

LLMs in 2025 crash-course

Agent(MCP(tools(rag(llm))))) ?????

Who am I?

- Jiaqi Chen
- 24 yr
- BSc AI, University Utrecht
- Start-ups / Entrepreneurship
 - Own startup
 - incubator board year
 - DialogueTrainer Product Manager
- Bouldering, live music, photography

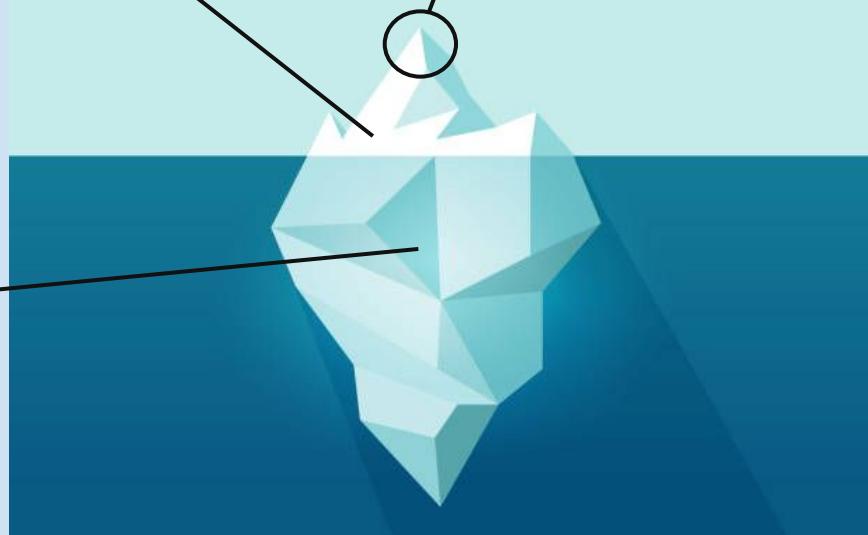


Preface

Techical

Societal,
mental/cognitive,
political,
Economic,
Cultural,
Academic,
Environmental,
Ethical,
Educational

This presentation



"Crash course in LLMs"

Techical

Societal,
Mental/cognitive
political,
Economic,
Cultural,
Academic,
Environmental,
Ethical,
Educational

Costs

This presentation

Effects

Effects

“Crash course in LLMs”

Costs Effects

Techicha

Costs Effects

Societal,

Costs / Effects

political,

Costs Effects

Cultural

Costs Effects

Environmental

Costs & Effects

Educa⁺ - Educational

Costs Effects

Costs Effects Costs

This presentation is part of the [OpenCourseWare](#) project.

Costs Effects Costs

—
—
—

Costs Effects Costs

© 2018 [Eduardo](#)

Costs Effects Costs

© 2010 Pearson Education, Inc.

Costs Effects Costs

www.earth-works.com

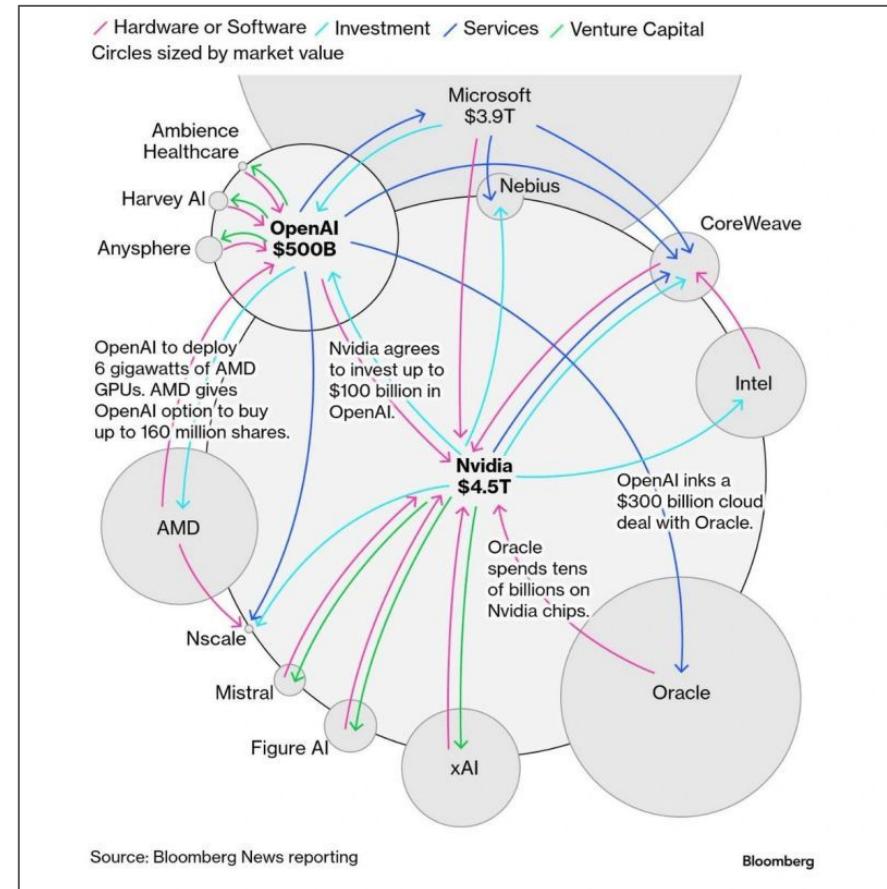
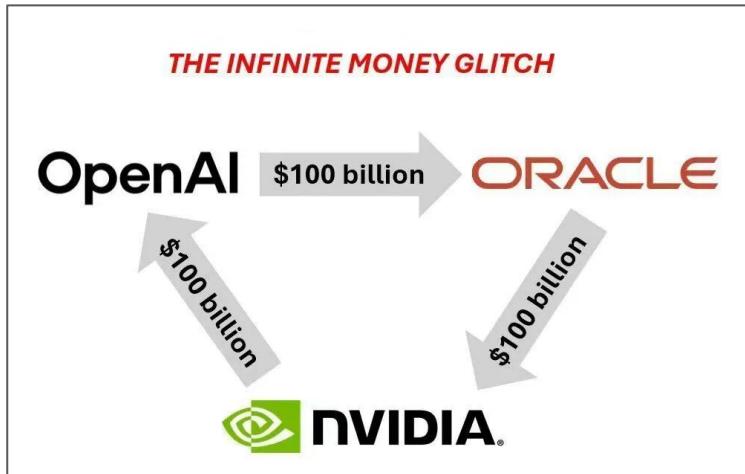
Costs Effects Costs

 Studydrive - Crash courses in life

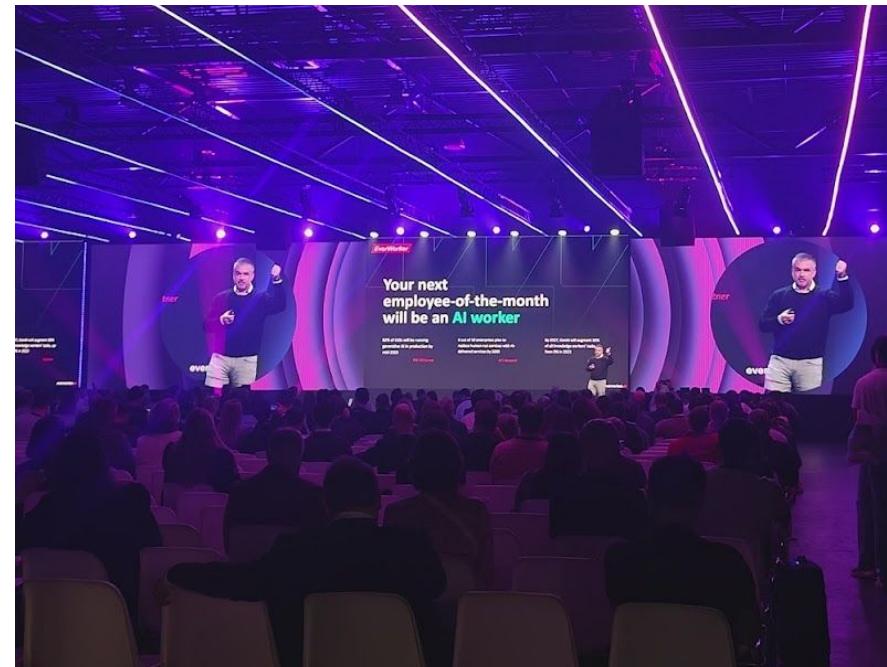
Costs Effects Costs

Financial 'dimension'

The investment space of tech in 2025



Human centric AI



❖ AI Overview

Tantalum is used primarily in electronics as a component for capacitors and resistors in devices like phones and computers. It is also used in high-temperature alloys for jet engines and nuclear reactors, and due to its corrosion resistance and biocompatibility, it is ideal for medical implants and surgical devices.

❖ AI Overview

Cobalt is used primarily in lithium-ion batteries for electronics and electric vehicles, in high-strength, heat-resistant alloys for jet engines and turbines, and in pigments to create vibrant blue and green colors. It is also used in magnetic materials, as a catalyst in the petroleum industry, and the radioactive isotope cobalt-60 is used for cancer treatment and as a gamma-ray source.



