5 Tags, Updated 8 months ago
- llama3.2**
Meta's Llama 3.2 goes small with 1B and 3B models.
tools 1b 3b
36.6M Pulls, 63 Tags, Updated 12 months ago
- mistral-small**
Mistral Small 3 sets a new benchmark in the "small" Large Language Models category below 70B.
tools 22b 24b
1.8M Pulls, 21 Tags, Updated 7 months ago
- mistral-small3.2**
An update to Mistral Small that improves on function calling, instruction following, and less repetition errors.
vision tools 24b

WHAT IS A TRANSFORMER?



"If only everything in life was as reliable as a..."

"If only everything in life was as reliable as a Volkswagen."

-VOLKSWAGEN UK ADVERTISEMENT CAMPAIGN, 1980S

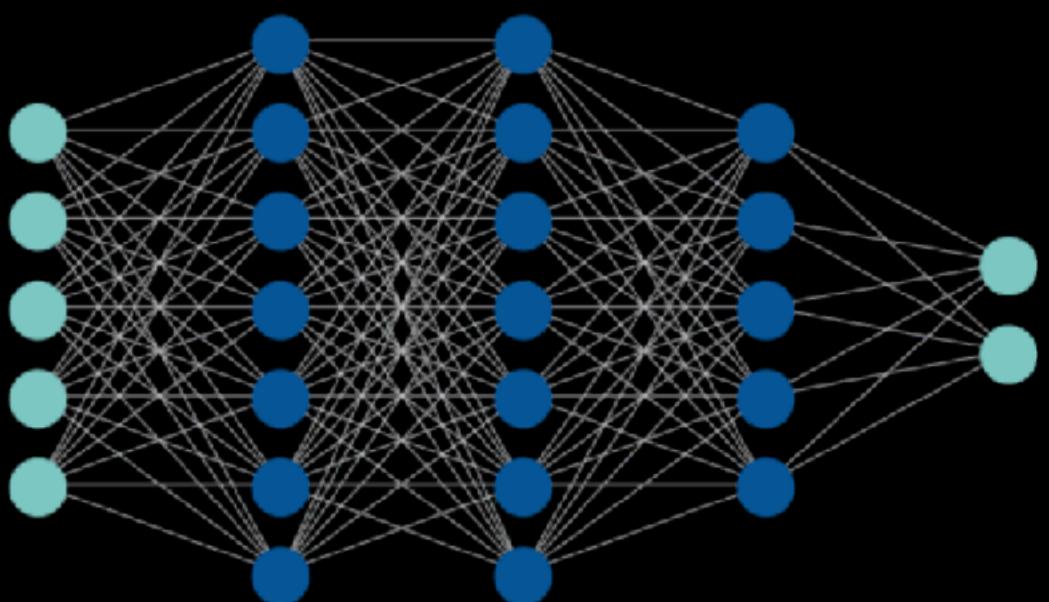


INSIDE AN LLM

If only everything in life
was as reliable as a



Input



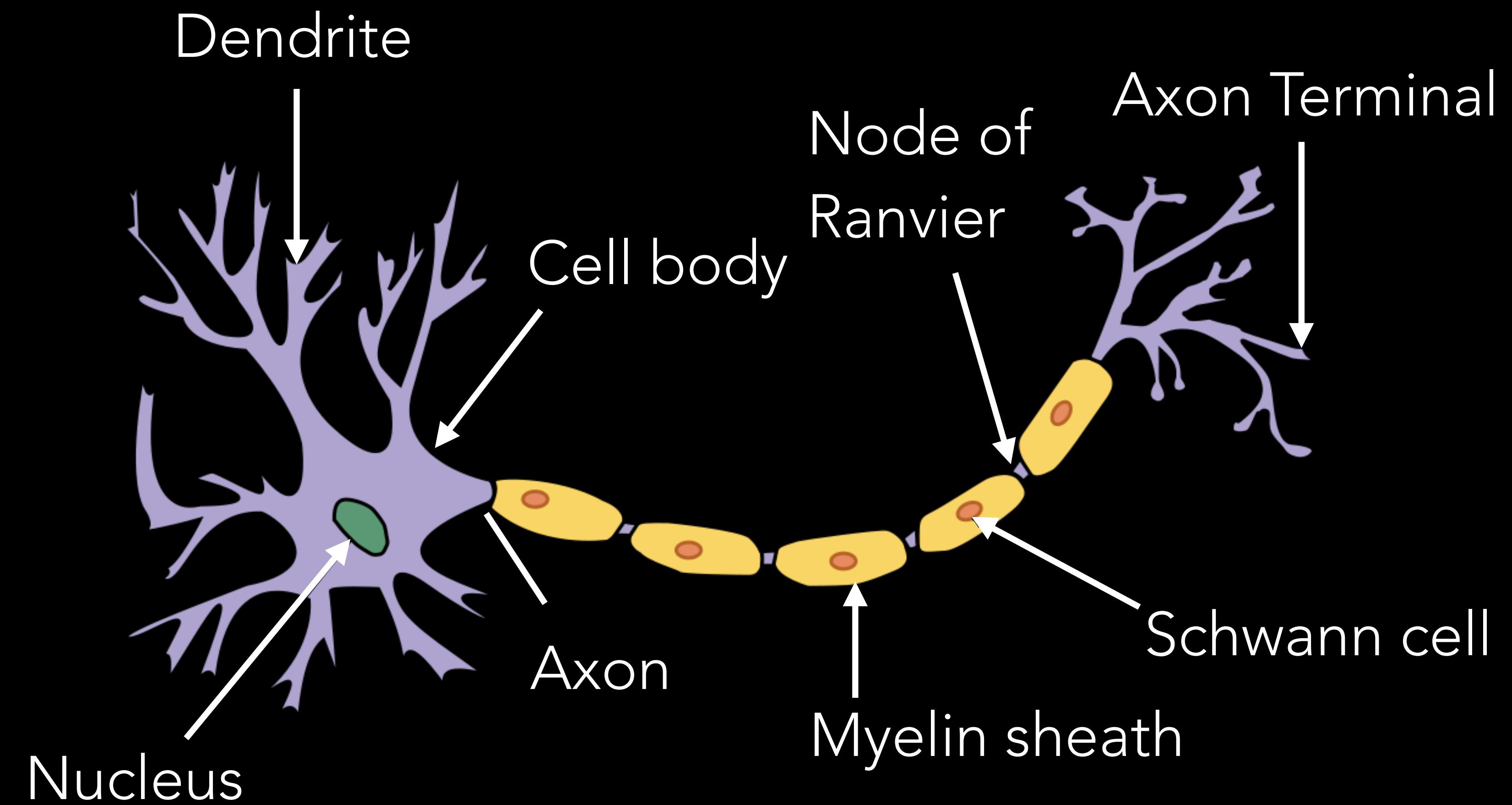
Neural Network
(LLM)



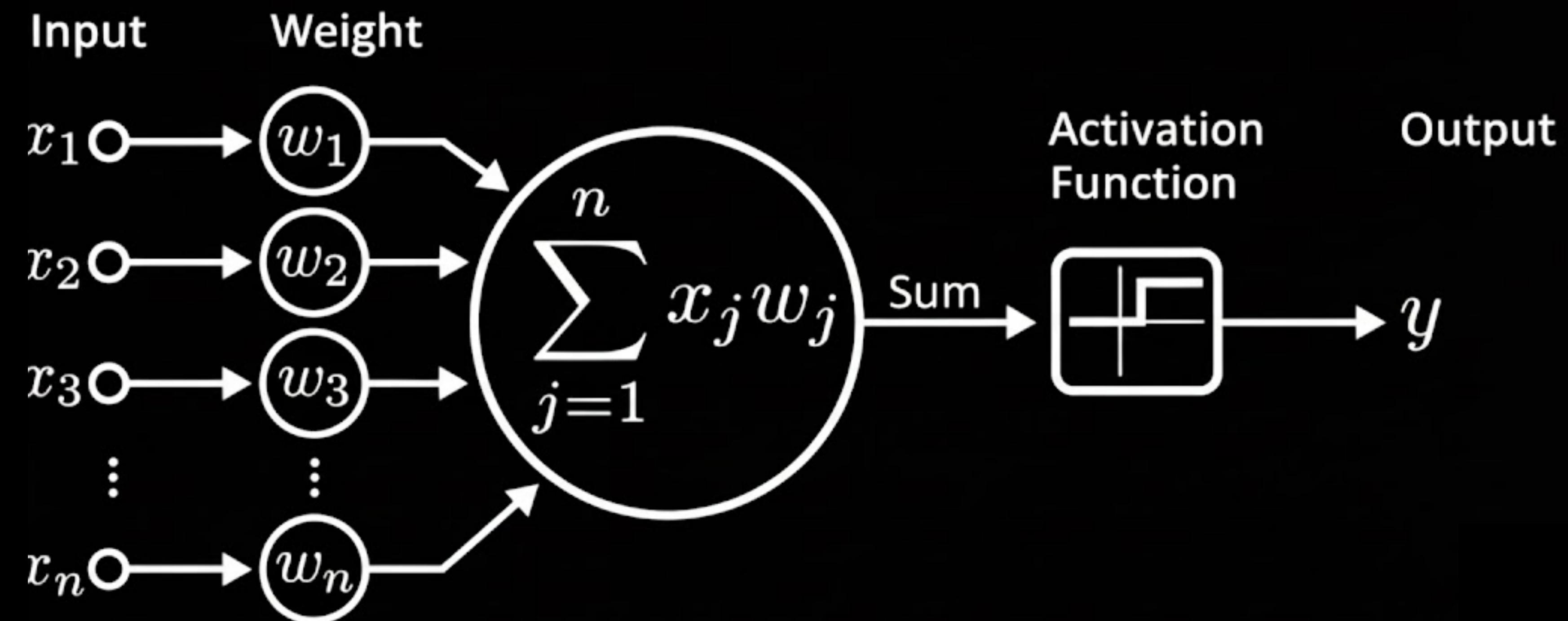
Output

WORD	PROBABILITY
0	Toyota 0.218
1	Dog 0.135
2	Volkswagen 0.892
3	Leaky Teapot 0.0001
4	BMW 0.023

WHAT IS A NEURON?



WHAT IS AN ARTIFICIAL NEURON?



PARAMETERS

Model ↑	License ↑	Parameters (B) ↑	Context ↕↑	Input \$/M ↕↑	Output \$/M ↕↑	MMLU ↕↑	HumanEval ↕↑	Multimodal ↑	Knowledge Cutoff ↑
DeepSeek-R1	Open ⓘ	671	131,072	\$0.55	\$2.19	-	-	×	-
DeepSeek-R1-0528	Open ⓘ	671	131,072	\$0.50	\$2.15	-	-	×	-
DeepSeek R1 Zero	Open ⓘ	671	0	-	-	-	-	×	-
DeepSeek-V3	Open ⓘ	671	131,072	\$0.27	\$1.10	88.5%	-	×	-
DeepSeek-V3 0324	Open ⓘ	671	163,840	\$0.28	\$1.14	-	-	×	-
DeepSeek-V3.1	Open ⓘ	671	163,840	\$0.27	\$1.00	-	-	×	-
Llama 3.1 405B Instruct	Open ⓘ	405	128,000	\$0.89	\$0.89	87.3%	89.0%	×	-
Llama 4 Maverick	Open ⓘ	400	1,000,000	\$0.17	\$0.60	85.5%	-	✓	-
DeepSeek-V2.5	Open ⓘ	236	8,192	\$0.14	\$0.28	80.4%	89.0%	×	-
Qwen3 235B A22B	Open ⓘ	235	128,000	\$0.10	\$0.10	87.8%	-	×	-
Qwen3-235B-A22B-Instruct-2507	Open ⓘ	235	131,072	\$0.15	\$0.80	-	-	×	-
Qwen3-235B-A22B-Thinking-2507	Open ⓘ	235	256,000	\$0.30	\$3.00	-	-	×	-
Pixtral Large	Open ⓘ	124	128,000	\$2.00	\$6.00	-	-	✓	-
Mistral Large 2	Open ⓘ	123	128,000	\$2.00	\$6.00	84.0%	92.0%	×	-
GPT OSS 120B	Open ⓘ	116.8	131,072	\$0.09	\$0.45	-	-	×	-
Llama 4 Scout	Open ⓘ	109	10,000,000	\$0.08	\$0.30	79.6%	-	✓	-
Llama 3.2 90B Instruct	Open ⓘ	90	128,000	\$0.35	\$0.40	86.0%	-	✓	-
Qwen3-Next-80B-A3B-Instruct	Open ⓘ	80	65,536	\$0.15	\$1.50	-	-	×	-
Qwen3-Next-80B-A3B-Thinking	Open ⓘ	80	65,536	\$0.15	\$1.50	-	-	×	-
QvQ-72B-Preview	Open ⓘ	73.4	0	-	-	-	-	✓	-
Qwen2-VL-72B-Instruct	Open ⓘ	73.4	0	-	-	-	-	✓	2023-06-30

TEMPERATURE

- A parameter that controls the randomness of the generated text ranging from 0 to 2
- Lower temperatures, like 0.2, lead to more deterministic outputs
- Higher temperatures, up to 2.0, result in more creative and random responses



BINGO CARD

AI

LLM

AGENTS

PROMPT
ENGINEERING

CONTEXT
ENGINEERING

VECTOR
SEARCH

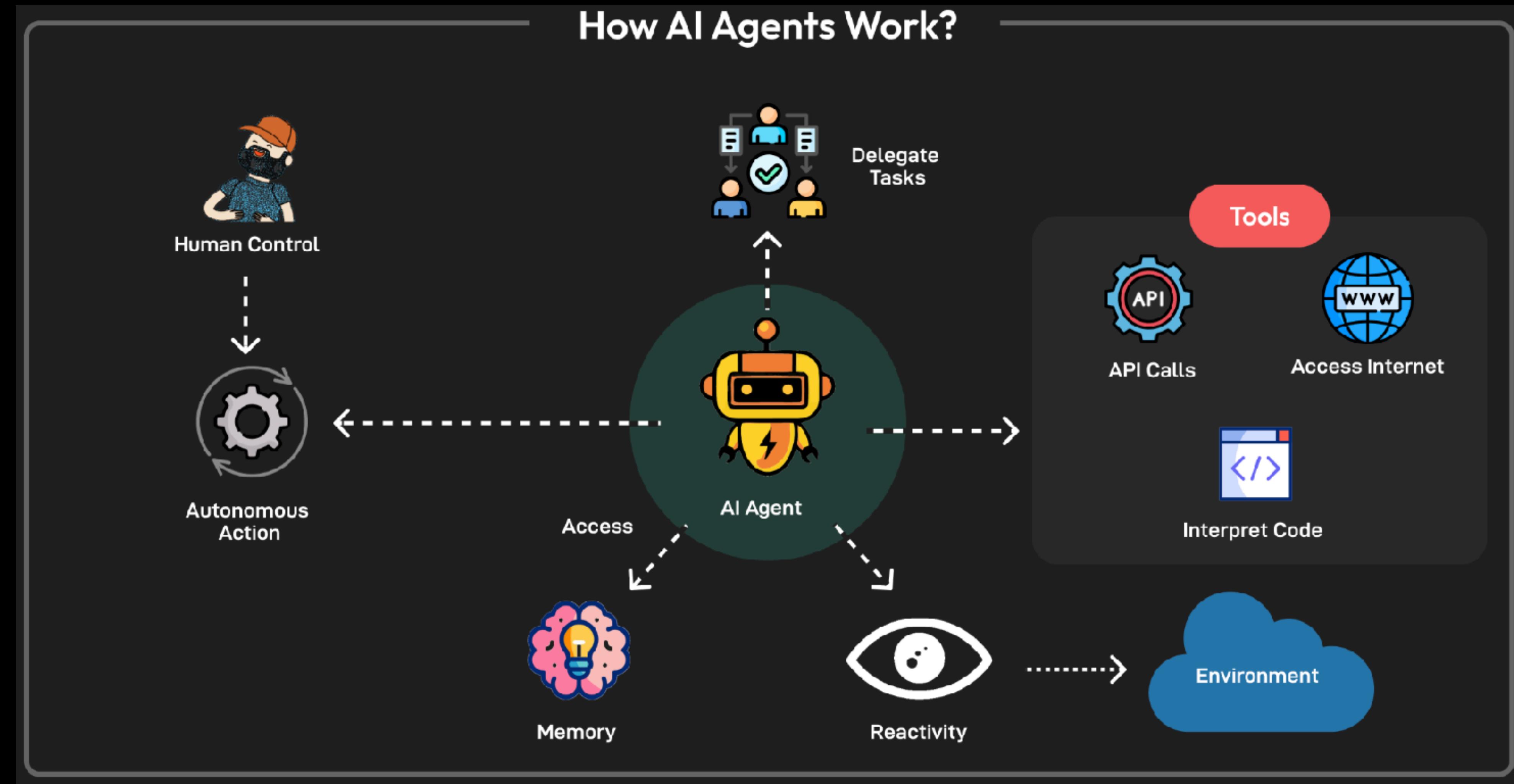
RAG

A2A

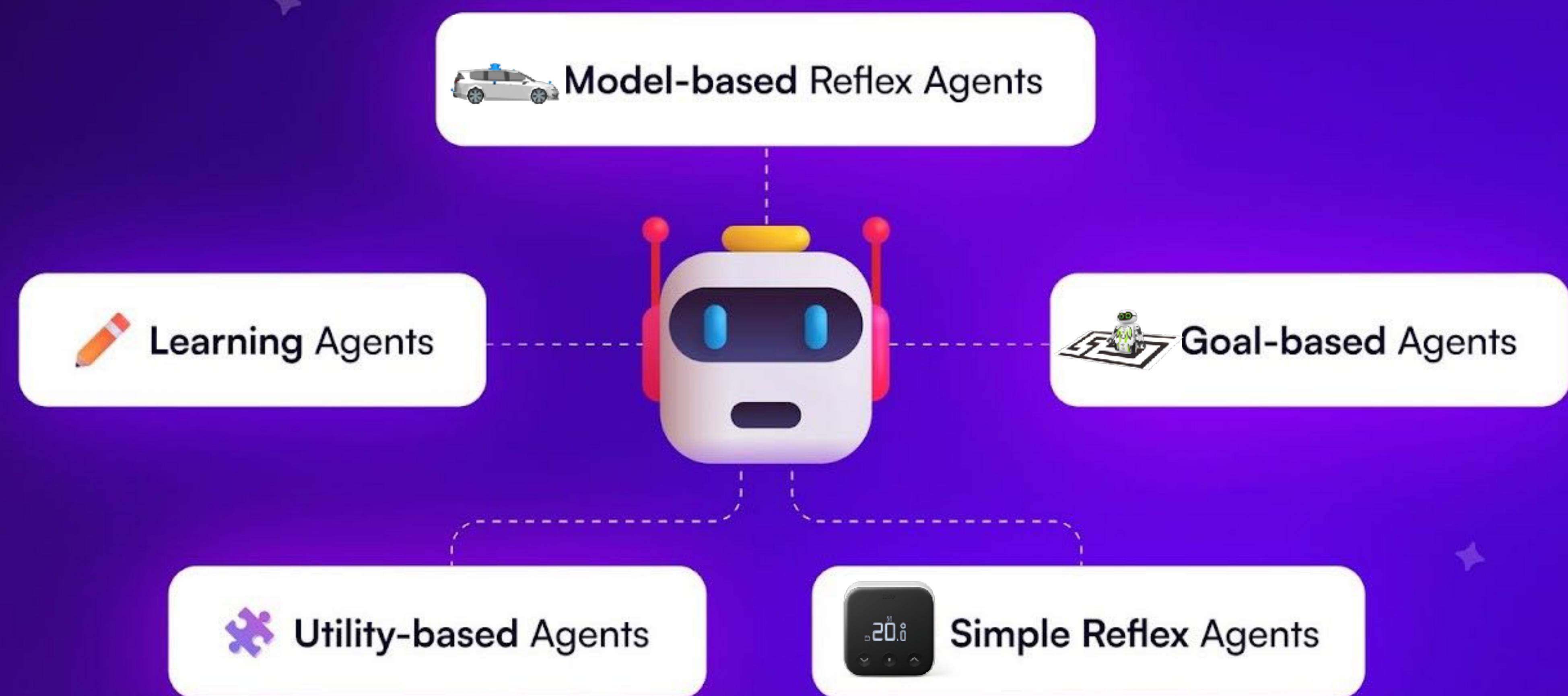
MCP

AGENTS

WHAT IS AN AGENT?