22 Ti Titanium	$-1.63(2)$ 1.3 0.41 $[Ar] 4d^2 5s^2$	-0.42^{+} 79^{+} $-1.175(2)$ 1.5 0.01 $[Kr] 4d^2 5s$	-124.9α 118 281^{+} 813^{+} $-0.744(3)$ 1.6 0.02 $[Kr] 4d^2 5s$	23 V Vanadium	$-1.63(2)$ 1.3 0.41 $[Ar] 4d^2 5s^2$	$-1.175(2)$ 1.5 0.01 $[Kr] 4d^2 5s$	24 Cr Chromium	$-0.744(3)$ 1.6 0.02 $[Kr] 4d^2 5s$	25 Mn Manganese	-1.3 1.3 0.02 $[Kr] 4d^2 5s$				
41 Nb Niobium	-1.45 59^{+} 72^{+} $-1.553(4)$ 1.2 0.02 $[Ar] 4d^2 5s^2$	$-1.42.9$ 134 64^{+} 72^{+} $1099(3)$ 1.2 10^{-3}	$-1.136.3$ 130 59^{+} 61^{+} $-0.200(3)$ 1.3 10^{-3}	42 Mo Molybdenum	-1.42 5.3 1.2 0.02 $[Xe] 4f^{14} 5d^4 6s^2$	-1.437 130 60^{+} 66^{+} $-0.09(4)$ 1.4 $6 \cdot 10^{-3}$	43 Tc Technetium	$-1.435.2$ 7 1.3 137.1 128 5	73 Ta Tantalum	$-1.48 \cdot 10^{-4}$ $8 \cdot 10^{-4}$ $[Rn] 5f^{14} 6d^3 7s^2$	74 W Tungsten	$-1.46 \cdot 10^{-3}$ $[Rn] 5f^{14} 6d^4 7s^2$	75 Re Rhenium	$-1.4[Rn] 5f^{14} 6d^5$
76 Os Osmium	-1.476^{+} $05(4)$ 1.2 0^{-4} $[Xe] 4f^{14} 5d^6 6s^2$	$-1.4143.0$ 13 1.2 0^{-4}	-1.4137 130 60^{+} 66^{+} $-0.09(4)$ 1.4 $6 \cdot 10^{-3}$	76 Os Osmium	-1.476^{+} $05(4)$ 1.2 0^{-4} $[Xe] 4f^{14} 5d^6 6s^2$	-1.4137 130 60^{+} 66^{+} $-0.09(4)$ 1.4 $6 \cdot 10^{-3}$	76 Os Osmium	-1.476^{+} $05(4)$ 1.2 0^{-4} $[Xe] 4f^{14} 5d^6 6s^2$	76 Os Osmium	-1.476^{+} $05(4)$ 1.2 0^{-4} $[Xe] 4f^{14} 5d^6 6s^2$	76 Os Osmium	-1.476^{+} $05(4)$ 1.2 0^{-4} $[Xe] 4f^{14} 5d^6 6s^2$		
77 Rh Rhodium	-1.476^{+} $05(4)$ 1.2 0^{-4} $[Xe] 4f^{14} 5d^7 6s^1$	$-1.4143.0$ 13 1.2 0^{-4}	-1.4137 130 60^{+} 66^{+} $-0.09(4)$ 1.4 $6 \cdot 10^{-3}$	77 Rh Rhodium	-1.476^{+} $05(4)$ 1.2 0^{-4} $[Xe] 4f^{14} 5d^7 6s^1$	-1.4137 130 60^{+} 66^{+} $-0.09(4)$ 1.4 $6 \cdot 10^{-3}$	77 Rh Rhodium	-1.476^{+} $05(4)$ 1.2 0^{-4} $[Xe] 4f^{14} 5d^7 6s^1$	77 Rh Rhodium	-1.476^{+} $05(4)$ 1.2 0^{-4} $[Xe] 4f^{14} 5d^7 6s^1$	77 Rh Rhodium	-1.476^{+} $05(4)$ 1.2 0^{-4} $[Xe] 4f^{14} 5d^7 6s^1$		



Congo is one of the world's main sources of cobalt and other critical minerals that power batteries, devices, and AI hardware.

It should have been one of the most prosperous regions on earth, yet it has remained one of the hardest places to be born for over a century.

Resource extraction has repeatedly led to conflict, exploitation, and systemic instability. It is a tragic illustration of how technological progress, global markets, and human suffering can become entangled when cost is measured only in money or performance.

Arthur C. Clarke's Third Law:

Any sufficiently advanced
technology is
indistinguishable
from magic

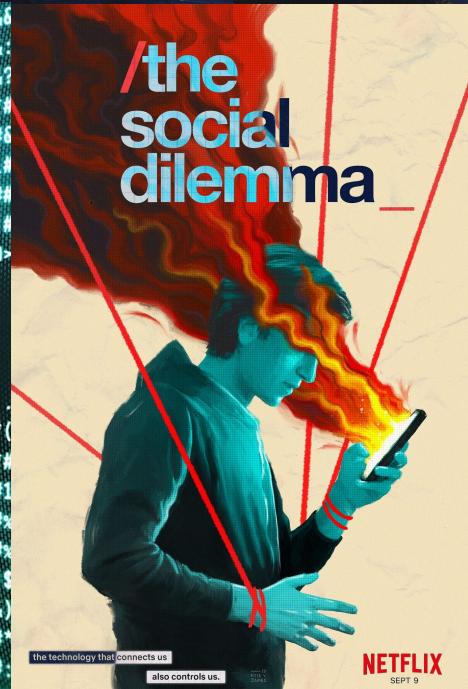
Techical

Societal,
mental/cognitive,
political,
Economic,
Cultural,
Academic,
Environmental,
Ethical,
Educational

Arthur C. Clarke's Third Law:

Any sufficiently advanced technology is indistinguishable from magic

Magic



AI Projects

AI projects

Clients

Students

Mentors

Bondfish

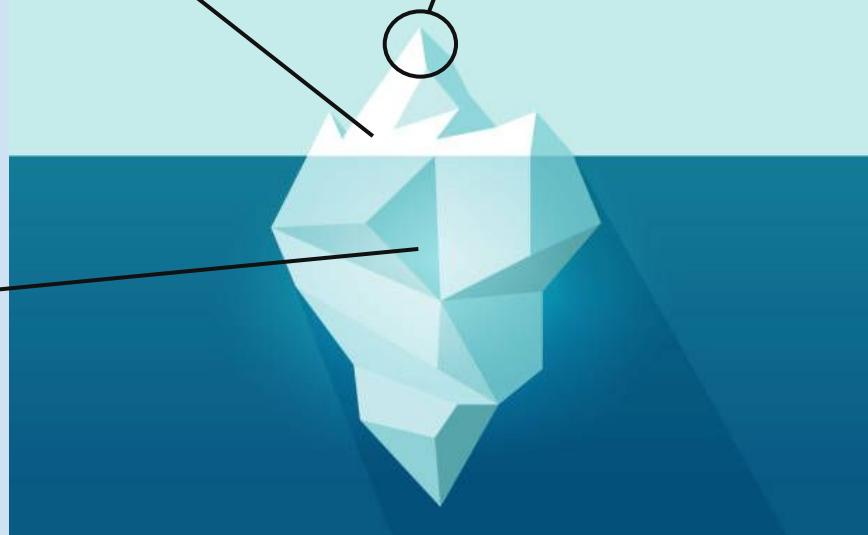
Gradient Metrics

Profectus

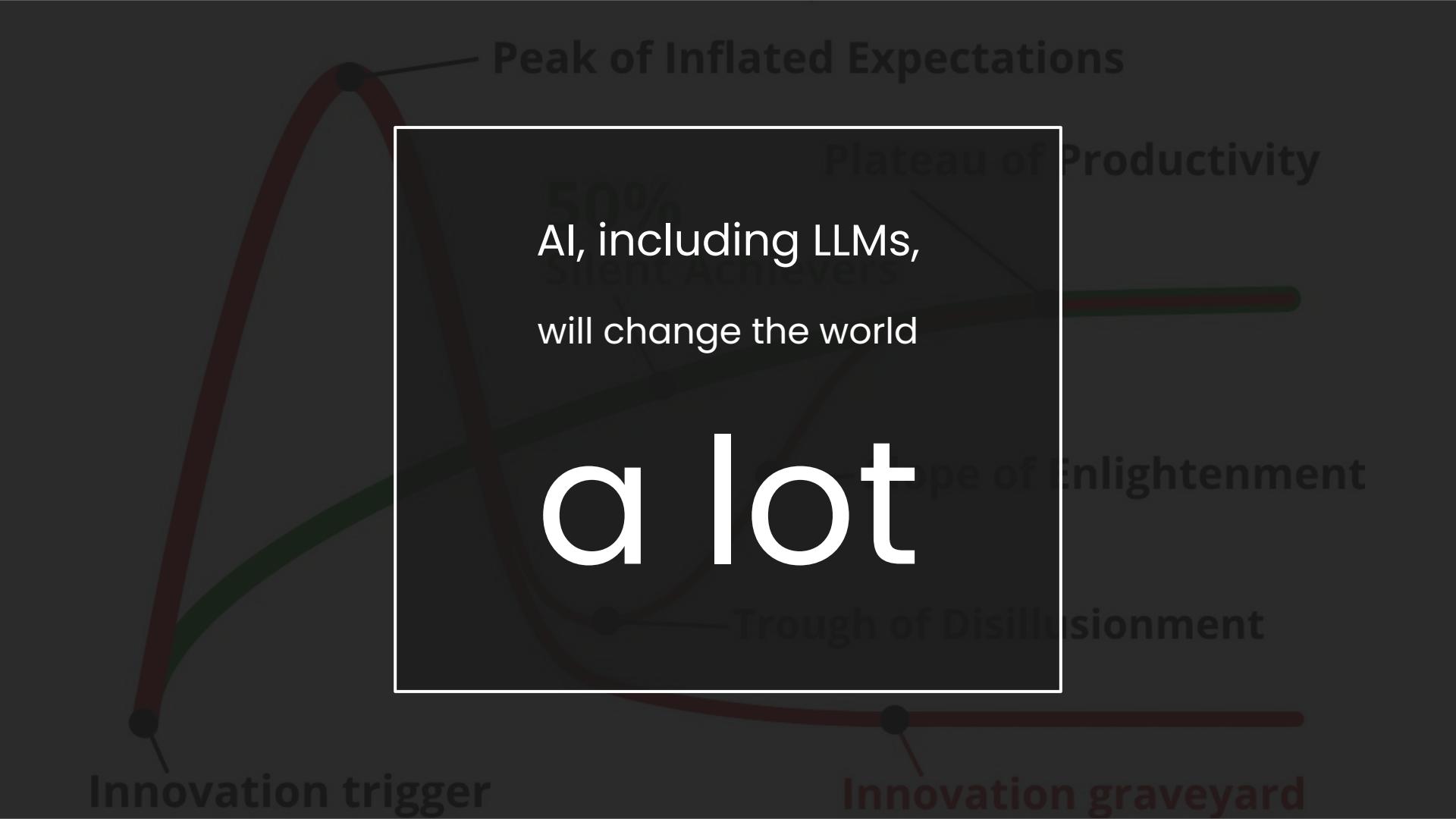
Techical

Societal,
mental/cognitive,
political,
Economic,
Cultural,
Academic,
Environmental,
Ethical,
Educational