Types of AI Agents

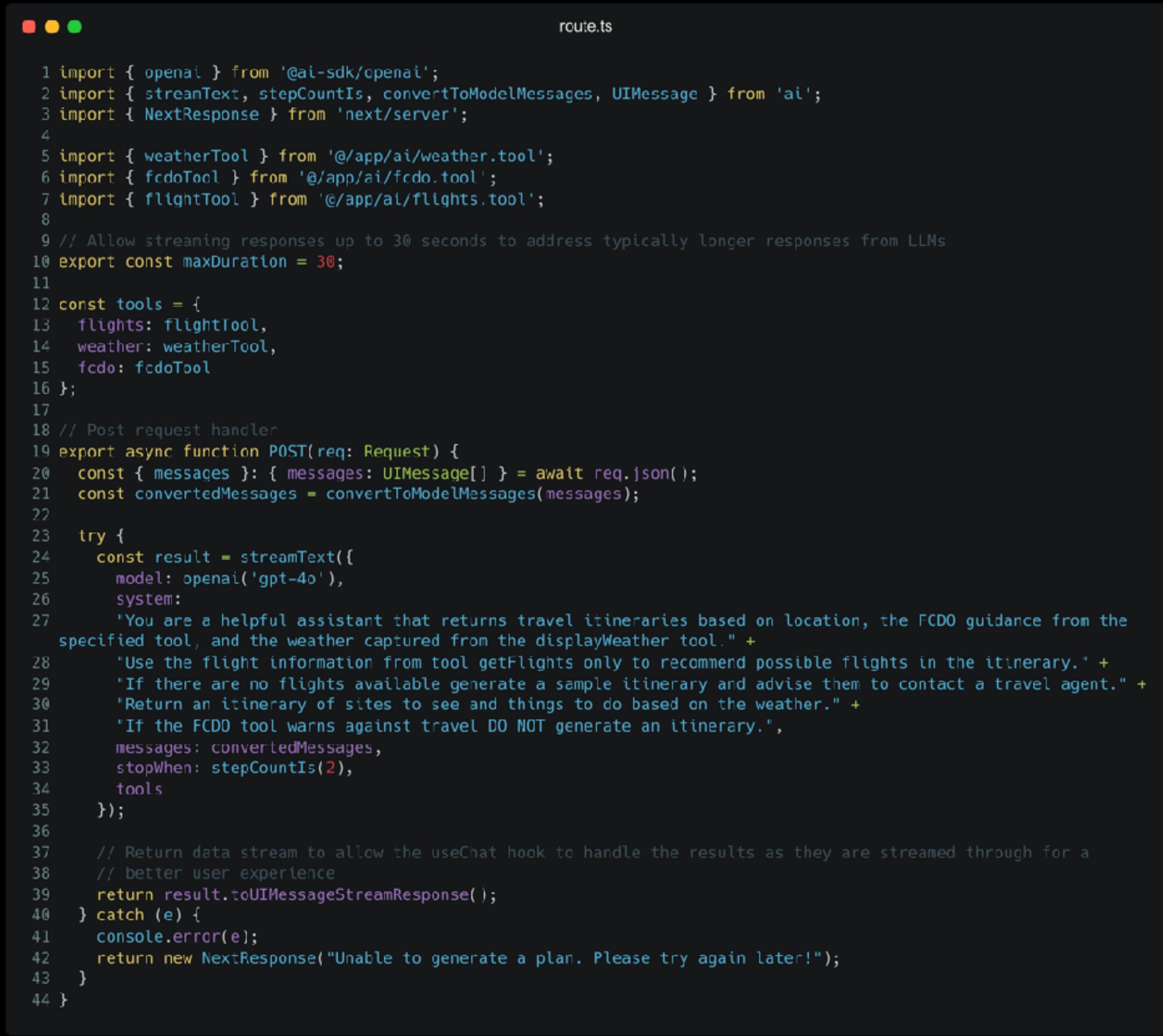


SIMPLE AGENT

 index.ts

```
1 import { Ai } from '@cloudflare/ai'
2 import { Hono } from 'hono'
3
4 export type Bindings = {
5   AI: Ai
6 }
7
8 const app = new Hono<{ Bindings: Bindings }>()
9
10 app.post('/ml/question', async (c) => {
11   const { text } = await c.req.json()
12   if (!text) {
13     return new Response('No text provided', { status: 400 })
14   }
15
16   const ai = new Ai(c.env.AI)
17   const prompt: string[] = [
18     'My name is Skippr and I help with Software Engineering problems.',
19     'I do not ask questions, I only respond to prompts. I never suggest follow
      up questions.',
20     'My answers are concise and to the point, I do not guess.'
21   ]
22
23   const messages = [
24     ...prompt.map((item) => ({ role: 'assistant', content: item })),
25     { role: 'user', content: text }
26   ]
27   const output = await ai.run('@cf/meta/llama-2-7b-chat-fp16', {
28     messages
29   })
30   return new Response(JSON.stringify(output))
31 })
```

EXAMPLE AGENT



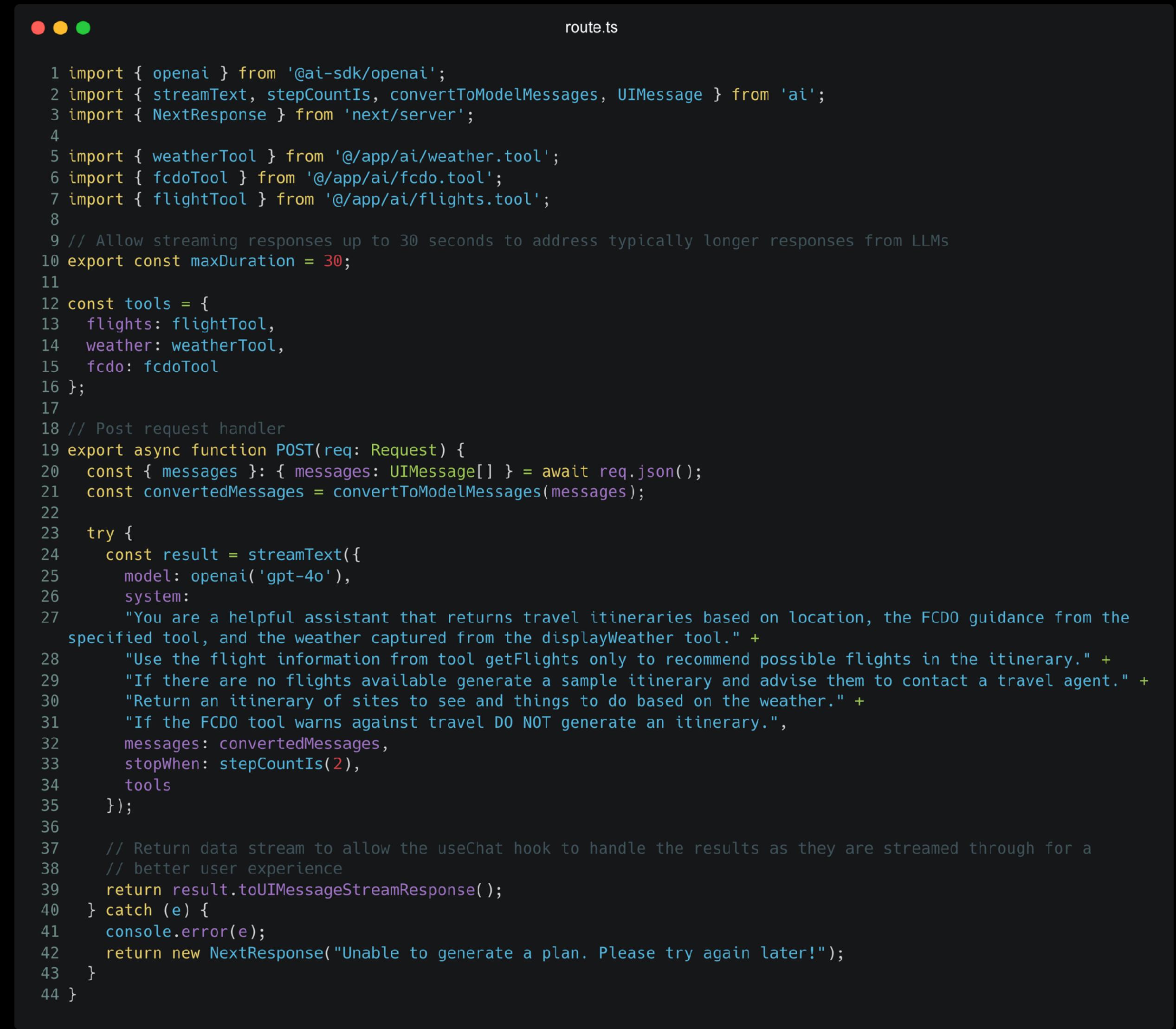
The image shows a terminal window with a dark background and light-colored text. The title bar of the terminal says "route.ts". The code itself is a TypeScript file with syntax highlighting. It imports various modules from the AI SDK and defines a POST request handler. The handler uses streamText to generate responses, specifying a model (openai('gpt-4o')), system message, and a series of steps (stopWhen: stepCountIs(2), tools). The system message includes a detailed description of the agent's capabilities and behavior. Error handling is included at the end.

```
route.ts

1 import { openai } from '@ai-sdk/openai';
2 import { streamText, stepCountIs, convertToModelMessages, UIMessage } from 'ai';
3 import { NextResponse } from 'next/server';
4
5 import { weatherTool } from '@/app/ai/weather.tool';
6 import { fcdoTool } from '@/app/ai/fcdo.tool';
7 import { flightTool } from '@/app/ai/flights.tool';
8
9 // Allow streaming responses up to 30 seconds to address typically longer responses from LLMs
10 export const maxDuration = 30;
11
12 const tools = {
13   flights: flightTool,
14   weather: weatherTool,
15   fcdo: fcdoTool
16 };
17
18 // Post request handler
19 export async function POST(req: Request) {
20   const { messages }: { messages: UIMessage[] } = await req.json();
21   const convertedMessages = convertToModelMessages(messages);
22
23   try {
24     const result = streamText({
25       model: openai('gpt-4o'),
26       system:
27         "You are a helpful assistant that returns travel itineraries based on location, the FCDO guidance from the specified tool, and the weather captured from the displayWeather tool." +
28         "Use the flight information from tool getFlights only to recommend possible flights in the itinerary." +
29         "If there are no flights available generate a sample itinerary and advise them to contact a travel agent." +
30         "Return an itinerary of sites to see and things to do based on the weather." +
31         "If the FCDO tool warns against travel DO NOT generate an itinerary.",
32       messages: convertedMessages,
33       stopWhen: stepCountIs(2),
34       tools
35     });
36
37   // Return data stream to allow the useChat hook to handle the results as they are streamed through for a
38   // better user experience
39   return result.toUIMessageStreamResponse();
40 } catch (e) {
41   console.error(e);
42   return new NextResponse("Unable to generate a plan. Please try again later!");
43 }
44 }
```

```
12 const tools = {
13   flights: flightTool,
14   weather: weatherTool,
15   fcdo: fcdoTool
16 };
17
18 // Post request handler
19 export async function POST(req: Request) {
20   const { messages }: { messages: UIMessage[] } = await req.json();
21   const convertedMessages = convertToModelMessages(messages);
22
23   try {
24     const result = streamText({
25       model: openai('gpt-4o'),
26       system:
27         "You are a helpful assistant that returns travel itineraries based on location, the FCDO guidance from the specified tool, and the weather captured from the displayWeather tool." +
28         "Use the flight information from tool getFlights only to recommend possible flights in the itinerary." +
29         "If there are no flights available generate a sample itinerary and advise them to contact a travel agent." +
30         "Return an itinerary of sites to see and things to do based on the weather." +
31         "If the FCDO tool warns against travel DO NOT generate an itinerary.",
32       messages: convertedMessages,
33       stopWhen: stepCountIs(2),
34       tools
35     });
36
37   // Return data stream to allow the useChat hook to handle the results as they are streamed through for a
38   // better user experience
39   return result.toUIMessageStreamResponse();
40 } catch (e) {
41   console.error(e);
```

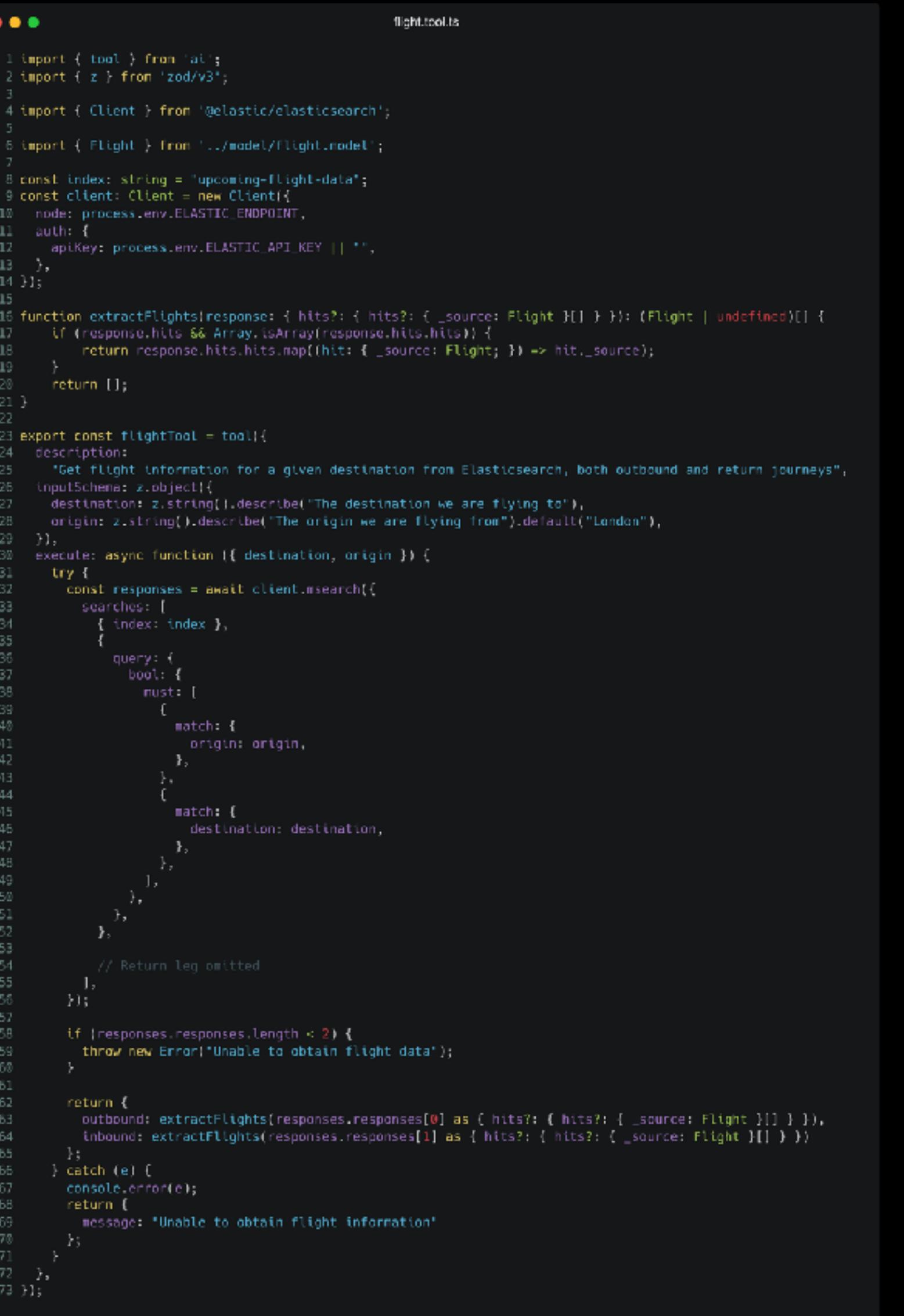
EXAMPLE SYSTEM PROMPT



```
route.ts

1 import { openai } from '@ai-sdk/openai';
2 import { streamText, stepCountIs, convertToModelMessages, UIMessage } from 'ai';
3 import { NextResponse } from 'next/server';
4
5 import { weatherTool } from '@/app/ai/weather.tool';
6 import { fcdoTool } from '@/app/ai/fcdo.tool';
7 import { flightTool } from '@/app/ai/flights.tool';
8
9 // Allow streaming responses up to 30 seconds to address typically longer responses from LLMs
10 export const maxDuration = 30;
11
12 const tools = {
13   flights: flightTool,
14   weather: weatherTool,
15   fcdo: fcdoTool
16 };
17
18 // Post request handler
19 export async function POST(req: Request) {
20   const { messages }: { messages: UIMessage[] } = await req.json();
21   const convertedMessages = convertToModelMessages(messages);
22
23   try {
24     const result = streamText({
25       model: openai('gpt-4o'),
26       system:
27         "You are a helpful assistant that returns travel itineraries based on location, the FCDO guidance from the specified tool, and the weather captured from the displayWeather tool." +
28         "Use the flight information from tool getFlights only to recommend possible flights in the itinerary." +
29         "If there are no flights available generate a sample itinerary and advise them to contact a travel agent." +
30         "Return an itinerary of sites to see and things to do based on the weather." +
31         "If the FCDO tool warns against travel DO NOT generate an itinerary.",
32       messages: convertedMessages,
33       stopWhen: stepCountIs(2),
34       tools
35     });
36
37   // Return data stream to allow the useChat hook to handle the results as they are streamed through for a
38   // better user experience
39   return result.toUIMessageStreamResponse();
40 } catch (e) {
41   console.error(e);
42   return new NextResponse("Unable to generate a plan. Please try again later!");
43 }
44 }
```

EXAMPLE AGENT TOOL



```
flighttools
1 import { tool } from 'ai';
2 import { z } from 'zod/v3';
3
4 import { Client } from '@elastic/elasticsearch';
5
6 import { Flight } from '../model/flight.model';
7
8 const index: string = 'upcoming-flight-data';
9 const client: Client = new Client({
10   node: process.env.ELASTIC_ENDPOINT,
11   auth: {
12     apiKey: process.env.ELASTIC_API_KEY || '',
13   },
14 });
15
16 function extractFlights(response: { hits?: { hits?: { _source: Flight[] | null } } } | (Flight | undefined)[]): Flight[] {
17   if (response.hits && Array.isArray(response.hits.hits)) {
18     return response.hits.hits.map(hit: { _source: Flight }) => hit._source;
19   }
20   return [];
21 }
22
23 export const flightTool = tool({
24   description:
25     "Get flight information for a given destination from Elasticsearch, both outbound and return journeys",
26   inputSchema: z.object({
27     destination: z.string().describe("The destination we are flying to"),
28     origin: z.string().describe("The origin we are flying from").default("London"),
29   }),
30   execute: async function ({ destination, origin }) {
31     try {
32       const responses = await client.search({
33         searches: [
34           { index: index },
35           {
36             query: {
37               bool: {
38                 must: [
39                   {
40                     match: {
41                       origin: origin,
42                     ...
43                   },
44                   {
45                     match: {
46                       destination: destination,
47                     ...
48                   },
49                   ...
50                 ],
51               },
52             },
53           ],
54           // Return leg omitted
55         ],
56       });
57
58       if (!responses.responses.length < 2) {
59         throw new Error('Unable to obtain flight data');
60       }
61
62       return {
63         outbound: extractFlights(responses.responses[0] as { hits?: { hits?: { _source: Flight[] } } }),
64         inbound: extractFlights(responses.responses[1] as { hits?: { hits?: { _source: Flight[] } } })
65       };
66     } catch (e) {
67       console.error(e);
68       return {
69         message: 'Unable to obtain flight information'
70       };
71     }
72   },
73 });
74 }
```

```
23 export const flightTool = tool({
24   description:
25     "Get flight information for a given destination from Elasticsearch, both outbound and return journeys",
26   inputSchema: z.object({
27     destination: z.string().describe("The destination we are flying to"),
28     origin: z.string().describe("The origin we are flying from").default("London"),
29   }),
30   execute: async function ({ destination, origin }) {
31     try {
32       const responses = await client.msearch({
33         searches: [
34           { index: index },
35           {
36             query: {
37               bool: {
38                 must: [
39                   {
40                     match: {
41                       origin: origin,
42                     },
43                   },
44                   {
45                     match: {
46                       destination: destination,
47                     },
48                   },
49                   ],
50                 },
51               },
52             ],
53           }
54         }
55       )
56       return responses
57     } catch (error) {
58       console.error(error)
59       throw error
60     }
61   }
62 }
```

```
45     match. {
46         destination: destination,
47         },
48         ],
49         },
50         },
51         },
52         },
53         // Return leg omitted
54     ],
55 });
56 });
57
58 if (responses.responses.length < 2) {
59     throw new Error("Unable to obtain flight data");
60 }
61
62 return {
63     outbound: extractFlights(responses.responses[0] as { hits?: { hits?: { _source: Flight }[] } }),
64     inbound: extractFlights(responses.responses[1] as { hits?: { hits?: { _source: Flight }[] } })
65 };
66 } catch (e) {
67     console.error(e);
68     return {
69         message: "Unable to obtain flight information"
70     };
71 }
72 },
73});
```

EXAMPLE AGENT TOOL

```
flight.tool.ts

1 import { tool } from 'ai';
2 import { z } from 'zod/v3';
3
4 import { Client } from '@elastic/elasticsearch';
5
6 import { Flight } from '../model/flight.model';
7
8 const index: string = "upcoming-flight-data";
9 const client: Client = new Client({
10   node: process.env.ELASTIC_ENDPOINT,
11   auth: {
12     apiKey: process.env.ELASTIC_API_KEY || "",
13   },
14 });
15
16 function extractFlights(response: { hits?: { hits?: { _source: Flight }[] } }): (Flight | undefined)[] {
17   if (response.hits && Array.isArray(response.hits.hits)) {
18     return response.hits.hits.map((hit: { _source: Flight }) => hit._source);
19   }
20   return [];
21 }
22
23 export const flightTool = tool({
24   description:
25     "Get flight information for a given destination from Elasticsearch, both outbound and return journeys",
26   inputSchema: z.object({
27     destination: z.string().describe("The destination we are flying to"),
28     origin: z.string().describe("The origin we are flying from").default("London"),
29   }),
30   execute: async function ({ destination, origin }) {
31     try {
32       const responses = await client.msearch({
33         searches: [
34           { index: index },
35           {
36             query: {
37               bool: {
38                 must: [
39                   {
40                     match: {
41                       origin: origin,
42                     },
43                   {
44                     match: {
45                       destination: destination,
46                     },
47                   },
48                 ],
49               },
50             },
51           },
52         ],
53       }, // Return leg omitted
54     ],
55   });
56
57   if (responses.responses.length < 2) {
58     throw new Error("Unable to obtain flight data");
59   }
60
61   return {
62     outbound: extractFlights(responses.responses[0] as { hits?: { hits?: { _source: Flight }[] } }),
63     inbound: extractFlights(responses.responses[1] as { hits?: { hits?: { _source: Flight }[] } })
64   };
65 } catch (e) {
66   console.error(e);
67   return {
68     message: "Unable to obtain flight information"
69   };
70 }
71 },
72 },
73});
```

←  D / [MCP Playground](#) / Agent Chat 🔍 ? [Give feedback](#)  CR

Elasticsearch

Home

Discover

Dashboards

Agents

Build

[Index Management](#)

Playground

Relevance

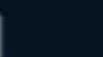
Synonyms

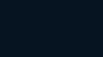
Query rules

Inference endpoints

>_ [Developer Tools](#)

⚙️ [Project settings](#)

Today 
LEGO Red Window Ide...

Last Week 
LEGO Red Window Ide...



Welcome to Elastic Agent Builder

Work interactively with your AI **agents** using the chat interface. Your selected agent answers questions by searching your data with its assigned **tools**.

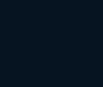
[Read the docs](#)

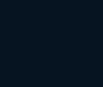
Ask anything 

LEGO Finder 

 [Create a new agent](#)
Build a custom agent tuned to your data and workflows.

 [Manage agents](#)
View, edit, and organize your existing agents.

 [Manage tools](#)
Add, remove, or edit the tools your agents can use.

 [Get started](#)
Learn how to start building agents and tools.

BINGO CARD

AI

LLM

AGENTS

PROMPT
ENGINEERING

CONTEXT
ENGINEERING

VECTOR
SEARCH

RAG

A2A

MCP

PROMPT ENGINEERING

You are a helpful AI designed to generate engaging, informative and concise descriptions for an automotive website in the UK. You will receive a list of facts about the vehicle, such as make, model, fuel type etc. You should use the features that the user provides and stay accurate to them.

Write a description for the vehicle in two coherent paragraphs with 3 or 4 sentences in each. Base the description on the following facts:

{...}

Co+Driver

WHAT IS A PROMPT?

Role

"You are a **helpful car salesperson** that recommends cars based on the attributes provided by the user." You can use the `getCars` tool to obtain the relevant car information from our inventory.

Directive

"Return a list of cars, specifying the make, model, colour, price, location and features in a **tabular format**."

Exemplars

"For example, if you are generating a list of small family cars in the colour silver, consider returning the details of a used 2024 Honda Civic in Silver with front airbags and a full service history."

Style Instructions

"Present the car details in the following format:

```
|Year|Make|Model|Colour|Features|Location|Price|  
|2024|Honda|Civic|Silver|Front airbags|Hull|£1,295|  
|2024|Honda|Civic|Silver|Airbags, GPS|London|£1,450|  
|2024|Honda|Civic|Silver|Airbags|Glasgow|£1,250|."
```

Additional Information

"If you are unable to find a set of matching cars, simply say 'Sorry, we are unable to find relevant cars' and offer to pass them to an agent."

Output Formatting

EXAMPLE STRATEGIES

- Shot-based prompting
- Chain of Thought
- Meta Prompting

Standard Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. X

Chain-of-Thought Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

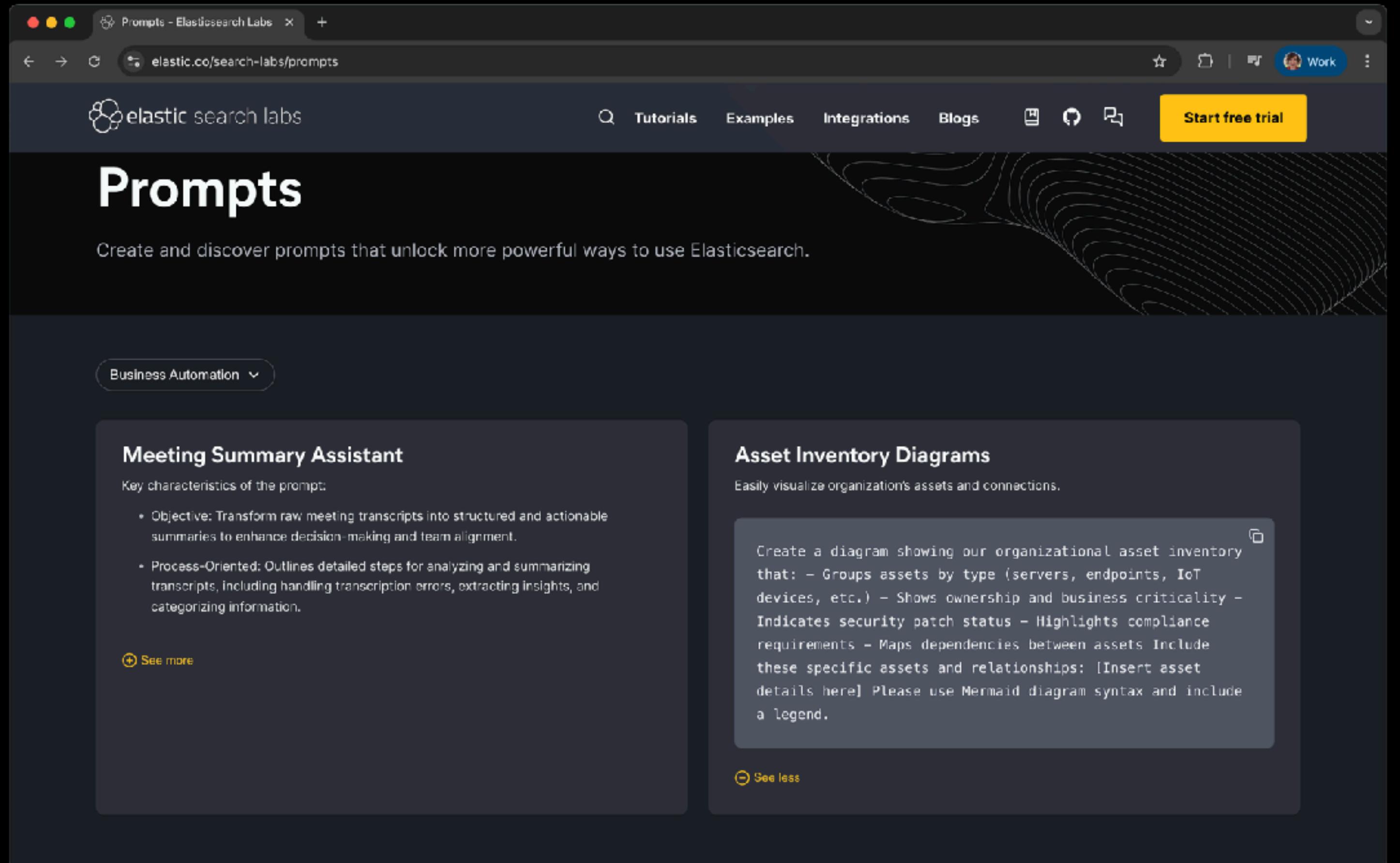
Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$. The answer is 9. ✓

HOW TO BUILD A GOOD PROMPT

- Task type
- Complexity and ambiguity
- Inputs and outputs
- Data format
- Persona to emulate
- LLM itself



BINGO CARD

AI

LLM

AGENTS

PROMPT
ENGINEERING

CONTEXT
ENGINEERING

VECTOR
SEARCH

RAG

A2A

MCP

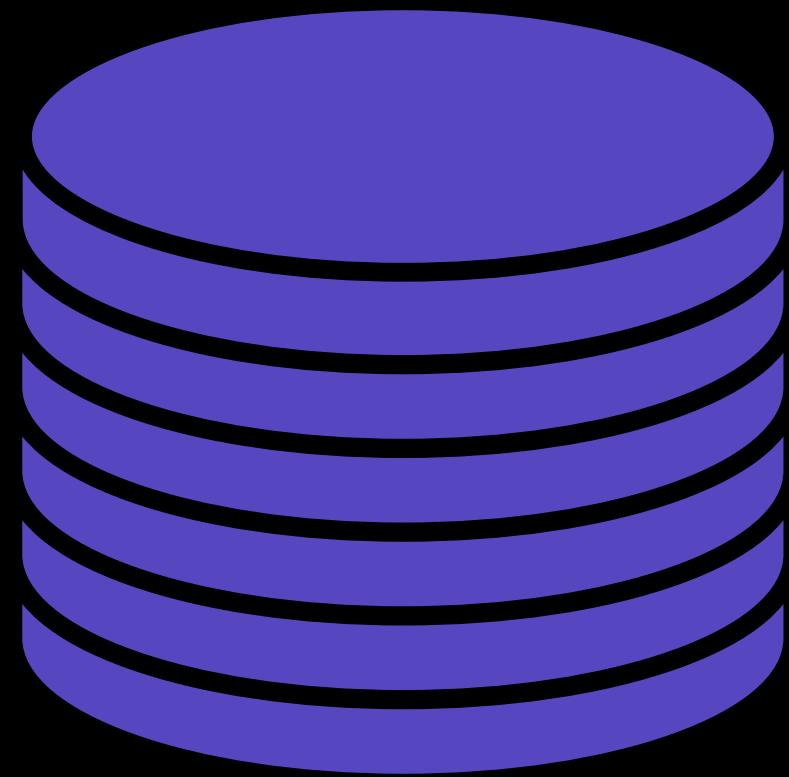
CONTEXT ENGINEERING

CONTEXT WINDOW

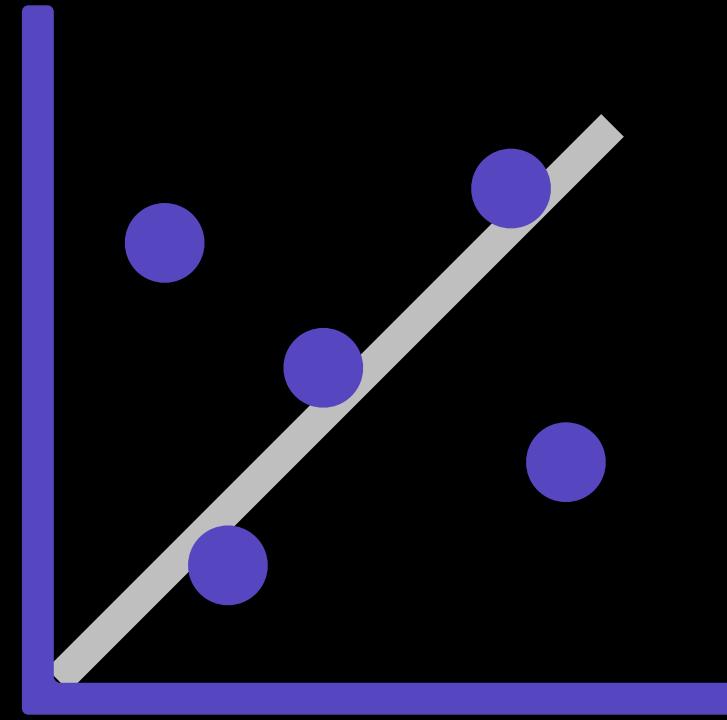
The maximum number of tokens an LLM can process at once.