This presentation



"Crash course in LLMs"



Peak of Inflated Expectations

Plateau of Productivity

AI, including LLMs,
will change the world

a lot

Innovation graveyard

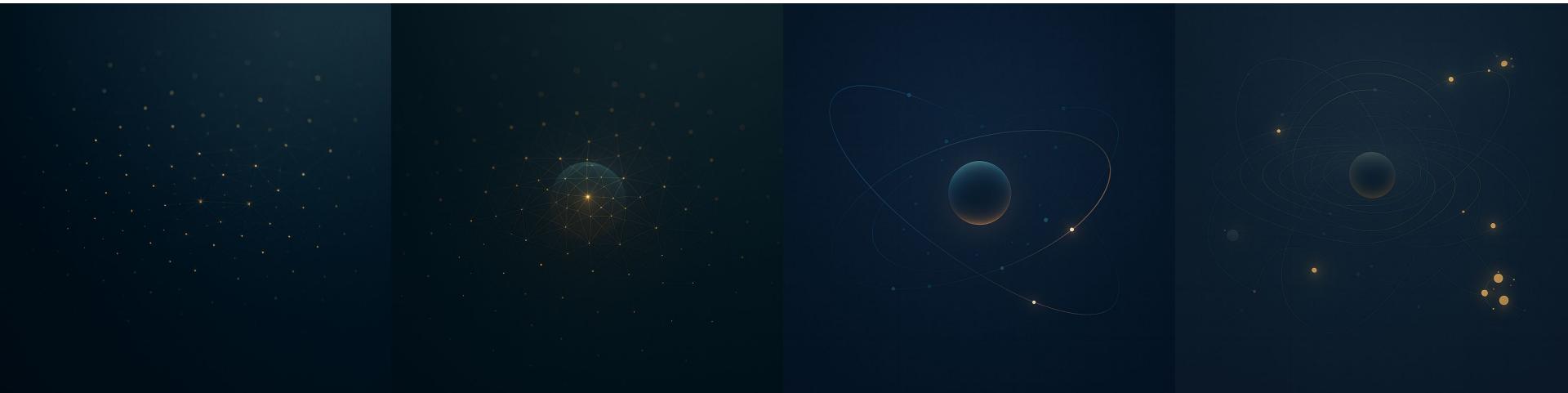
Innovation trigger

Pre-LLM

LLM

RAG

Agentic



Terms don't strictly follow timeline

Pre-LLM

LLM

RAG

Agentic

Symbolic AI

Attention mechanism

System Prompting

Autonomous Workflow

Transformer Architecture

Hybrid Retrieval

Feature Engineering

Vector Databases & Embeddings

Goal Decomposition

Word Embeddings

Knowledge Distillation

Multi-Agent systems

Semantic search

LLM to SQL

Instruction-Tuning

Tool Use & Function Calling

Meta-Prompting

MoE

Tokenization

Chunking & Indexing

MCP

Sequence Models (RNN/LSTM)

Context Window

Context Compression / Engineering

RLHF

AGI

Prompt engineering

Human-centered AI

NLP (Natural Language Processing)

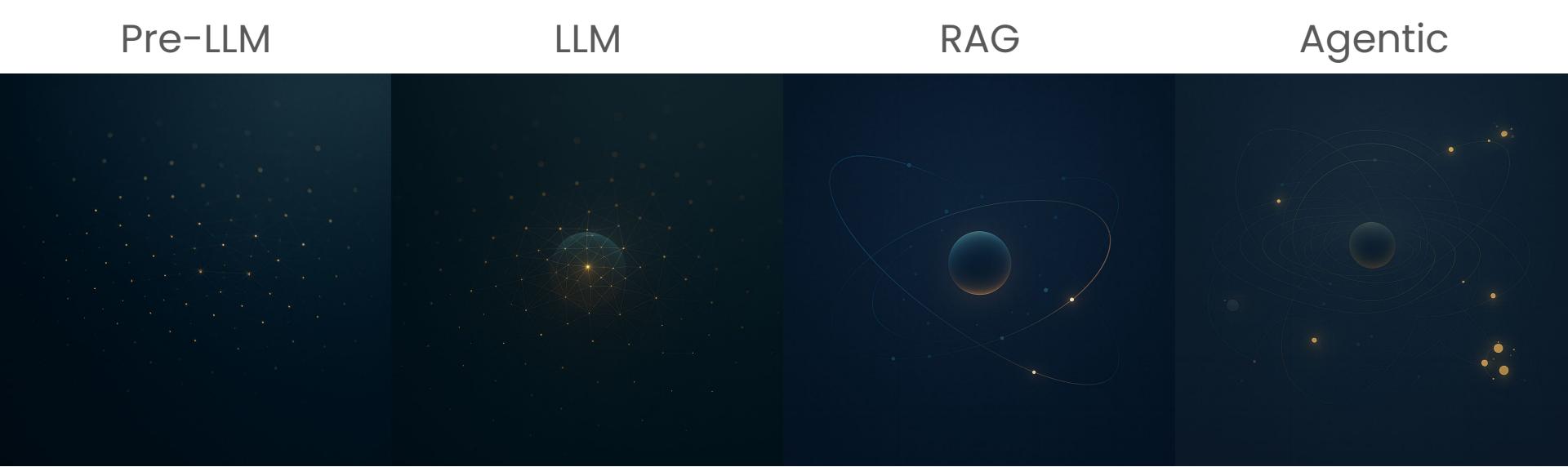
Human-Computer Interaction (HCI)