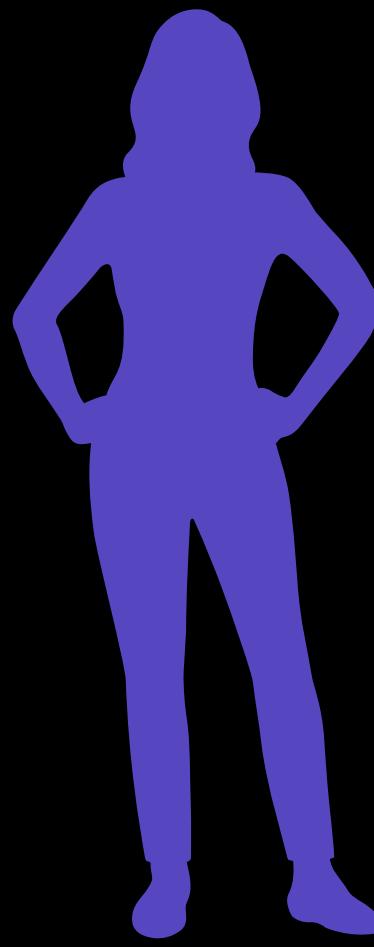
WHY DO LLMS HALLUCINATE?



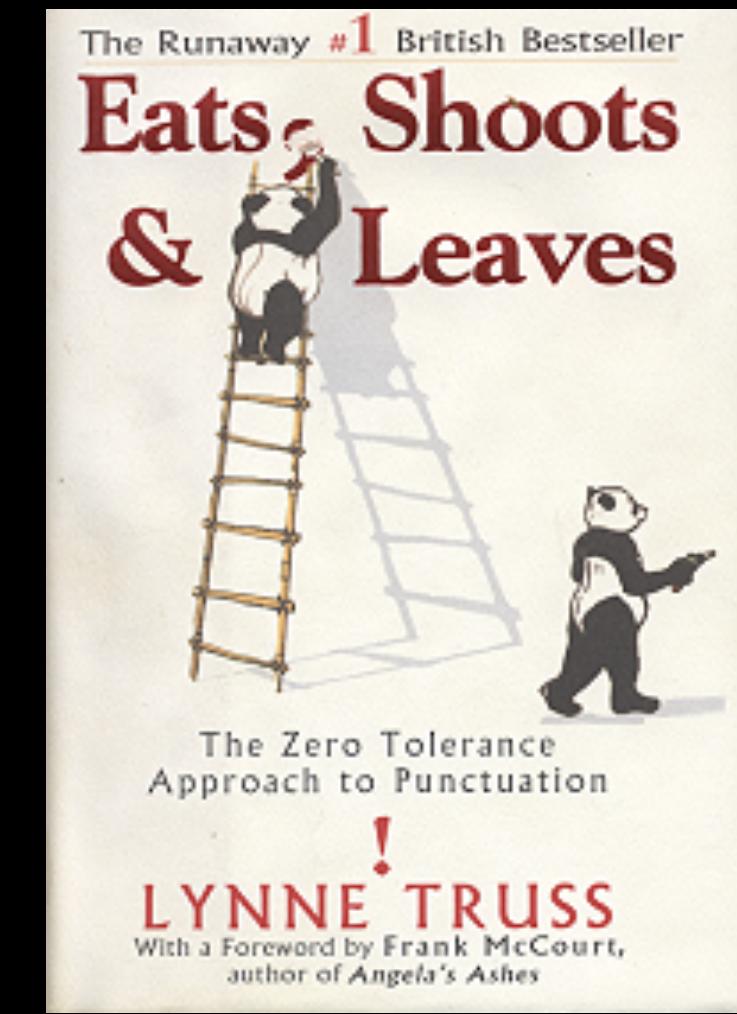
Frozen or Limited
Knowledge



Overfitting



Biases



Language
Ambiguity

Why Language Models Hallucinate

Adam Tauman Kalai*
OpenAI

Ofir Nachum
OpenAI

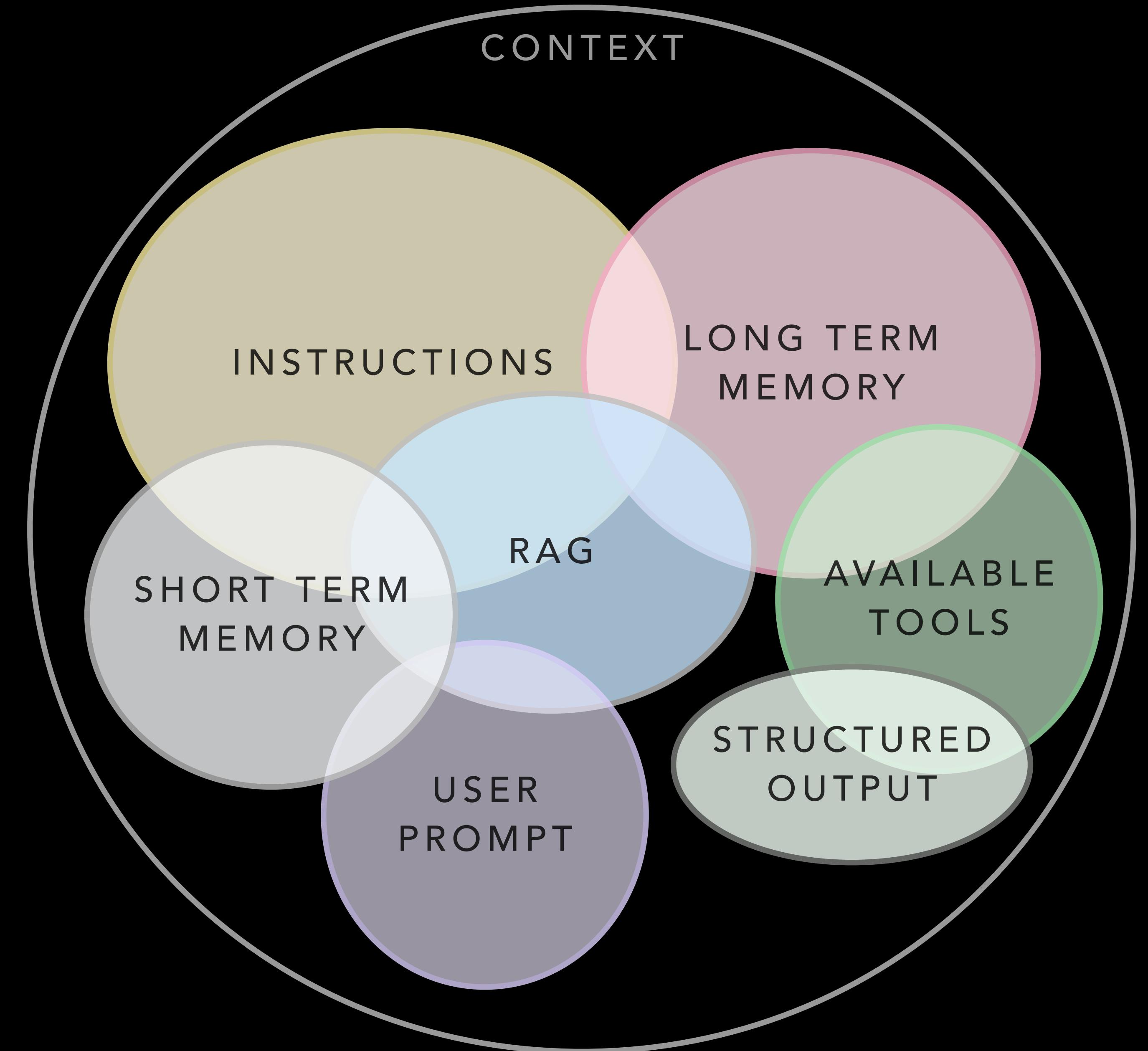
Santosh S. Vempala†
Georgia Tech

Edwin Zhang
OpenAI

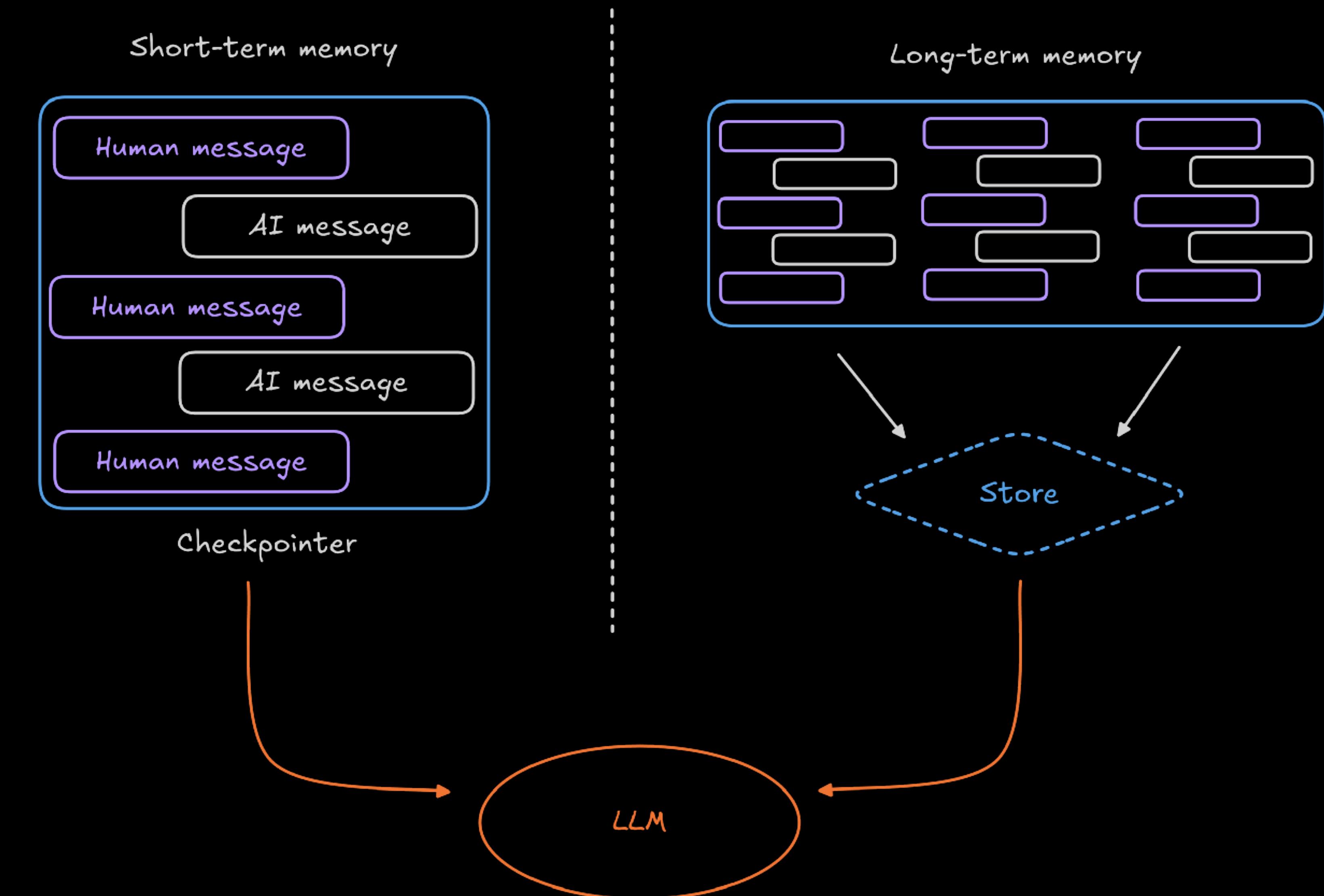
September 4, 2025

Abstract

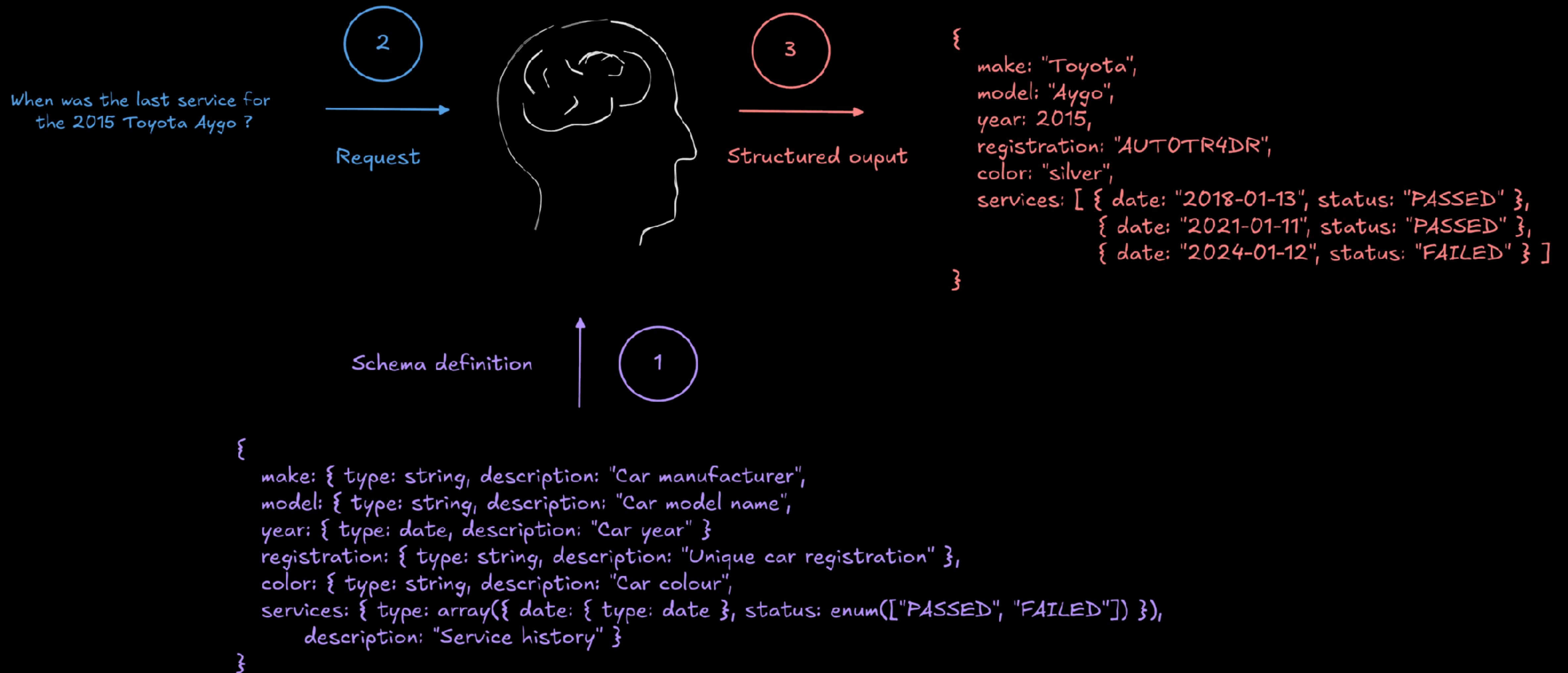
Like students facing hard exam questions, large language models sometimes guess when uncertain, producing plausible yet incorrect statements instead of admitting uncertainty. Such “hallucinations” persist even in state-of-the-art systems and undermine trust. We argue that language models hallucinate because the training and evaluation procedures reward guessing over acknowledging uncertainty and we analyze the statistical causes of hallucinations in the modern training pipeline. Hallucinations need not be mysterious—they originate simply as errors in binary classification. If incorrect statements cannot be distinguished from facts, then hallucinations in pretrained language models will arise through natural statistical pressures. We then argue that hallucinations persist due to the way most evaluations are graded—language models are optimized to be good test-takers, and guessing when uncertain improves test performance. This “epidemic” of penalizing uncertain responses can only be addressed through a socio-technical mitigation: modifying the scoring of existing benchmarks that are misaligned but dominate leaderboards, rather than introducing additional hallucination evaluations. This change may steer the field toward more trustworthy AI systems.



MEMORY



STRUCTURED OUTPUTS



EXAMPLE



skipper.schema.yaml

```
1  output:
2    required:
3      - detection
4      - impact
5      - mitigation
6      - nextSteps
7  properties:
8    detection:
9      type: string
10   description: A brief summary of how the problem was detected
11    impact:
12      type: string
13      description: The business impact, both to our teams and customers
14    mitigation:
15      type: string
16      description: Any mitigating steps that were taken in the incident
17    nextSteps:
18      type: string
19      description: The next steps, following mitigation
```

EXAMPLE TOOL OUTPUT SCHEMA

```
weather.tool.ts

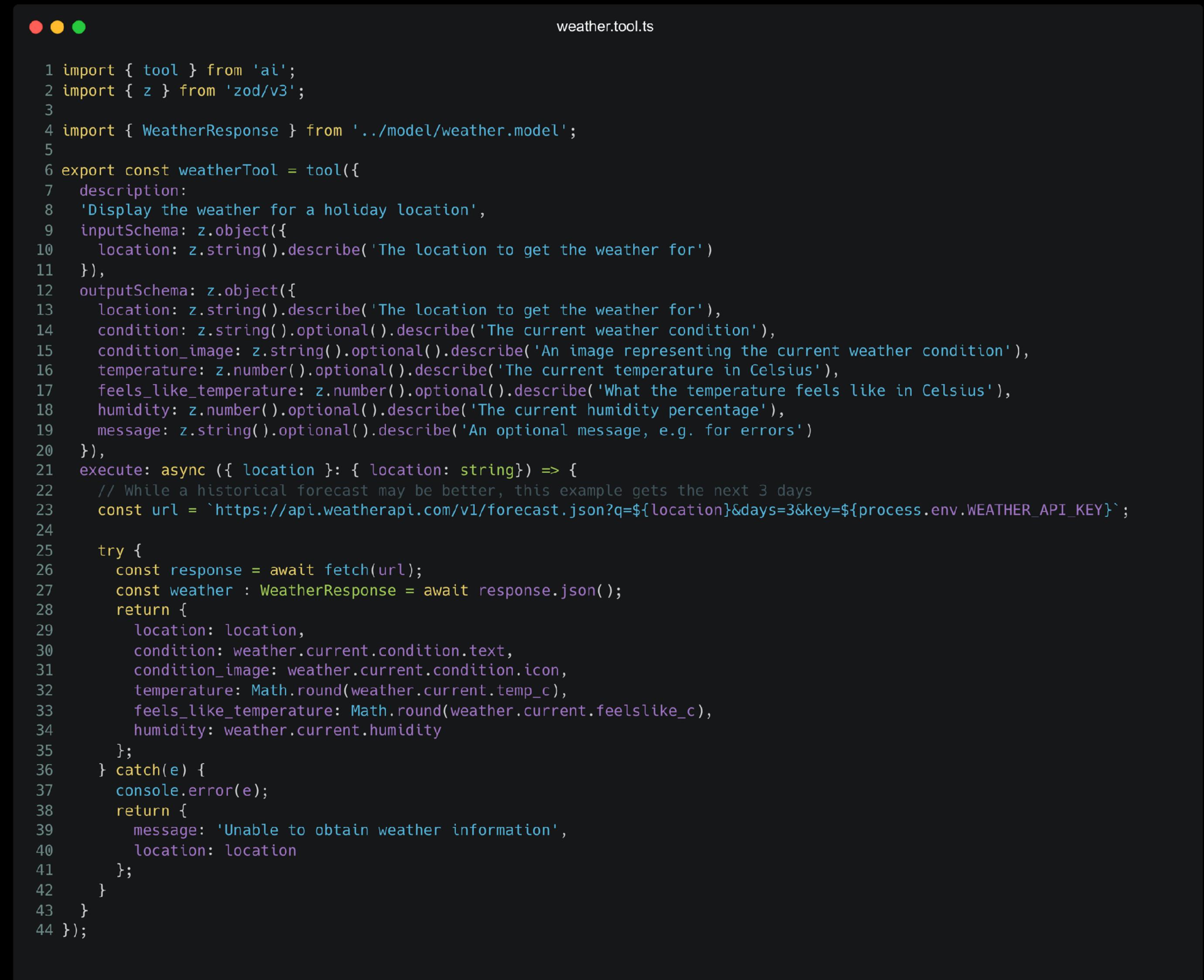
1 import { tool } from 'ai';
2 import { z } from 'zod/v3';
3
4 import { WeatherResponse } from '../model/weather.model';
5
6 export const weatherTool = tool({
7   description:
8     'Display the weather for a holiday location',
9   inputSchema: z.object({
10     location: z.string().describe('The location to get the weather for')
11   }),
12   outputSchema: z.object({
13     location: z.string().describe('The location to get the weather for'),
14     condition: z.string().optional().describe('The current weather condition'),
15     condition_image: z.string().optional().describe('An image representing the current weather condition'),
16     temperature: z.number().optional().describe('The current temperature in Celsius'),
17     feels_like_temperature: z.number().optional().describe('What the temperature feels like in Celsius'),
18     humidity: z.number().optional().describe('The current humidity percentage'),
19     message: z.string().optional().describe('An optional message, e.g. for errors')
20   },
21   execute: async ({ location }: { location: string }) => {
22     // While a historical forecast may be better, this example gets the next 3 days
23     const url = `https://api.weatherapi.com/v1/forecast.json?q=${location}&days=3&key=${process.env.WEATHER_API_KEY}`;
24
25     try {
26       const response = await fetch(url);
27       const weather: WeatherResponse = await response.json();
28       return {
29         location: location,
30         condition: weather.current.condition.text,
31         condition_image: weather.current.condition.icon,
32         temperature: Math.round(weather.current.temp_c),
33         feels_like_temperature: Math.round(weather.current.feelslike_c),
34         humidity: weather.current.humidity
35       };
36     } catch (e) {
37       console.error(e);
38       return {
39         message: 'Unable to obtain weather information',
40         location: location
41       };
42     }
43   }
44 });

```



```
1 import { tool } from 'ai';
2 import { z } from 'zod/v3';
3
4 import { WeatherResponse } from '../model/weather.model';
5
6 export const weatherTool = tool({
7   description:
8     'Display the weather for a holiday location',
9   inputSchema: z.object({
10     location: z.string().describe('The location to get the weather for')
11   }),
12   outputSchema: z.object({
13     location: z.string().describe('The location to get the weather for'),
14     condition: z.string().optional().describe('The current weather condition'),
15     condition_image: z.string().optional().describe('An image representing the current weather condition'),
16     temperature: z.number().optional().describe('The current temperature in Celsius'),
17     feels_like_temperature: z.number().optional().describe('What the temperature feels like in Celsius'),
18     humidity: z.number().optional().describe('The current humidity percentage'),
19     message: z.string().optional().describe('An optional message, e.g. for errors')
20   }),
21   execute: async ({ location }: { location: string }) => {
22     // While a historical forecast may be better, this example gets the next 3 days
23     const url = `https://api.weatherapi.com/v1/forecast.json?q=${location}&days=3&key=${process.env.WEATHER_API_KEY}`;
24
25     try {
26       const response = await fetch(url);
27       const weather: WeatherResponse = await response.json();
28       return {
29         location: location,
```

EXAMPLE TOOL OUTPUT SCHEMA



```
weather.tool.ts

1 import { tool } from 'ai';
2 import { z } from 'zod/v3';
3
4 import { WeatherResponse } from '../model/weather.model';
5
6 export const weatherTool = tool({
7   description:
8     'Display the weather for a holiday location',
9   inputSchema: z.object({
10     location: z.string().describe('The location to get the weather for')
11   }),
12   outputSchema: z.object({
13     location: z.string().describe('The location to get the weather for'),
14     condition: z.string().optional().describe('The current weather condition'),
15     condition_image: z.string().optional().describe('An image representing the current weather condition'),
16     temperature: z.number().optional().describe('The current temperature in Celsius'),
17     feels_like_temperature: z.number().optional().describe('What the temperature feels like in Celsius'),
18     humidity: z.number().optional().describe('The current humidity percentage'),
19     message: z.string().optional().describe('An optional message, e.g. for errors')
20   }),
21   execute: async ({ location }: { location: string }) => {
22     // While a historical forecast may be better, this example gets the next 3 days
23     const url = `https://api.weatherapi.com/v1/forecast.json?q=${location}&days=3&key=${process.env.WEATHER_API_KEY}`;
24
25     try {
26       const response = await fetch(url);
27       const weather: WeatherResponse = await response.json();
28       return {
29         location: location,
30         condition: weather.current.condition.text,
31         condition_image: weather.current.condition.icon,
32         temperature: Math.round(weather.current.temp_c),
33         feels_like_temperature: Math.round(weather.current.feelslike_c),
34         humidity: weather.current.humidity
35       };
36     } catch (e) {
37       console.error(e);
38       return {
39         message: 'Unable to obtain weather information',
40         location: location
41       };
42     }
43   }
44 });