Pre-LLM

LLM

RAG

Agentic



Data pipelines,
Sklearn pipelines
Embeddings
Tokenization

LLM/SLM
LLM
pipelines

Basic RAG
(Hybrid)
indexing

LLM
interactions
overview

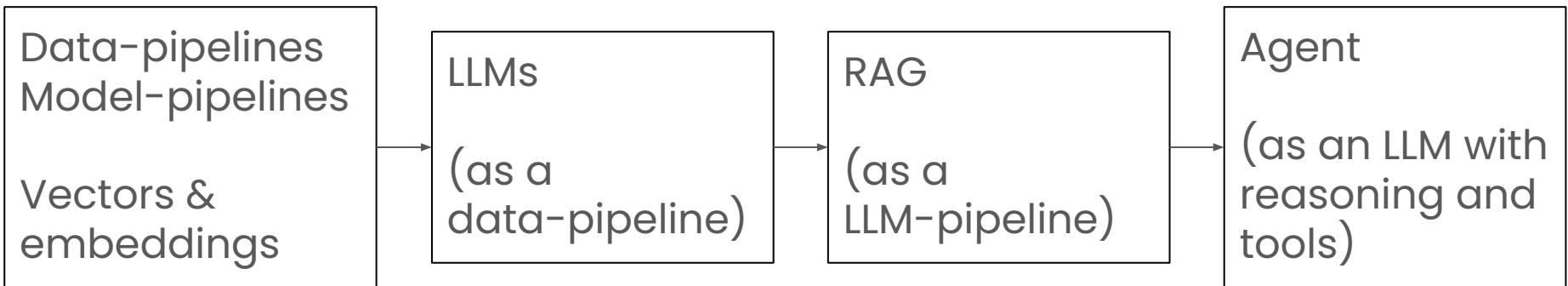
MCP

Pre-LLM

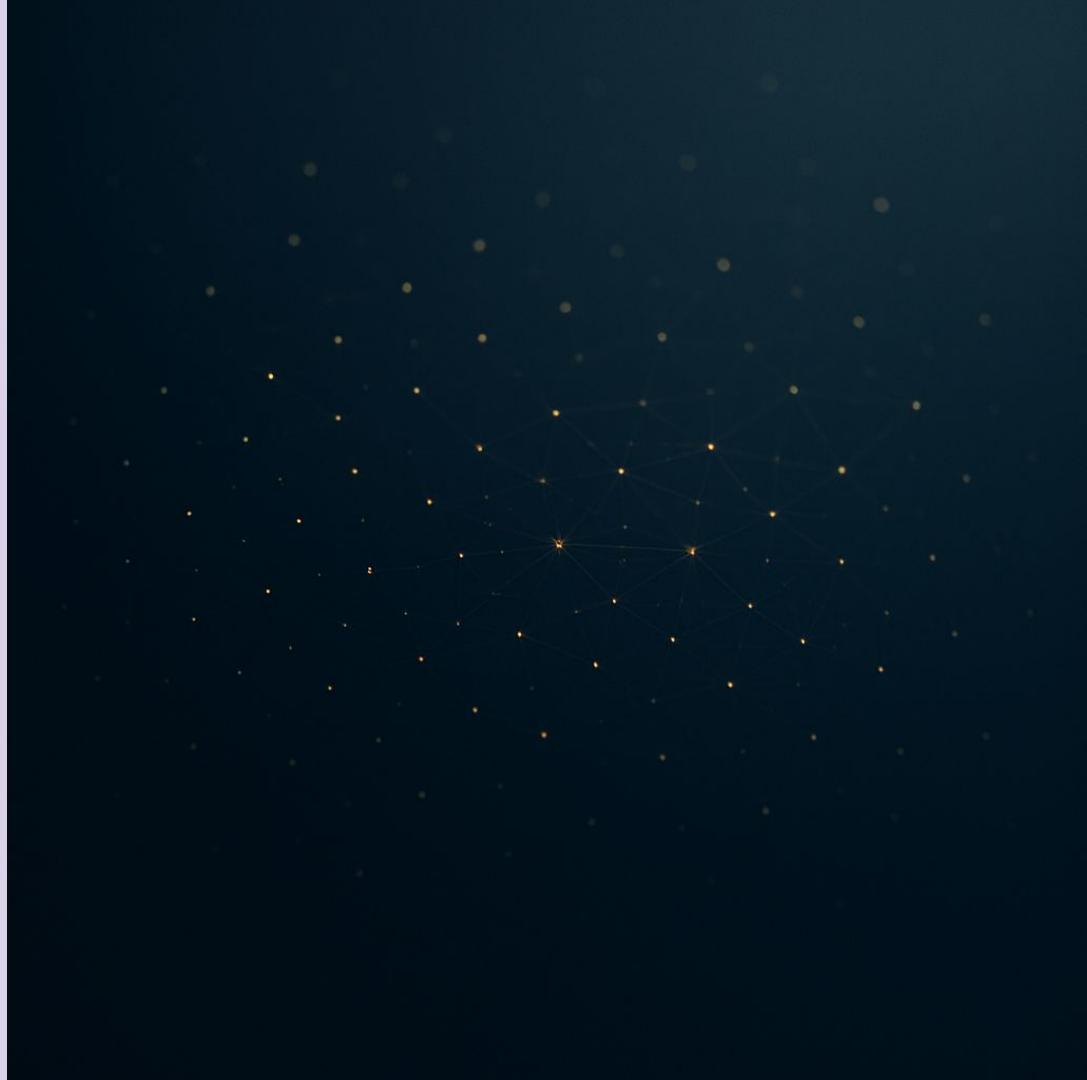
LLM

RAG

Agentic



Pre
LLM

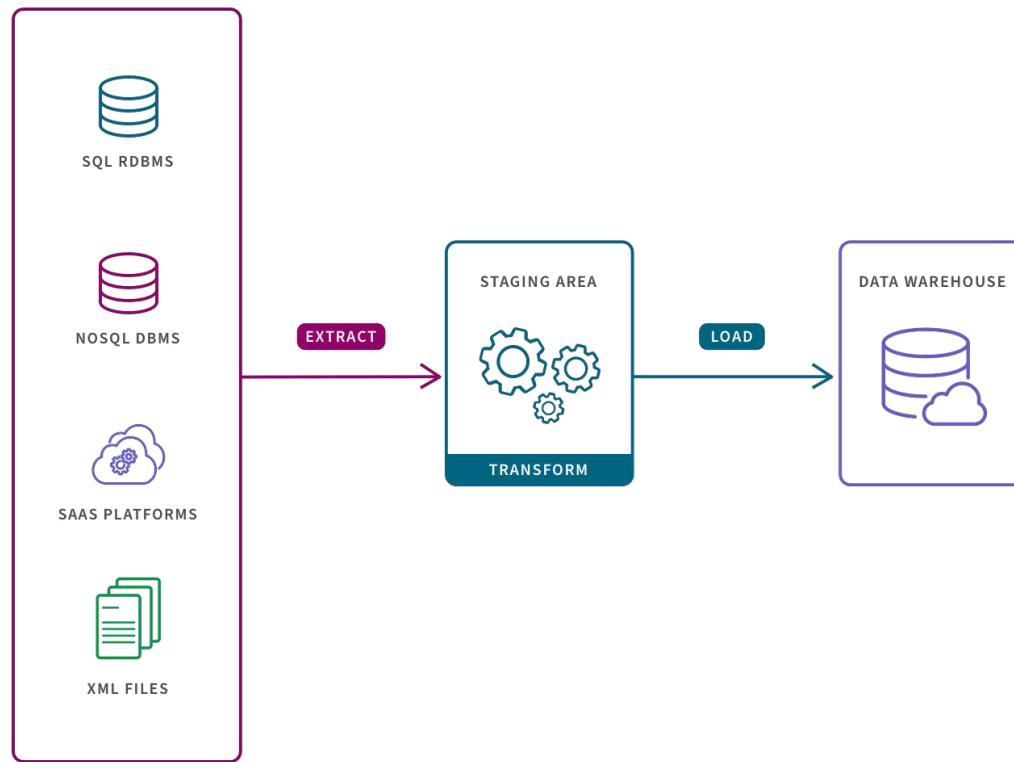


Data Pipelines

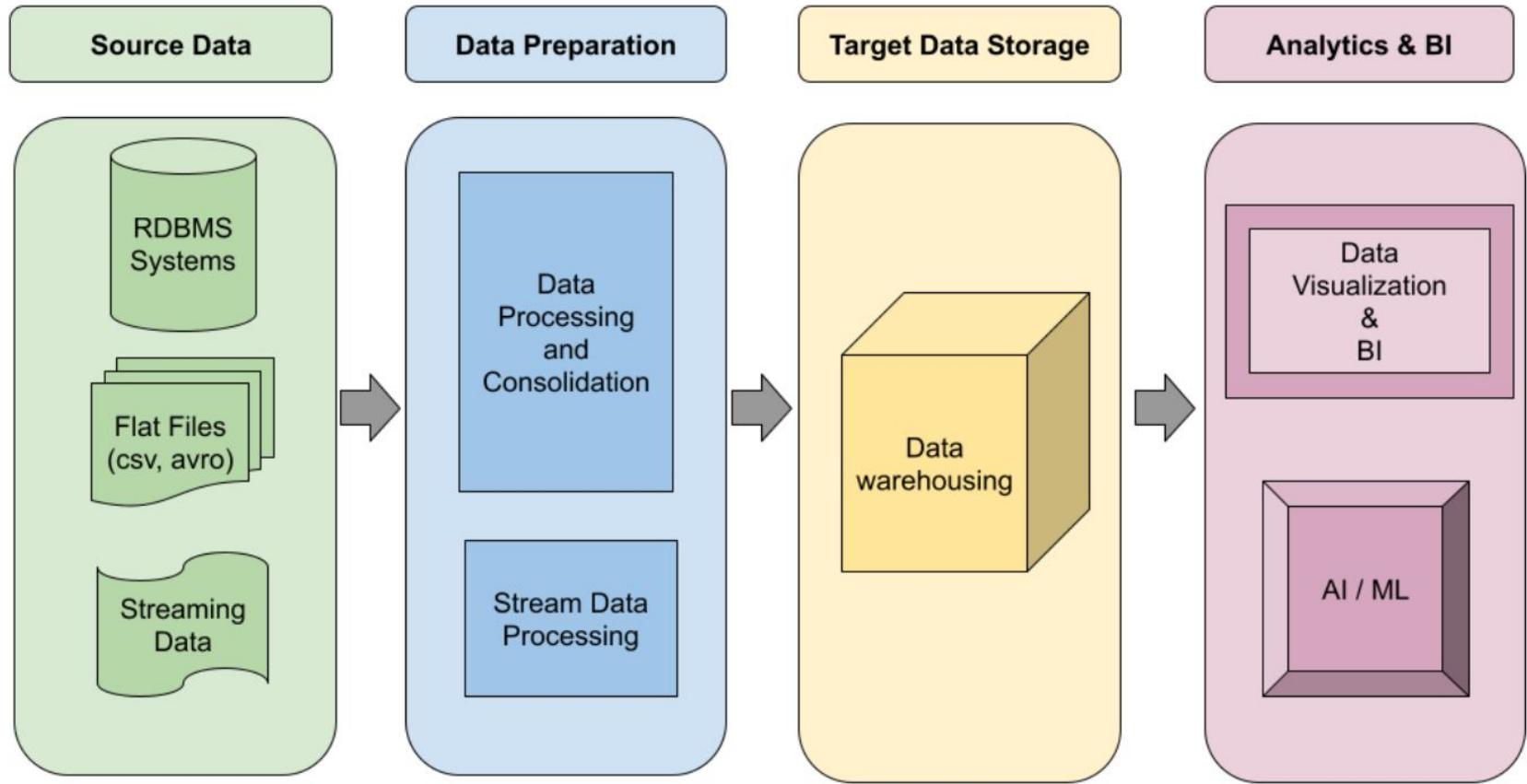


Level 1: input-output

ETL (Extract, Transform, Load)