```

BINGO CARD

AI

LLM

AGENTS

PROMPT
ENGINEERING

CONTEXT
ENGINEERING

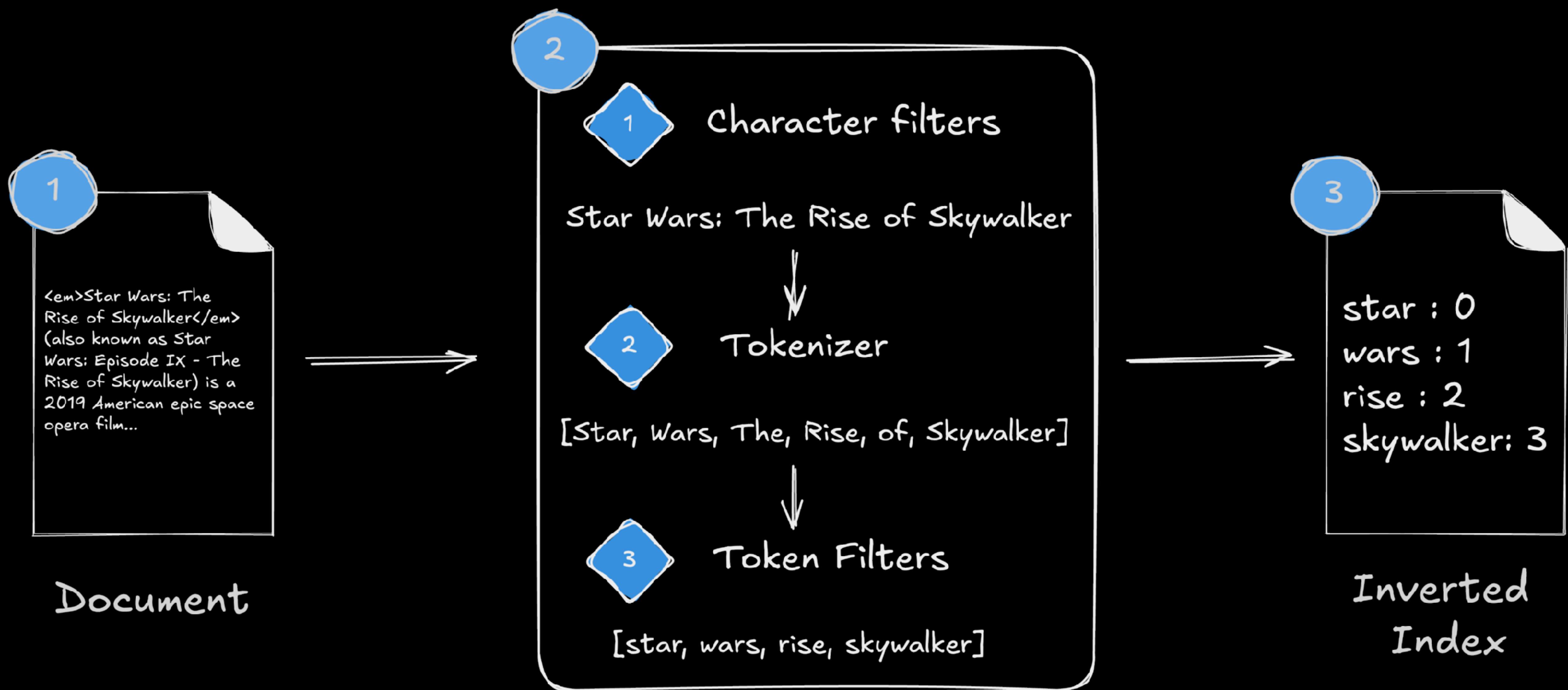
VECTOR
SEARCH

RAG

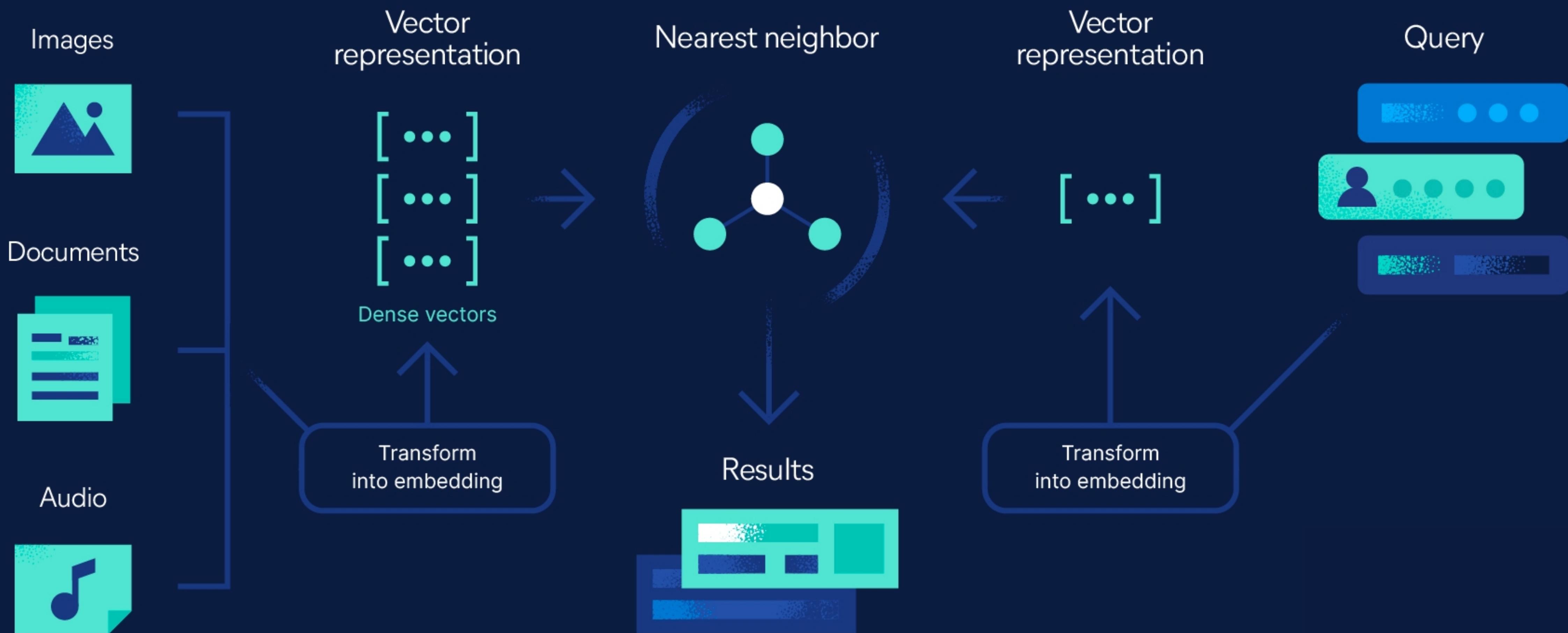
A2A

MCP

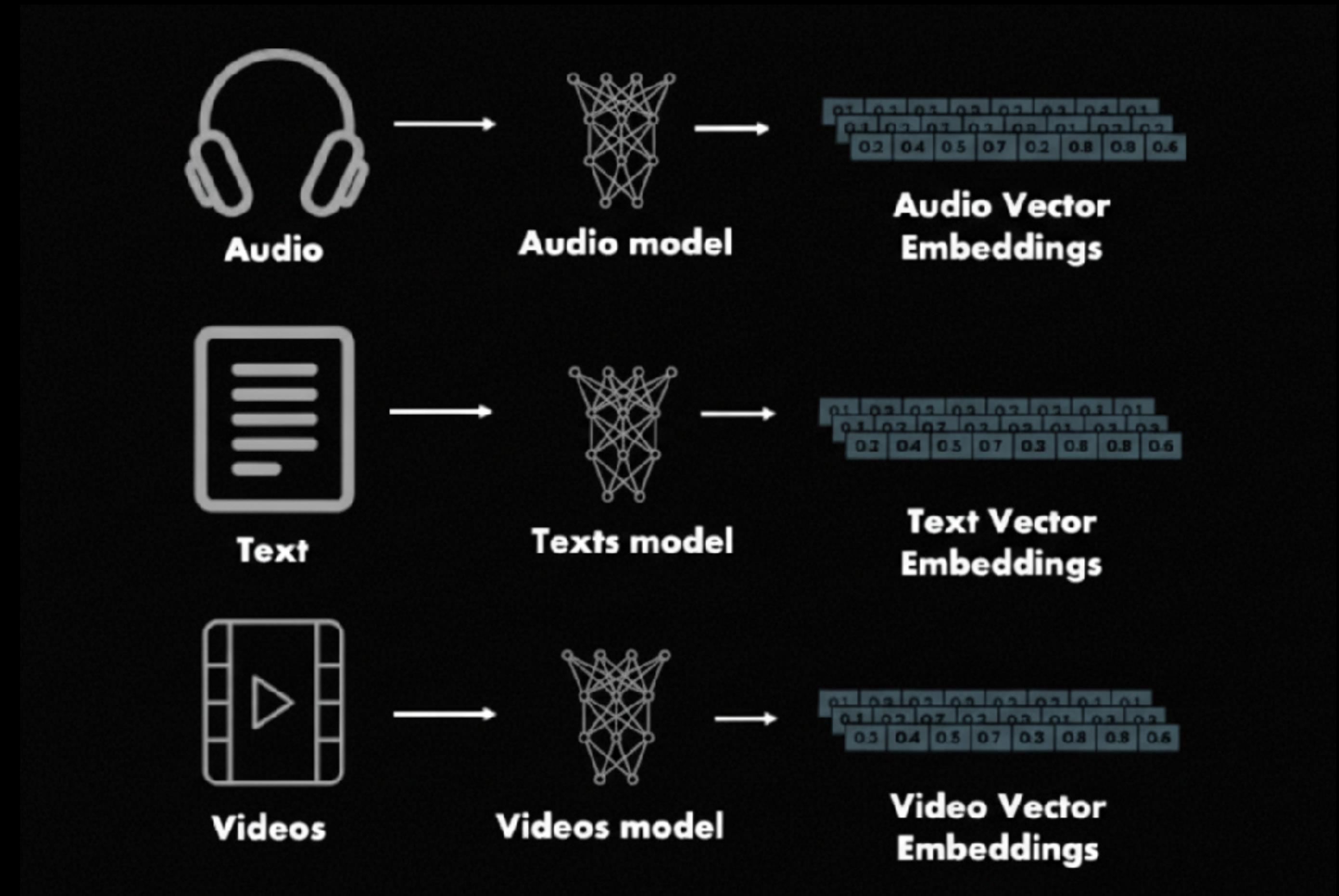
VECTOR SEARCH



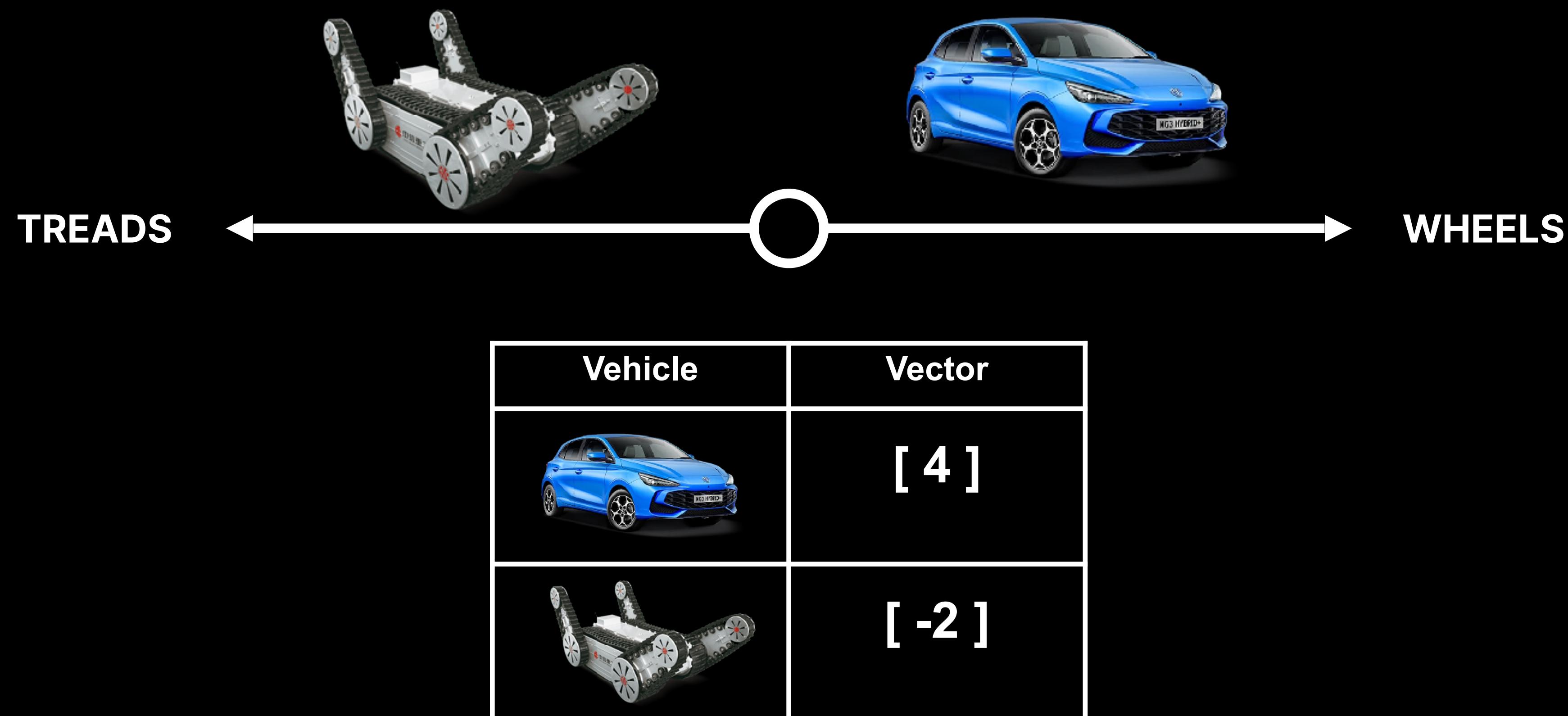
VECTOR SEARCH



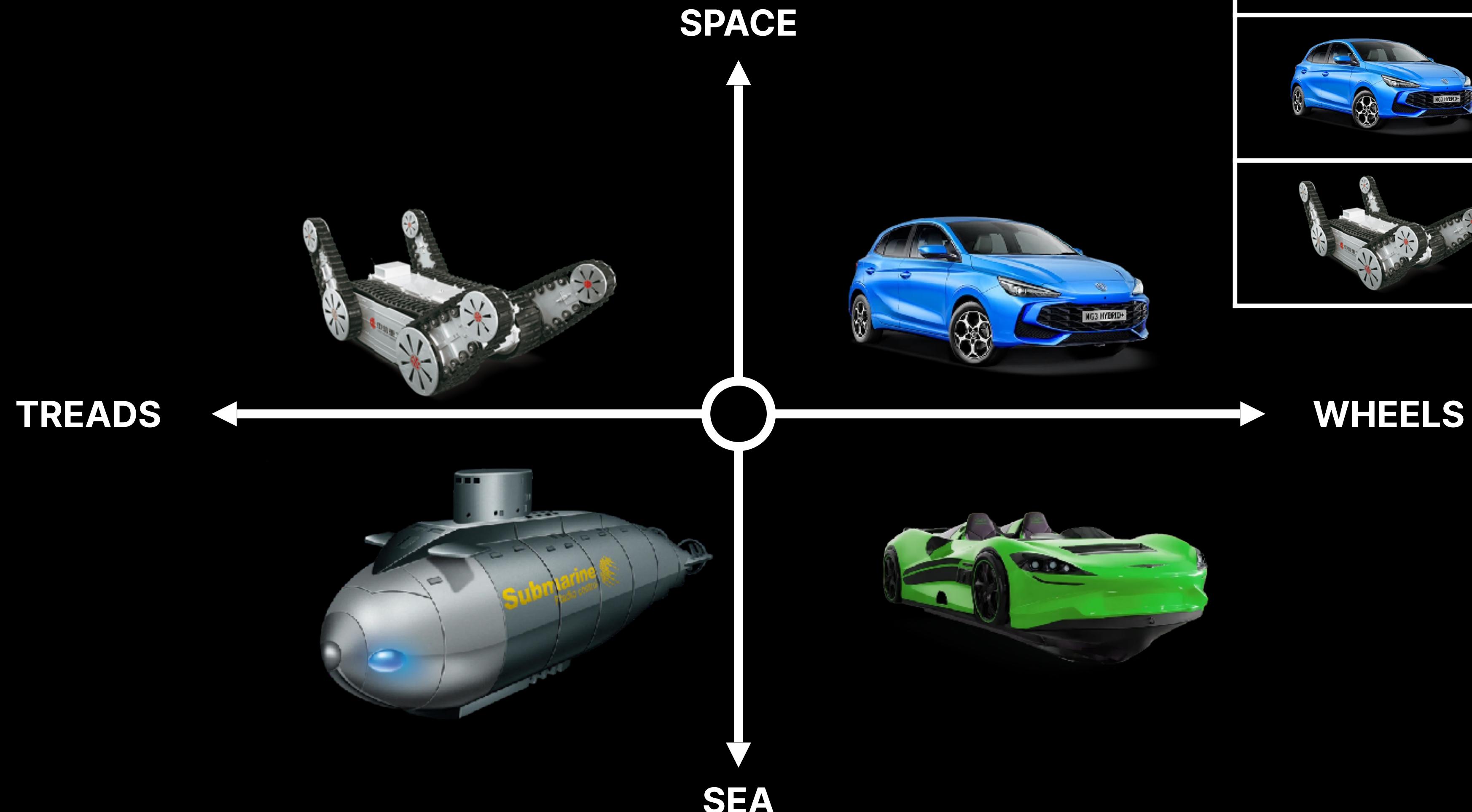
WHERE DO EMBEDDINGS COME FROM?



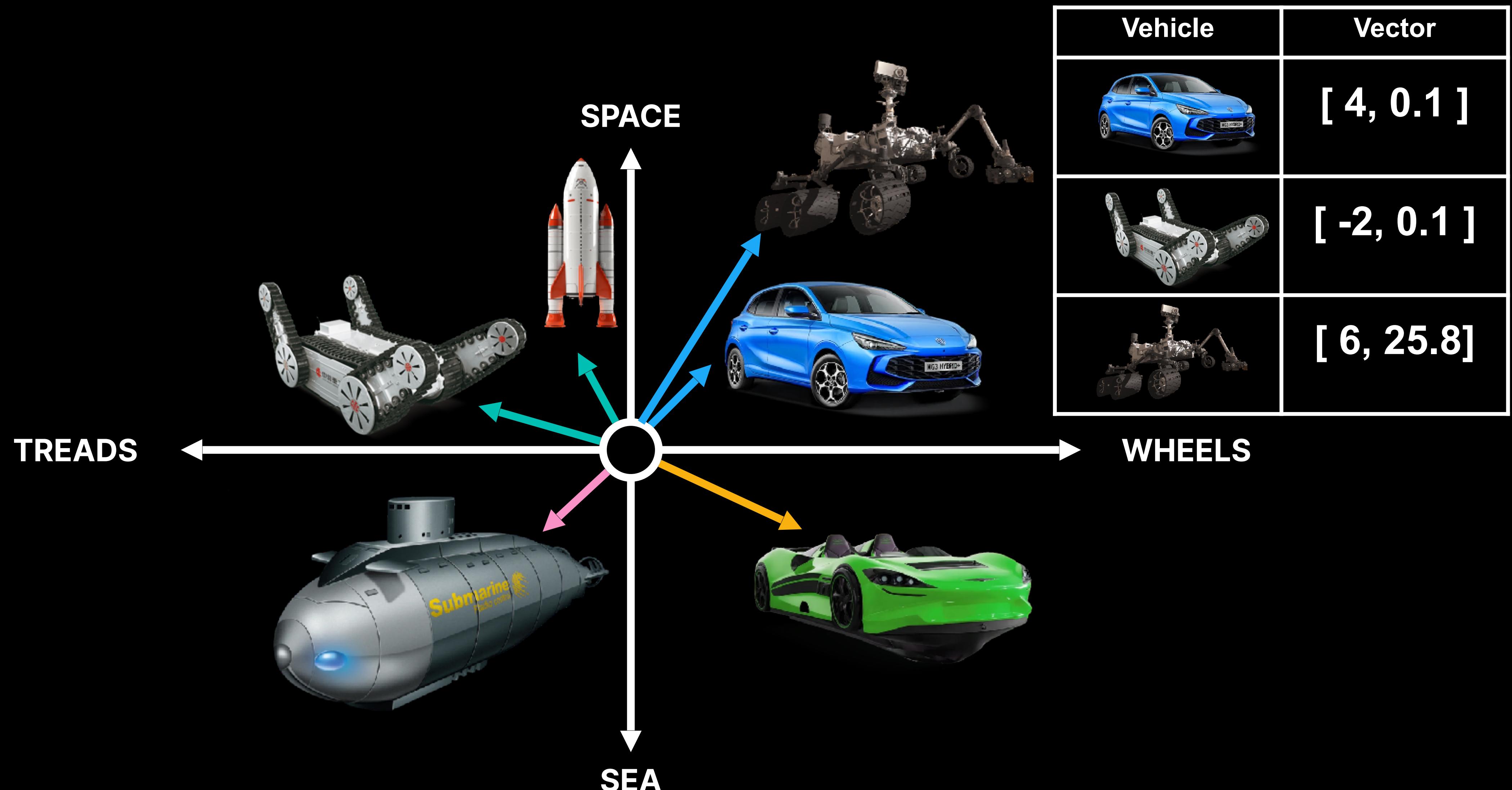
EMBEDDINGS REPRESENT YOUR DATA



MULTIPLE DIMENSIONS REPRESENT DIFFERENT DATA ASPECTS



SIMILAR DATA IS GROUPED TOGETHER

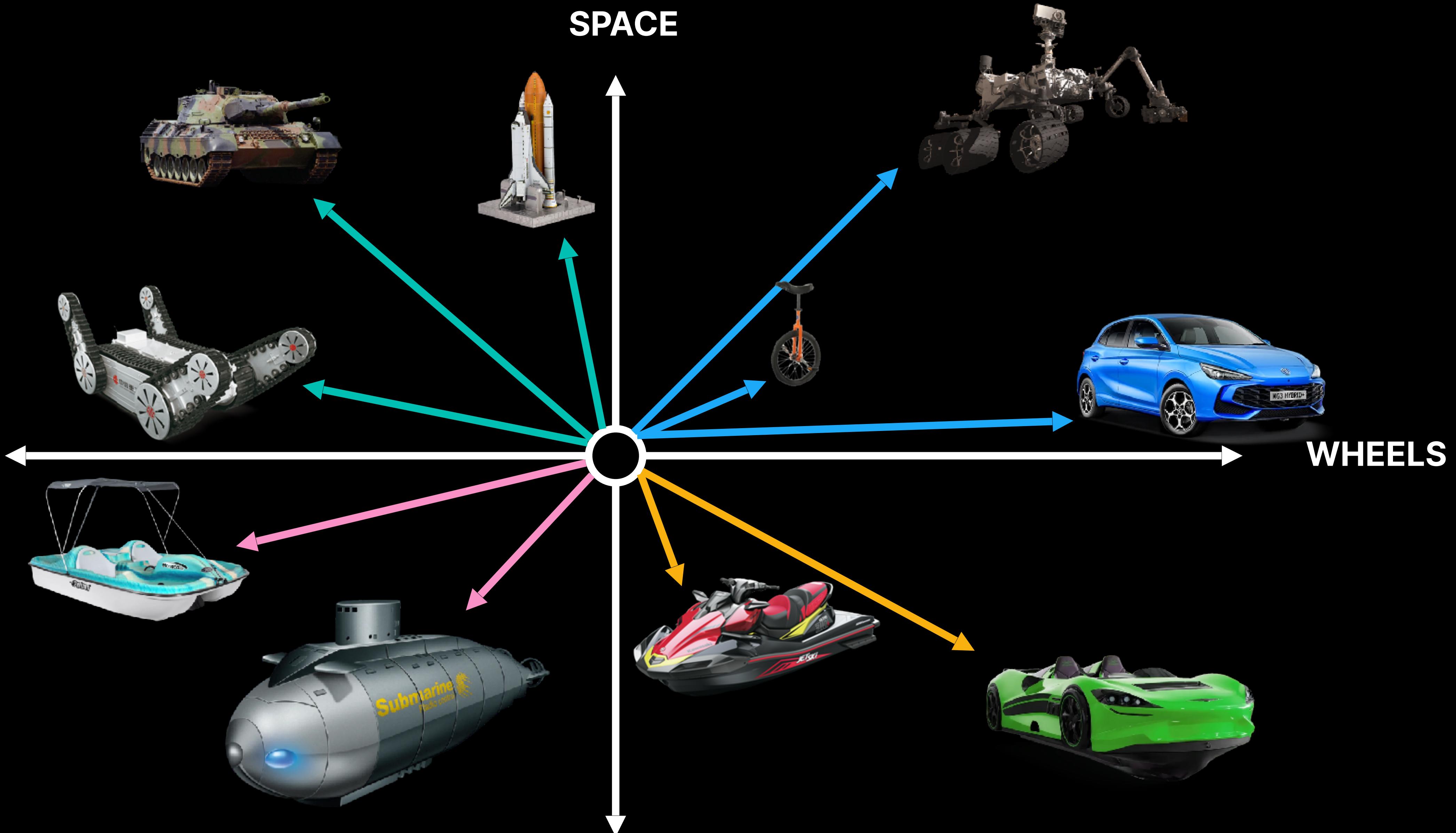


TREADS

SPACE

WHEELS

SEA



VECTOR SEARCH RANKS OBJECTS BY
SIMILARITY (RELEVANCE) TO THE QUERY



COSINE SIMILARITY

$$\theta \approx 0$$

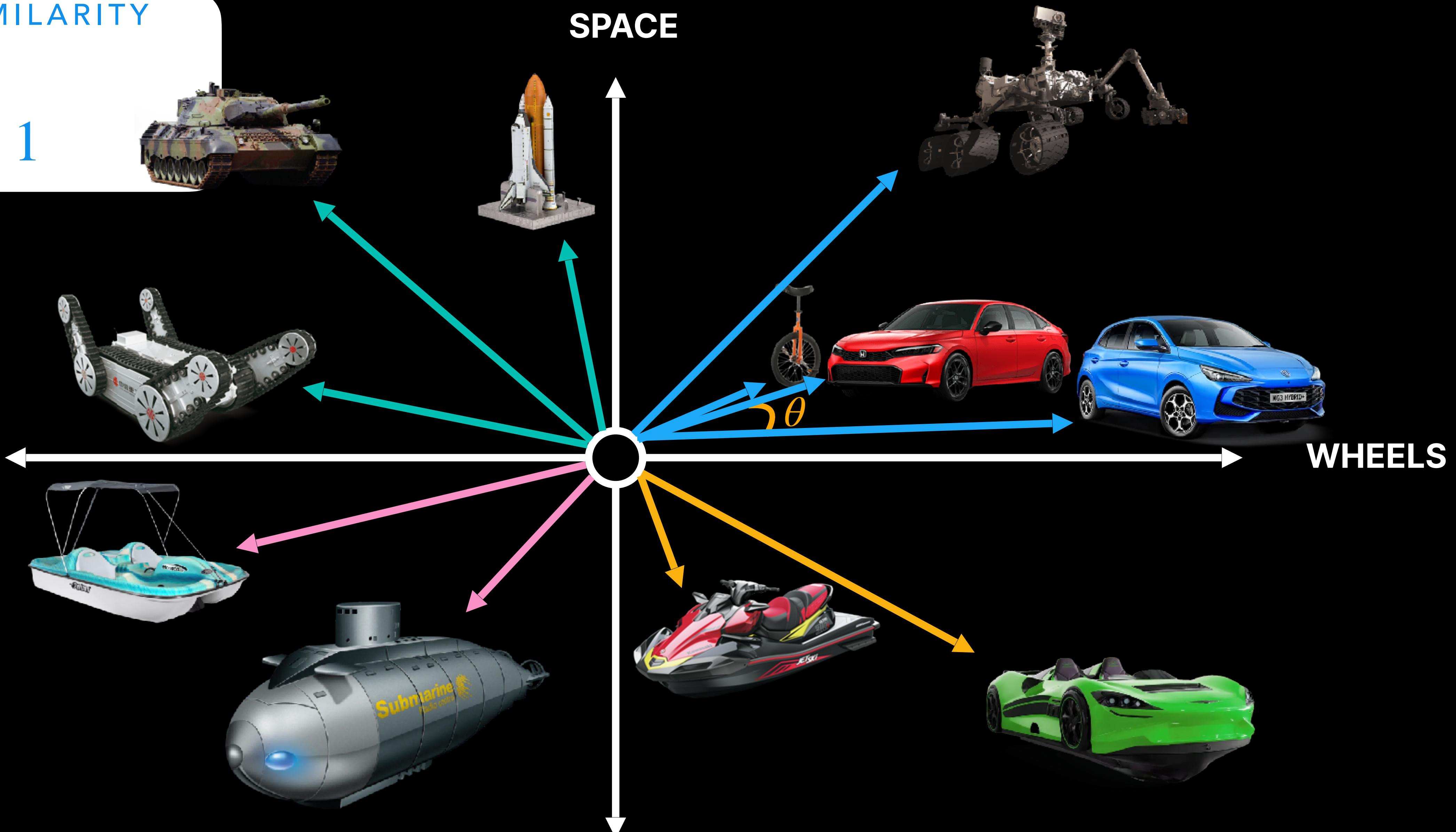
$$\cos(\theta) \approx 1$$

SPACE

TREADS

WHEELS

SEA



HOW VECTORS ARE INDEXED FOR SEARCH: GRAPHS



Hierarchical Navigable Small Worlds:
a layered approach that simplifies
access to the nearest neighbour



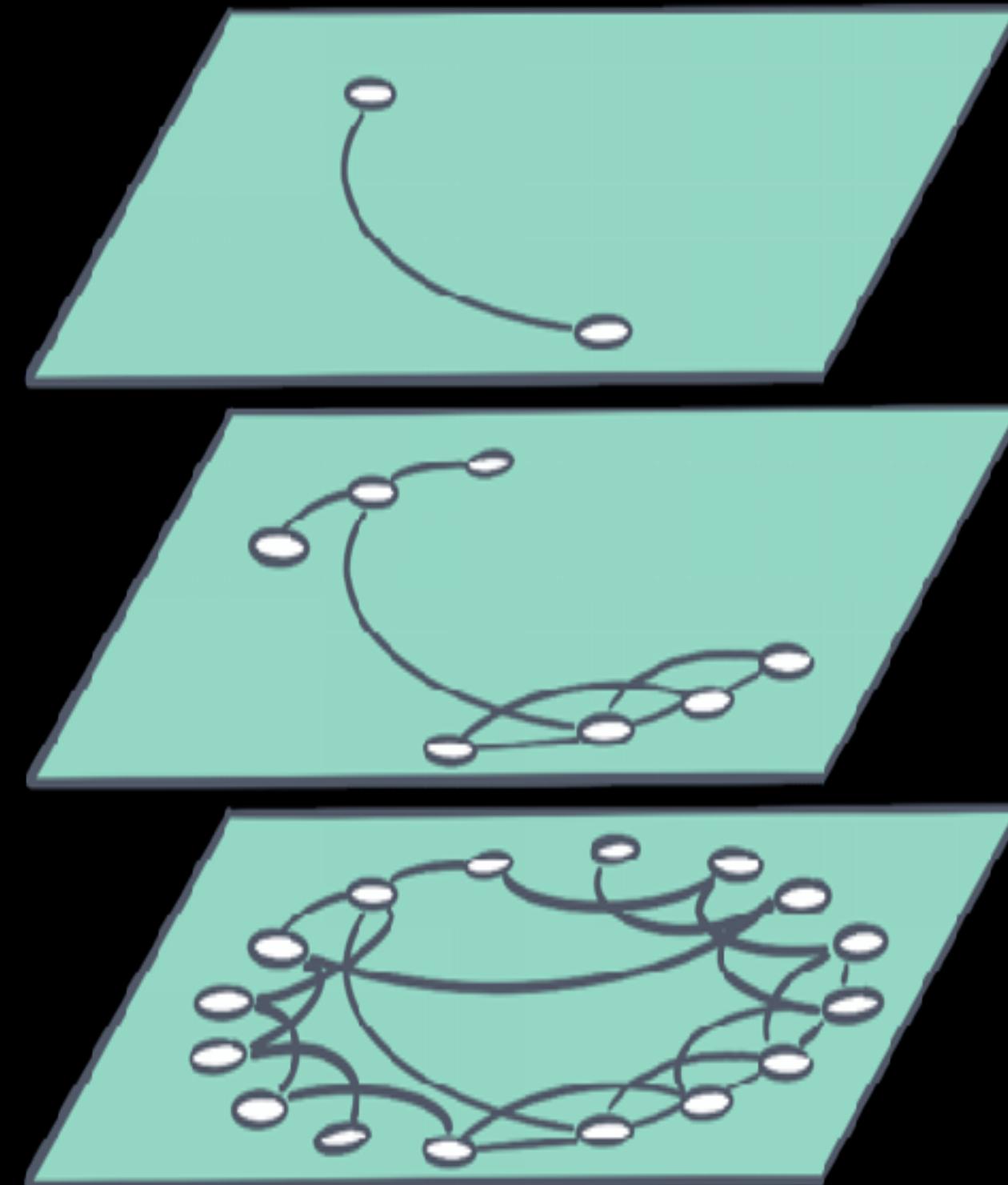
Tiered: from coarse to fine
approximation over a few steps



Balance: Bartering a little accuracy
for a lot of scalability



Speed: Excellent query latency on
large scale indices



VECTOR SEARCH

```
GET vector-search-star-wars-films/_search
{
  "query" : {
    "knn": {
      "field": "text_embedding.predicted_value",
      "k": 10,
      "num_candidates": 100,
      "query_vector": [-0.0078373895958599, -0.215008884668022],
    }
  }
}
```

YES!



BINGO CARD

AI

LLM

AGENTS

PROMPT
ENGINEERING

CONTEXT
ENGINEERING

VECTOR
SEARCH

RAG

A2A

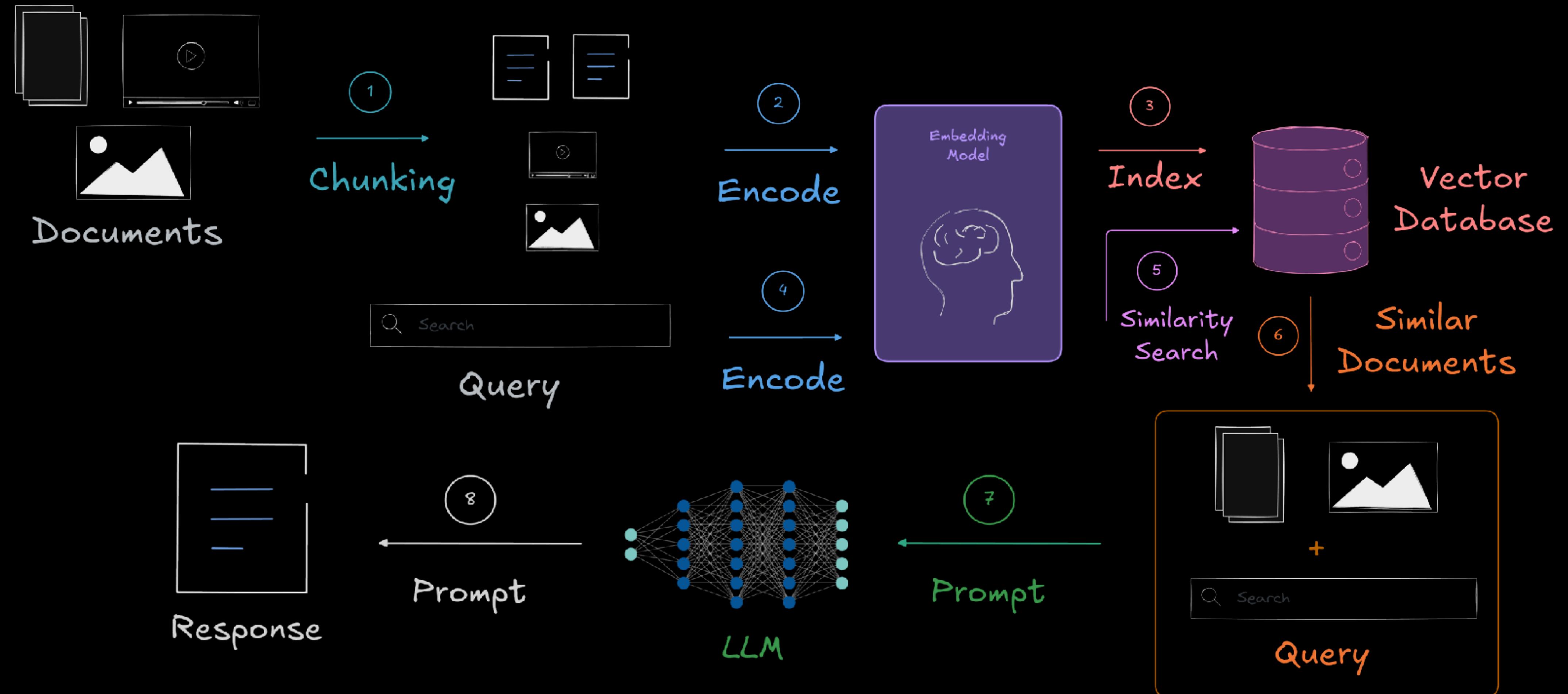
MCP



A close-up photograph of a person's hands interacting with a wire ball lottery cage. The cage is filled with numerous small, colorful balls. One hand is reaching into the cage, while the other holds a blue plastic tray underneath it. In the background, there are more wire cages and some blue plastic containers. A white cloth or paper is draped over one of the cages.

RAG

WHAT IS RAG?



You are a helpful AI designed to generate engaging, informative and concise descriptions for an automotive website in the UK. You will receive a list of facts about the vehicle, such as make, model, fuel type etc. You should use the features that the user provides and stay accurate to them.

Write a description for the vehicle in two coherent paragraphs with 3 or 4 sentences in each. Base the description on the following facts:

{...}

RAG-MCP: Mitigating Prompt Bloat in LLM Tool Selection via Retrieval-Augmented Generation

Tiantian Gan^{1,2} and Qiyao Sun^{1,2}

¹ Beijing University of Post and Communications, Beijing, China

² Queen Mary University of London, London, UK

jp2022213034@qmul.ac.uk, jp2022213402@qmul.ac.uk

Abstract. Large language models (LLMs) struggle to effectively utilize a growing number of external tools, such as those defined by the Model Context Protocol (MCP)[1], due to prompt bloat and selection complexity. We introduce RAG-MCP, a Retrieval-Augmented Generation framework that overcomes this challenge by offloading tool discovery. RAG-MCP uses semantic retrieval to identify the most relevant MCP(s) for a given query from an external index before engaging the LLM. Only the selected tool descriptions are passed to the model, drastically reducing prompt size and simplifying decision-making. Experiments, including an MCP stress test, demonstrate RAG-MCP significantly cuts prompt tokens (e.g., by over 50%) and more than triples tool selection accuracy (43.13% vs 13.62% baseline) on benchmark tasks. RAG-MCP enables scalable and accurate tool integration for LLMs.

BINGO CARD

AI

LLM

AGENTS

PROMPT
ENGINEERING

CONTEXT
ENGINEERING

VECTOR
SEARCH

RAG

A2A

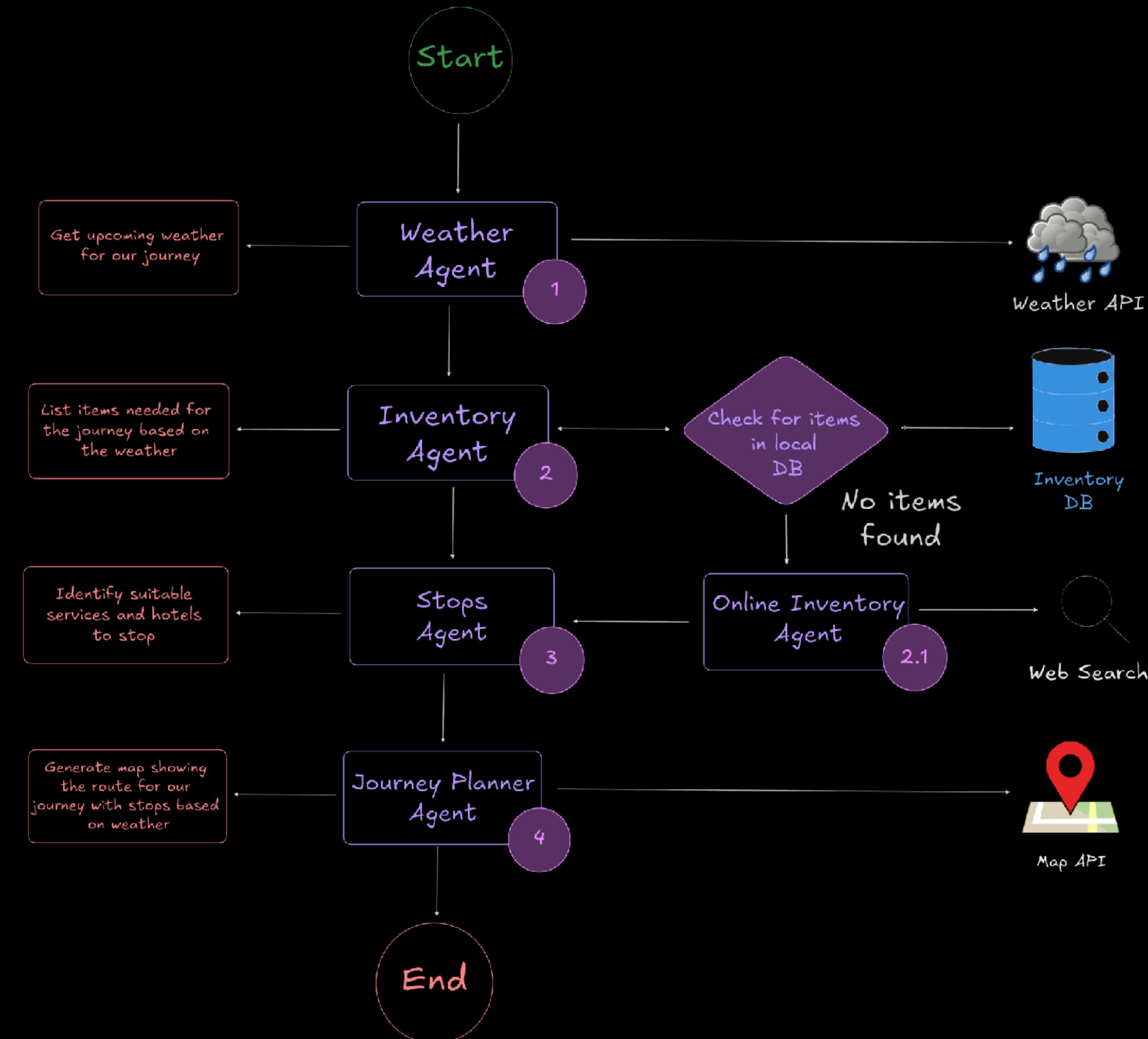
MCP



A2A

Example Multi-Agent Workflow

Tools & Functions



How A2A Works

descope

Client Agent

Remote Agent

Request Agent Card

01

Agent Card (JSON)

Task Assignment

02

Task Processing

Messages (Parts)

03

Messages (Parts)

Receive Updates

04

Status Updates

05

A2A Communication Flow

Discovery

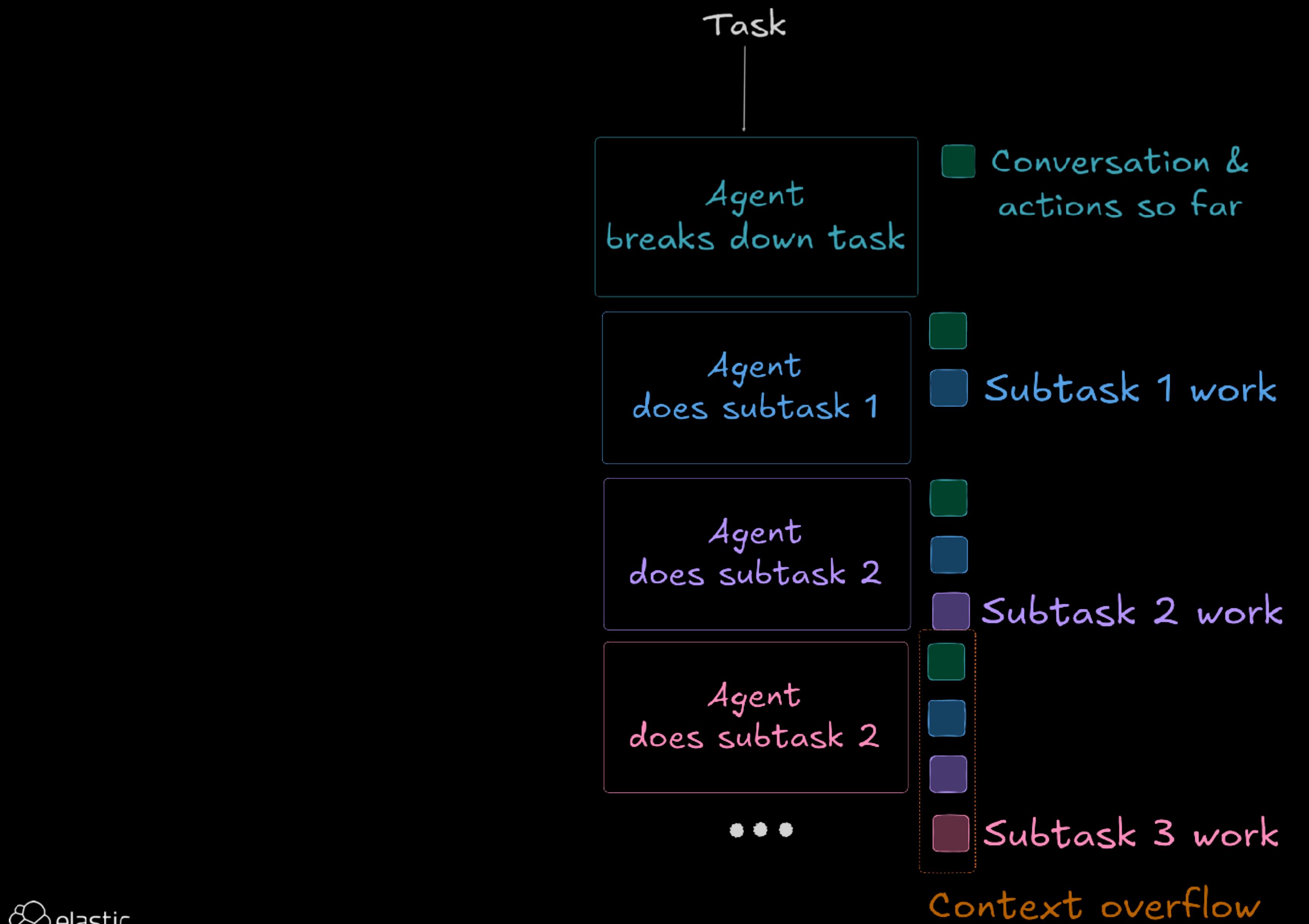
Task
Assignment

Communication

Task Progress

Completion
(Artifacts)