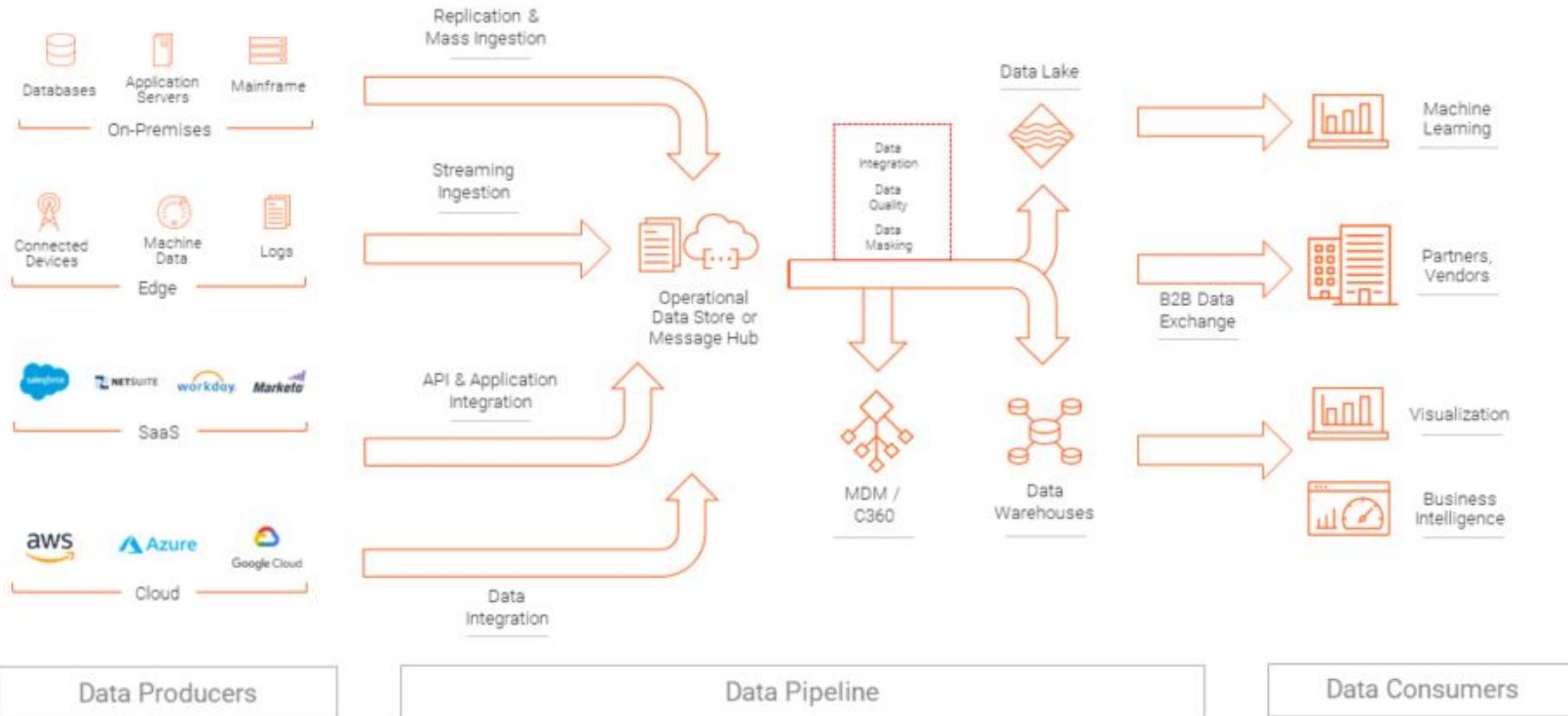


Level 2: ETL



Level 3: high level end to end

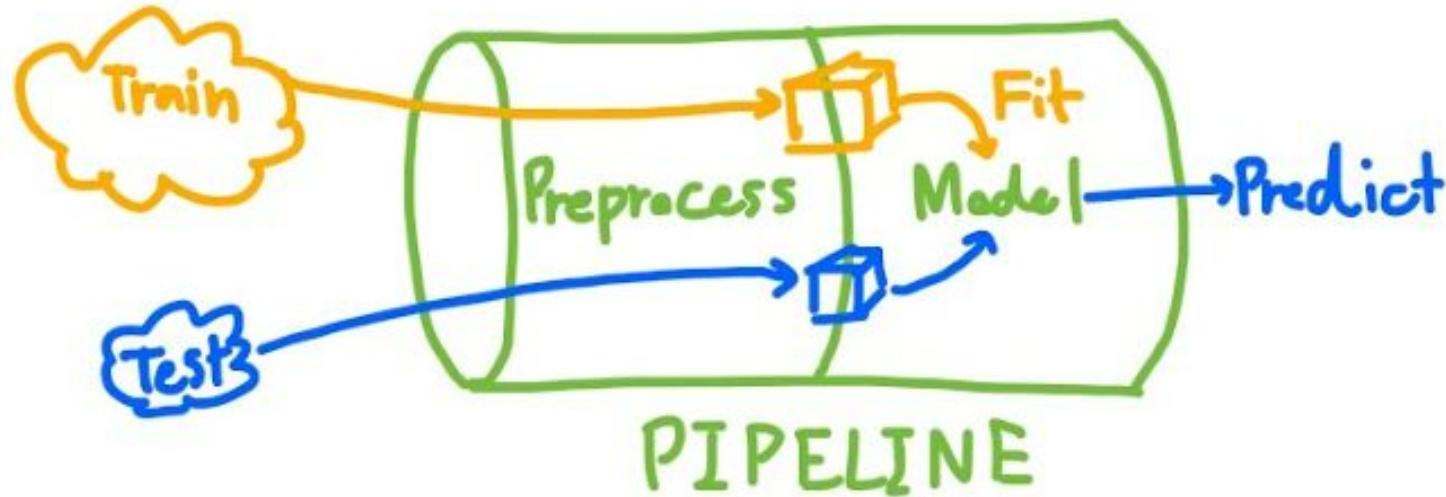
Data Pipeline Patterns



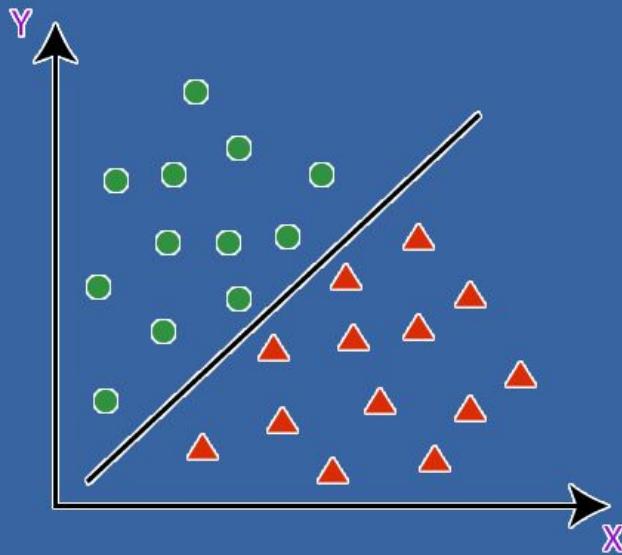
Level 4: integrations and more choices

Pre-ilm models

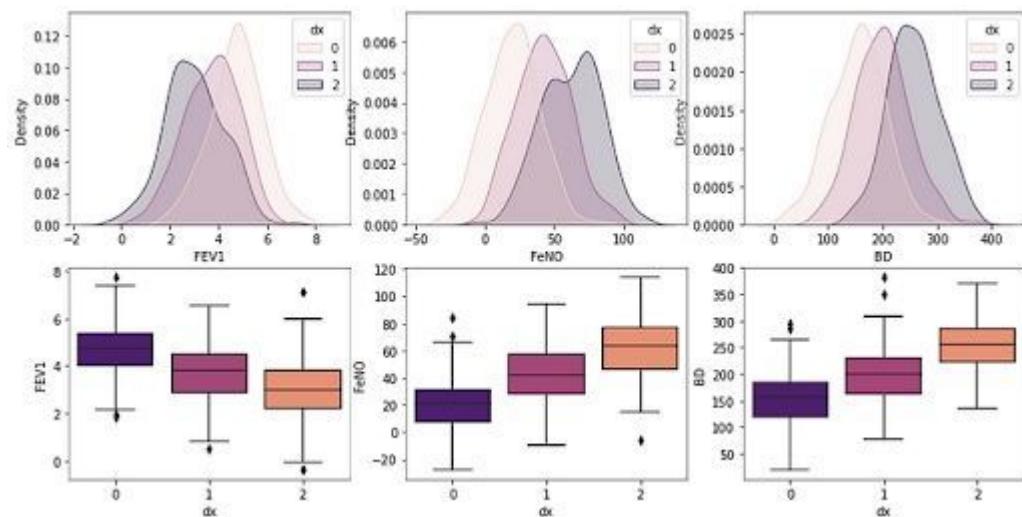
(AI pre-covid)