MULTI-AGENT CONTEXT



BINGO CARD

AI

LLM

AGENTS

PROMPT
ENGINEERING

CONTEXT
ENGINEERING

VECTOR
SEARCH

RAG

A2A

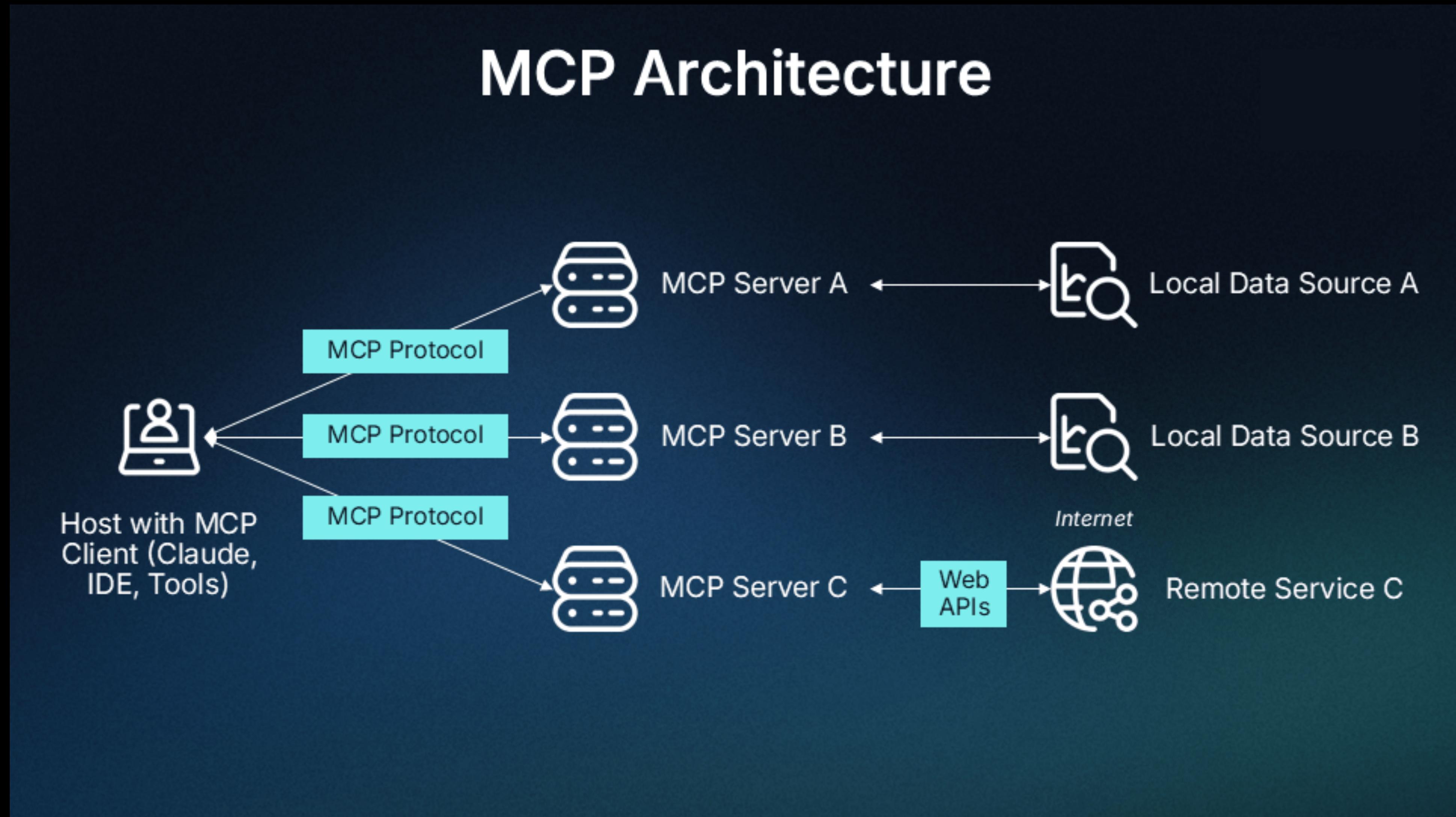
MCP



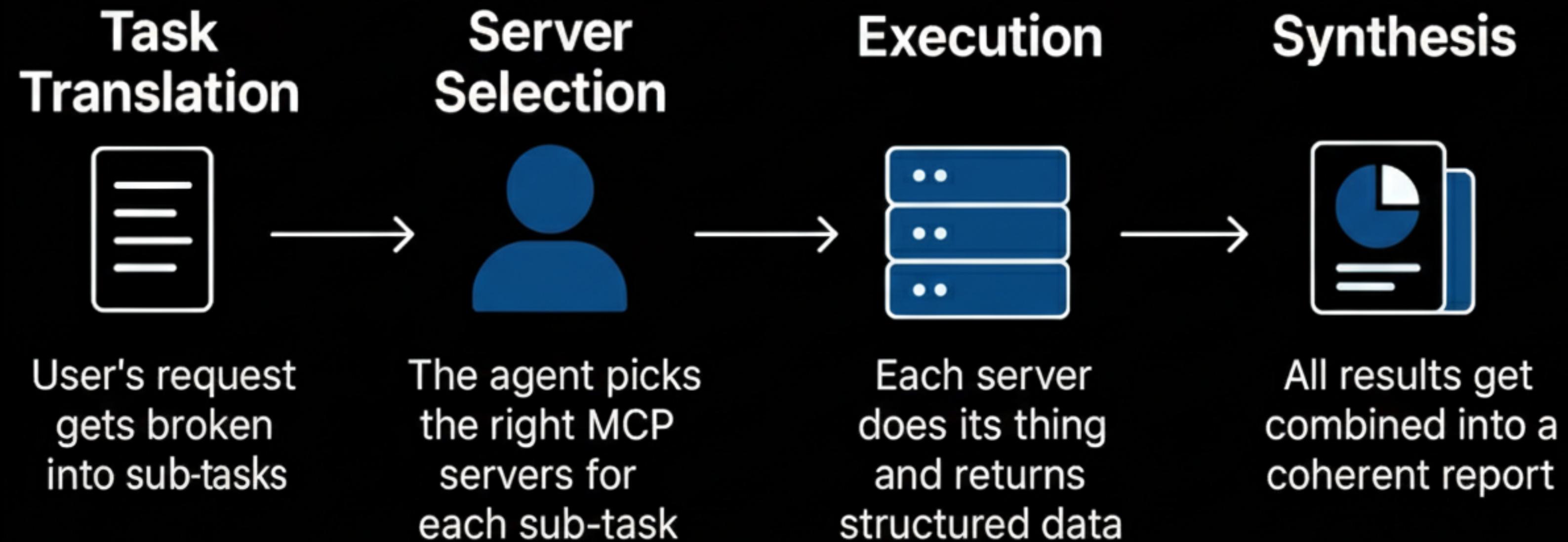
A photograph of a person's hand reaching into a metal lottery cage. The cage is filled with numerous small, colorful balls, each marked with a number. The background is slightly blurred, showing other people and what might be a festive or event setting. A large, solid purple circle is overlaid on the lower-left portion of the image, containing the text "MCP".

MCP

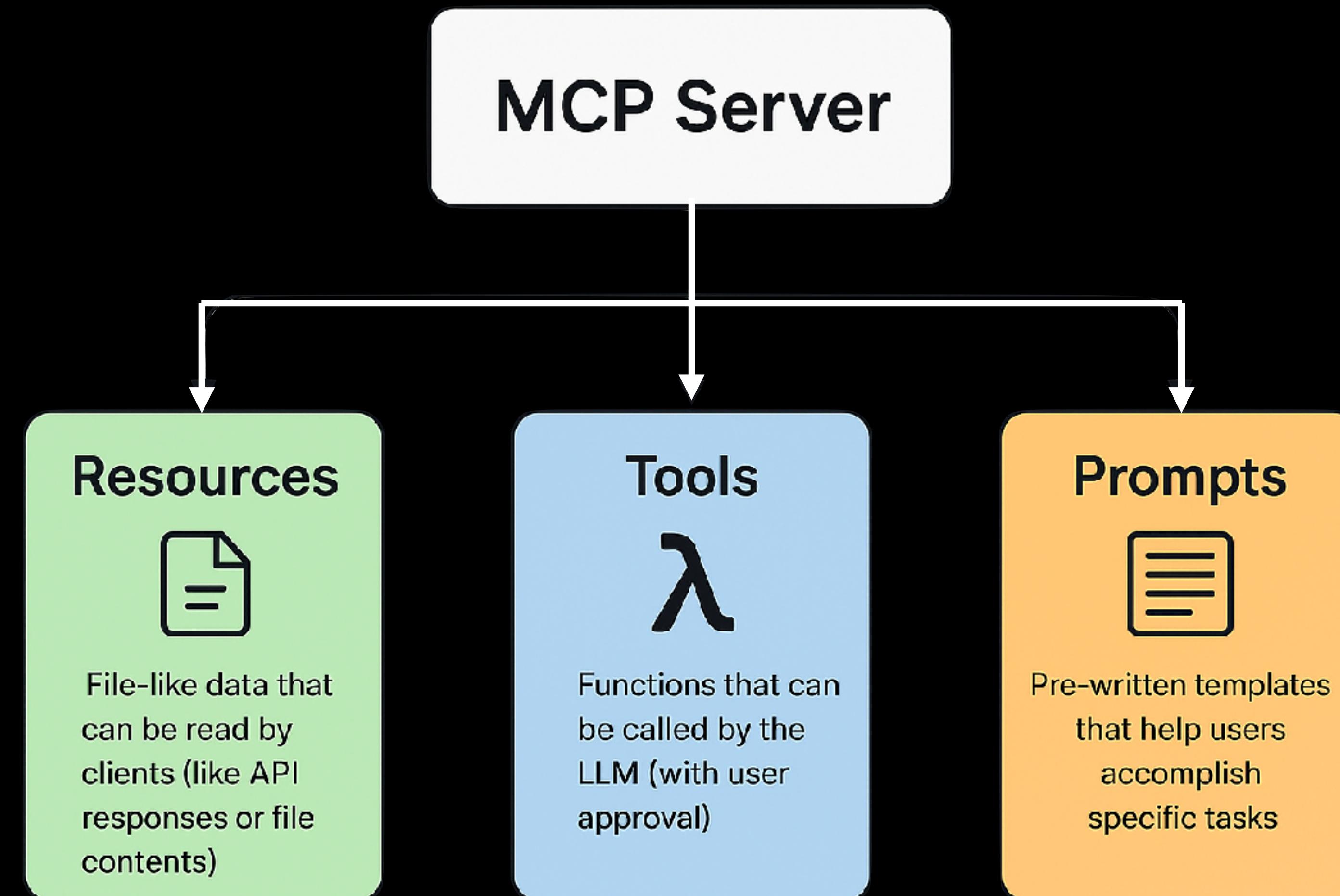
MCP



HOW DOES MCP WORK?



MCP COMPONENTS



EXAMPLE MCP RESOURCE

```
apple_watch_mcp.py

1 #!/usr/bin/env python3
2 """
3 Apple Watch Health Data MCP Server
4 A Model Context Protocol server for querying Apple HealthKit step data stored in Elasticsearch.
5
6 Credit: Alex Salgado
7 """
8 import os
9 from typing import Any, Optional
10 import json
11 from datetime import datetime
12 from pydantic import BaseModel, field_validator, ValidationError
13 from mcp.server.fastmcp import FastMCP
14 from elasticsearch import AsyncElasticsearch
15 from contextlib import asynccontextmanager
16
17 ## Other server initialization code omitted
18
19 @mcp.resource("health://steps/latest")
20 async def get_latest_steps() -> str:
21     """Gets the most recent step records"""
22     query = {
23         "query": {
24             "match_all": {}
25         },
26         "sort": [
27             {"endDate": {"order": "desc"}}
28         ],
29         "size": 10
30     }
31
32     data = await query_elasticsearch(query)
33     if not data:
34         return json.dumps({"error": "Unable to query Elasticsearch"}, indent=2)
35
36     results = []
37     for hit in data["hits"]["hits"]:
38         source = hit["_source"]
39         results.append({
40             "startDate": source.get("startDate"),
41             "endDate": source.get("endDate"),
42             "value": source.get("value"),
43             "device": source.get("device"),
44             "sourceName": source.get("sourceName"),
45             "dayOfWeek": source.get("dayOfWeek"),
46             "hour": source.get("hour")
47         })
48
49     return json.dumps({
50         "latest_steps": results
51     }, indent=2)
```



```
1 #!/usr/bin/env python3
2 """
3 Apple Watch Health Data MCP Server
4 A Model Context Protocol server for querying Apple HealthKit step data stored in Elasticsearch.
5
6 Credit: Alex Salgado
7 """
8 import os
9 from typing import Any, Optional
10 import json
11 from datetime import datetime
12 from pydantic import BaseModel, field_validator, ValidationError
13 from mcp.server.fastmcp import FastMCP
14 from elasticsearch import AsyncElasticsearch
15 from contextlib import asynccontextmanager
16
17 ## Other server initialization code comitted
18
19 @mcp.resource("health://steps/latest")
20 async def get_latest_steps() -> str:
21     """Gets the most recent step records"""
22     query = {
23         "query": {
24             "match_all": {}
```

```
19 @mcp.resource("health://steps/latest")
20 async def get_latest_steps() -> str:
21     """Gets the most recent step records"""
22     query = {
23         "query": {
24             "match_all": {}
25         },
26         "sort": [
27             {"endDate": {"order": "desc"}}
28         ],
29         "size": 10
30     }
31
32     data = await query_elasticsearch(query)
33     if not data:
34         return json.dumps({"error": "Unable to query Elasticsearch"}, indent=2)
35
36     results = []
37     for hit in data["hits"]["hits"]:
38         source = hit["_source"]
39         results.append({
40             "startDate": source.get("startDate"),
41             "endDate": source.get("endDate"),
42             "value": source.get("value"),
43             "device": source.get("device"),
44             "sourceName": source.get("sourceName"),
45             "isCSV": source.get("isCSV")})
```

```
26         "sort": [
27             {"endDate": {"order": "desc"}}
28         ],
29         "size": 10
30     }
31
32     data = await query_elasticsearch(query)
33     if not data:
34         return json.dumps({"error": "Unable to query Elasticsearch"}, indent=2)
35
36     results = []
37     for hit in data["hits"]["hits"]:
38         source = hit["_source"]
39         results.append({
40             "startDate": source.get("startDate"),
41             "endDate": source.get("endDate"),
42             "value": source.get("value"),
43             "device": source.get("device"),
44             "sourceName": source.get("sourceName"),
45             "dayOfWeek": source.get("dayOfWeek"),
46             "hour": source.get("hour")
47         })
48
49     return json.dumps(
50         {"latest_steps": results
51 }, indent=2)
```

EXAMPLE MCP RESOURCE

```
apple_watch_mcp.py

1 #!/usr/bin/env python3
2 """
3 Apple Watch Health Data MCP Server
4 A Model Context Protocol server for querying Apple HealthKit step data stored in Elasticsearch.
5
6 Credit: Alex Salgado
7 """
8 import os
9 from typing import Any, Optional
10 import json
11 from datetime import datetime
12 from pydantic import BaseModel, field_validator, ValidationError
13 from mcp.server.fastmcp import FastMCP
14 from elasticsearch import AsyncElasticsearch
15 from contextlib import asynccontextmanager
16
17 ## Other server initialization code committed
18
19 @mcp.resource("health://steps/latest")
20 async def get_latest_steps() -> str:
21     """Gets the most recent step records"""
22     query = {
23         "query": {
24             "match_all": {}
25         },
26         "sort": [
27             {"endDate": {"order": "desc"}}
28         ],
29         "size": 10
30     }
31
32     data = await query_elasticsearch(query)
33     if not data:
34         return json.dumps({"error": "Unable to query Elasticsearch"}, indent=2)
35
36     results = []
37     for hit in data["hits"]["hits"]:
38         source = hit["_source"]
39         results.append({
40             "startDate": source.get("startDate"),
41             "endDate": source.get("endDate"),
42             "value": source.get("value"),
43             "device": source.get("device"),
44             "sourceName": source.get("sourceName"),
45             "dayOfWeek": source.get("dayOfWeek"),
46             "hour": source.get("hour")
47         })
48
49     return json.dumps({
50         "latest_steps": results
51     }, indent=2)
```

HOUSE!

AI

LLM

AGENTS

PROMPT
ENGINEERING

CONTEXT
ENGINEERING

VECTOR
SEARCH

RAG

A2A

MCP

RESOURCES

- [Building an LLM-Powered Tool: From Concept to Production | Ioanna Nteka @ Autotrader Engineering](#)
- [Demystifying Large Language Models \(LLM101\) | Mahmoud Oshagh @ Autotrader Engineering](#)
- [Incident Summaries using LLMs | Karl Stoney](#)
- [What is Context Engineering | Carly Richmond @ Elasticsearch Labs](#)
- [12-Factor Agents - Principles for building reliable LLM applications | Dexter Horthy](#)
- [How to Fix Your Context | Drew Breunig](#)
- [The Prompt Report: A Systematic Survey of Prompt Engineering Techniques | Schulhoff et al.](#)
- [What is Memory | LangChain](#)
- [Don't Build Multi-Agents | Cognition](#)
- [Structured Outputs | LangChain](#)
- [AI Travel Agent Planner | Carly Richmond on GitHub](#)
- [RAG-MCP: Mitigating Prompt Bloat in LLM Tool Selection via Retrieval-Augmented Generation | Gan and Sun](#)





SCAN ME