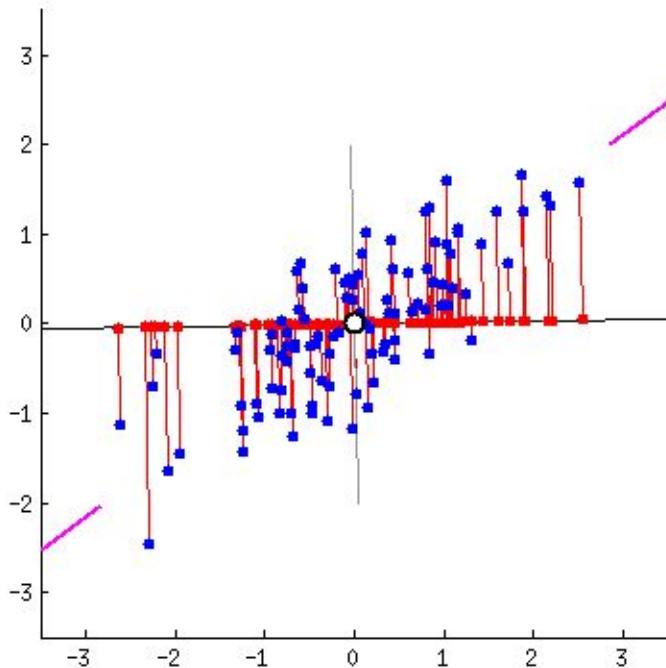
SVM



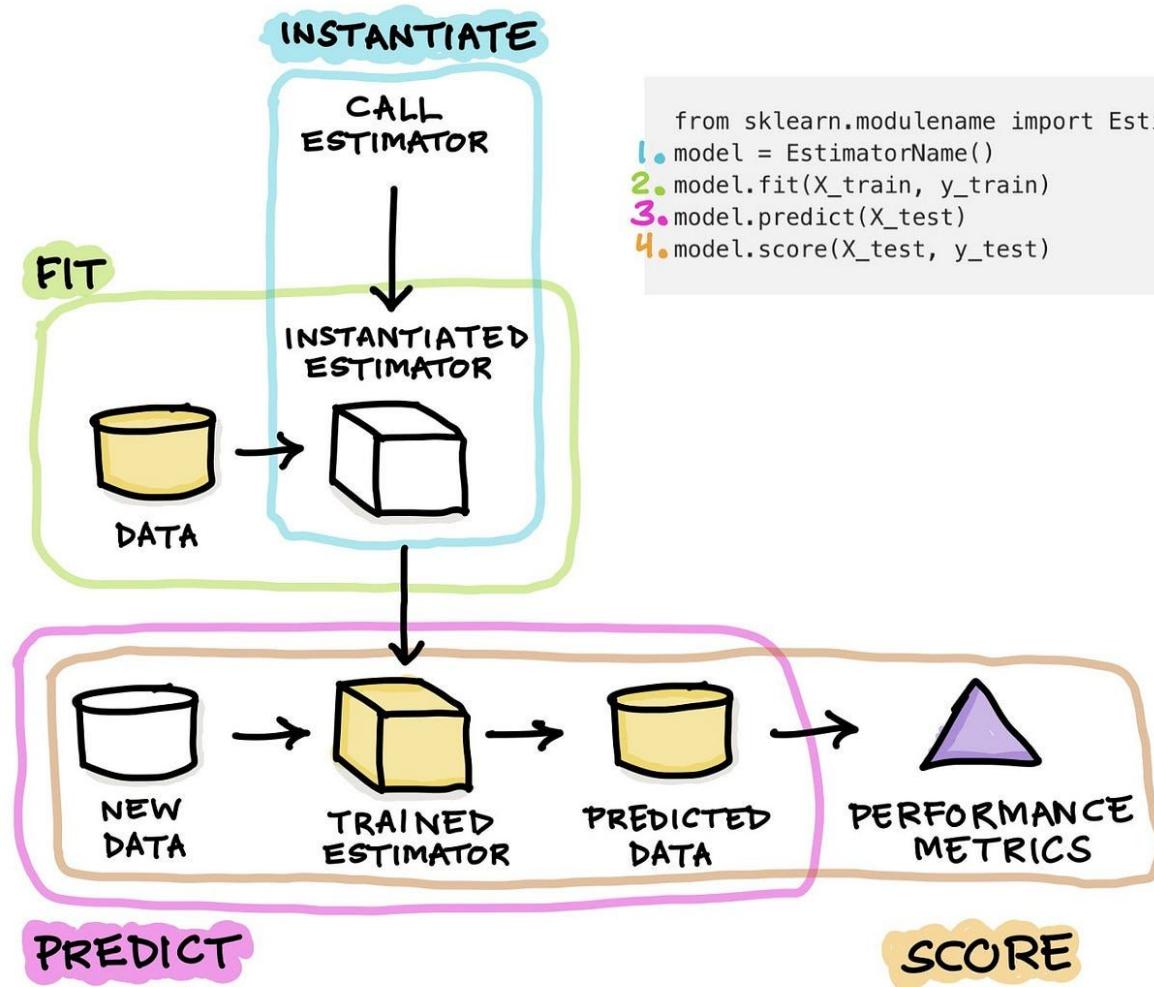
Bayes



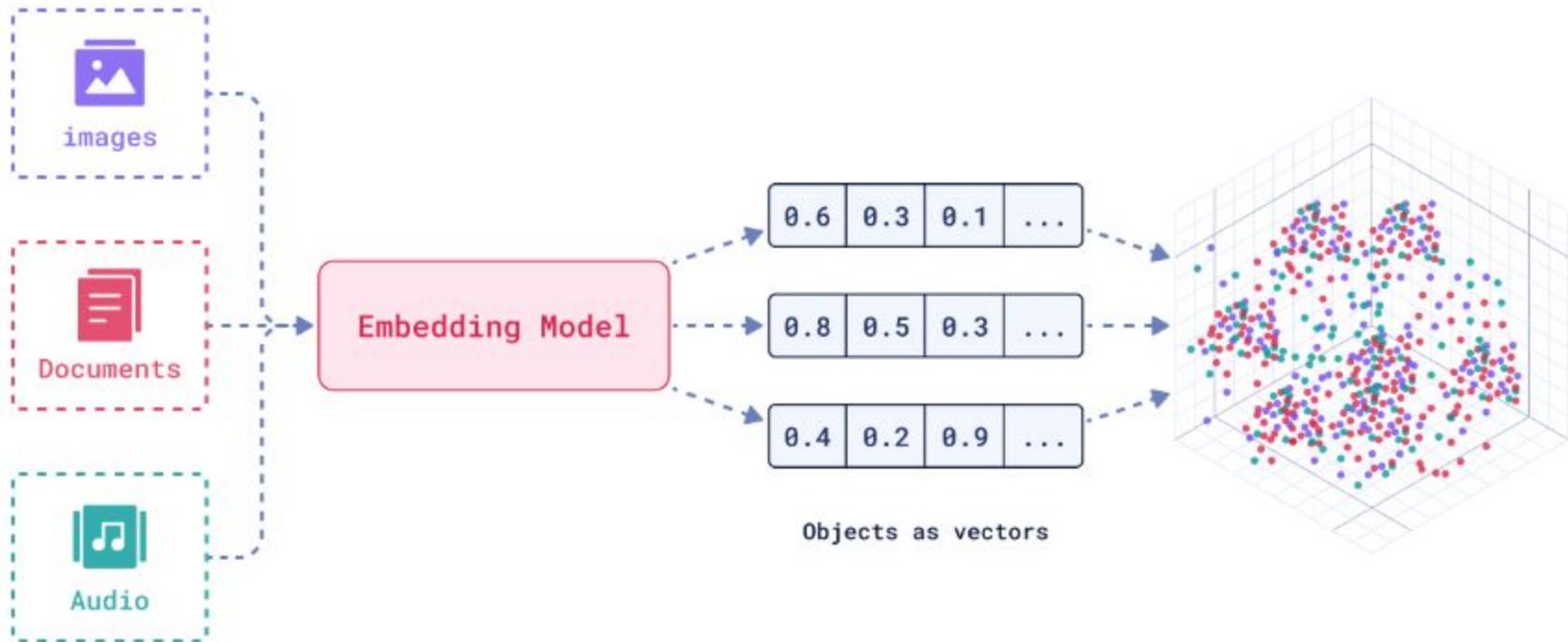
PCA



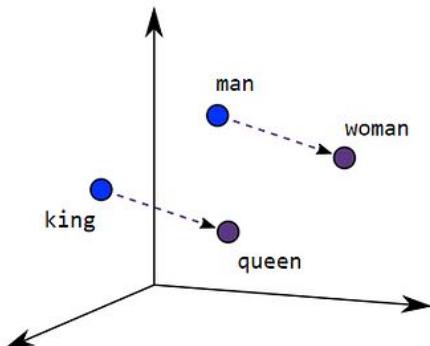
Sk-Learn



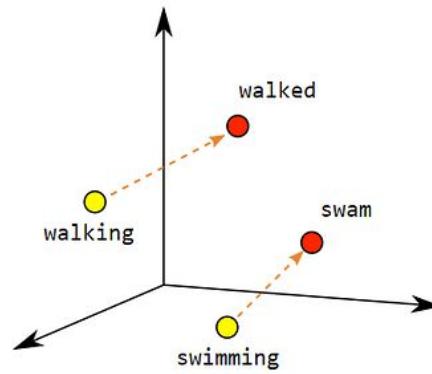
Embeddings



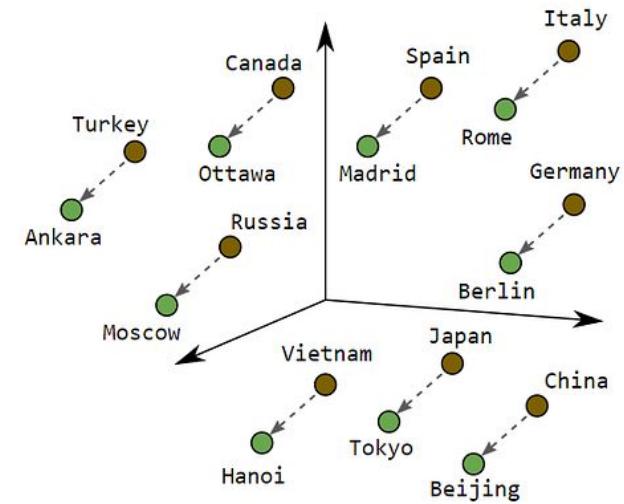
Embeddings can capture relations between the vectors (in a latent space)



Male-Female

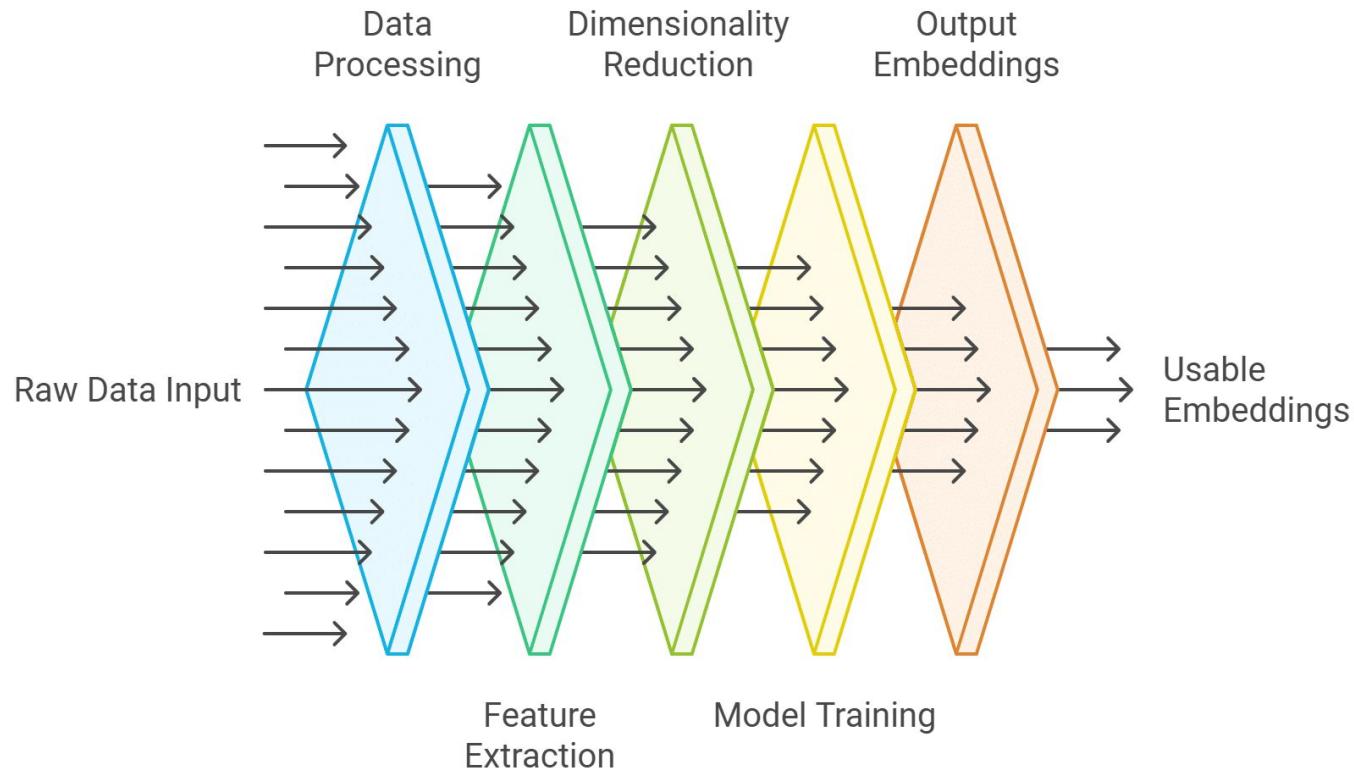


Verb Tense



Country-Capital

Embedding Model Process



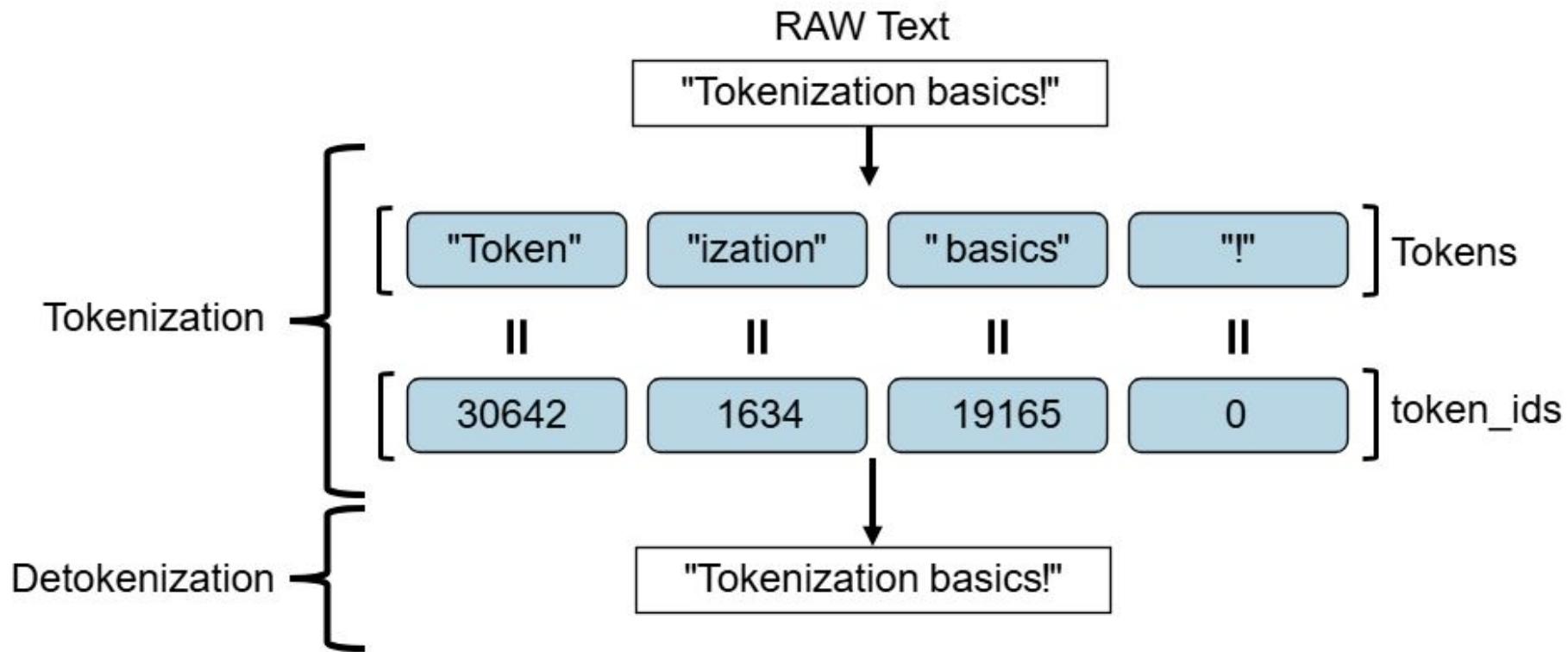
Tokenization

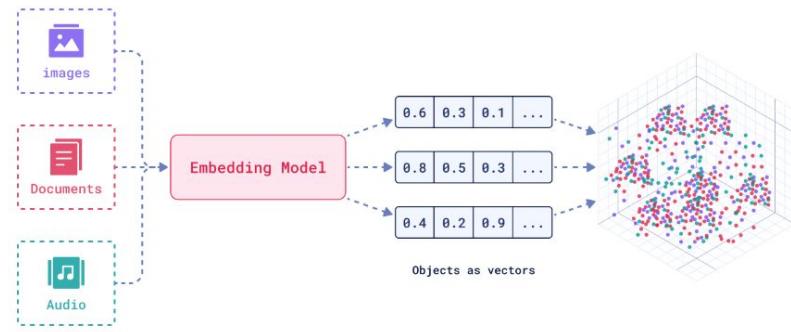
Tiktokenerizer

```
def fizz():
    for i in range(1, 101):
        if i % 5 == 0 and i % 3 == 0:
            print("fizzbuzz")
        elif i % 5 == 0:
            print("buzz")
        elif i % 3 == 0:
            print("fizz")
        else:
            print(i)
```

Token count
77

```
def.fizz():\n....for.i.in.range(1,.101):\n.....if.i.%5==0.and.i.%3==0:\n.....print("fizzbuzz")\n.....elif.i.%5==0:\n.....print("buzz")\n.....elif.i.%3==0:\n.....print("fizz")\n.....else:\n.....print(i)
```





Tiktokenizer

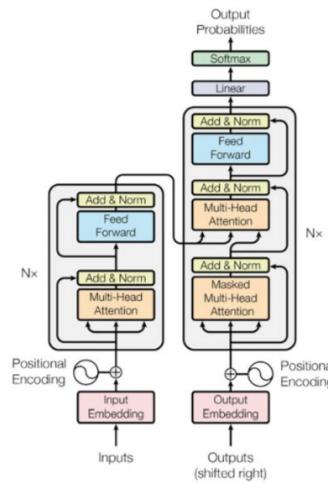
```
def fizz():
    for i in range(1, 101):
        if i % 5 == 0 and i % 3 == 0:
            print("fizzbuzz")
        elif i % 5 == 0:
            print("buzz")
        elif i % 3 == 0:
            print("fizz")
        else:
            print(i)
```

Token count
77

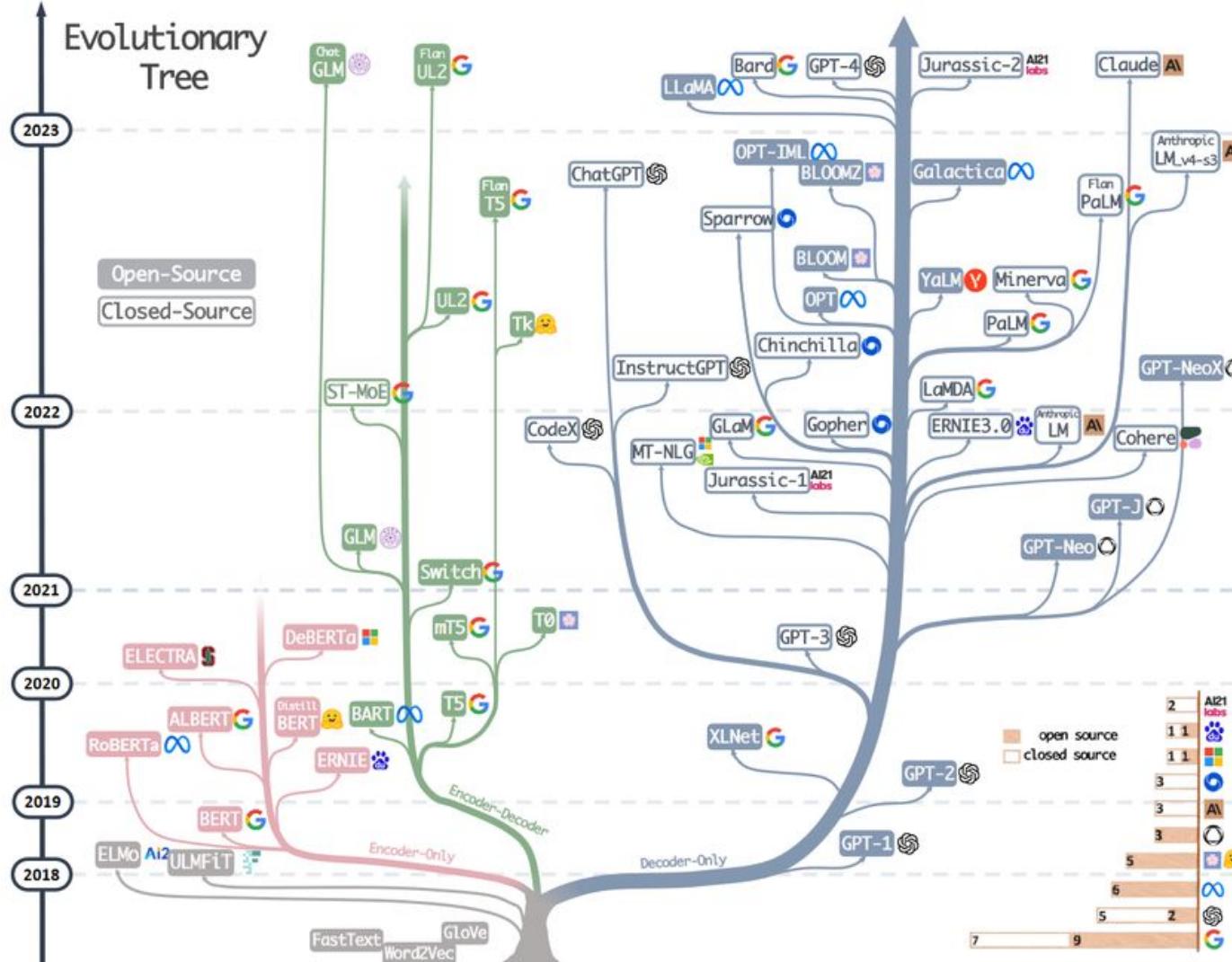
```
def fizz():\n    for i in range(1, 101):\n        if i % 5 == 0 and i % 3 == 0:\n            print("fizzbuzz")\n        elif i % 5 == 0:\n            print("buzz")\n        elif i % 3 == 0:\n            print("fizz")\n        else:\n            print(i)
```

Transformer

Attention Is All You Need

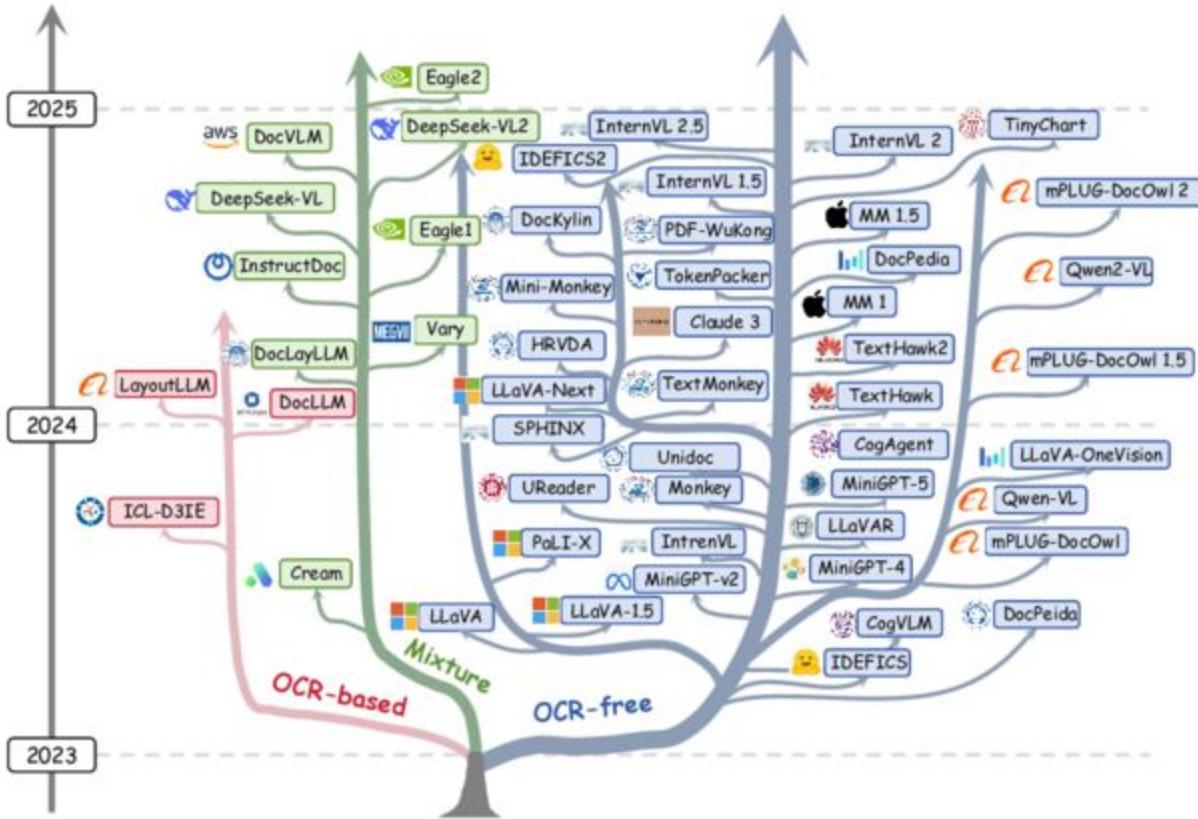


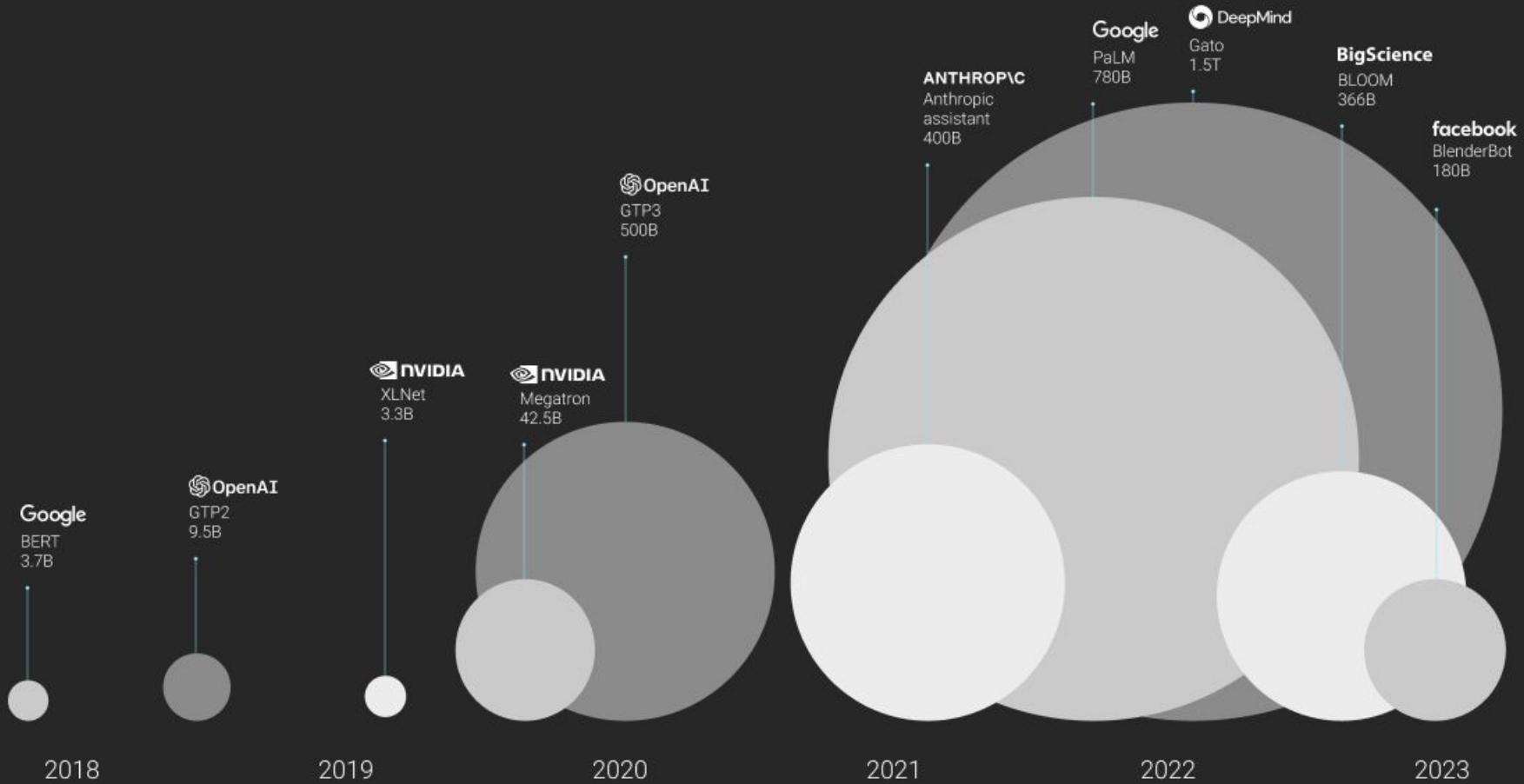
Evolutionary Tree



LLM







SLM: Small (or smart) Language model



r/OpenAI • 2 mo. ago
bgboy089

...

GPT-5 is actually a much smaller model

Discussion

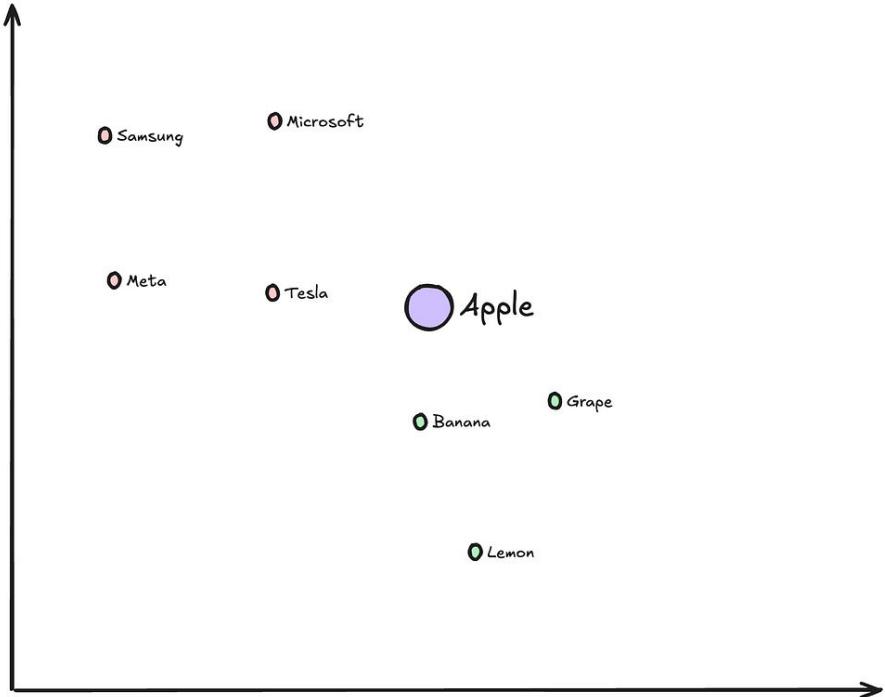
Another sign that GPT-5 is actually a much smaller model: just days ago, OpenAI's O3 model, arguably the best model ever released, was limited to 100 messages per week because they couldn't afford to support higher usage. That's with users paying \$20 a month. Now, after backlash, they've suddenly increased GPT-5's cap from 200 to 3,000 messages per week, something we've only seen with lightweight models like O4 mini.

If GPT-5 were truly the massive model they've been trying to present it as, there's no way OpenAI could afford to give users 3,000 messages when they were struggling to handle just 100 on O3. The economics don't add up. Combined with GPT-5's noticeably faster token output speed, this all strongly suggests GPT-5 is a smaller, likely distilled model, possibly trained on the thinking patterns of O3 or O4, and the knowledge base of 4.5.

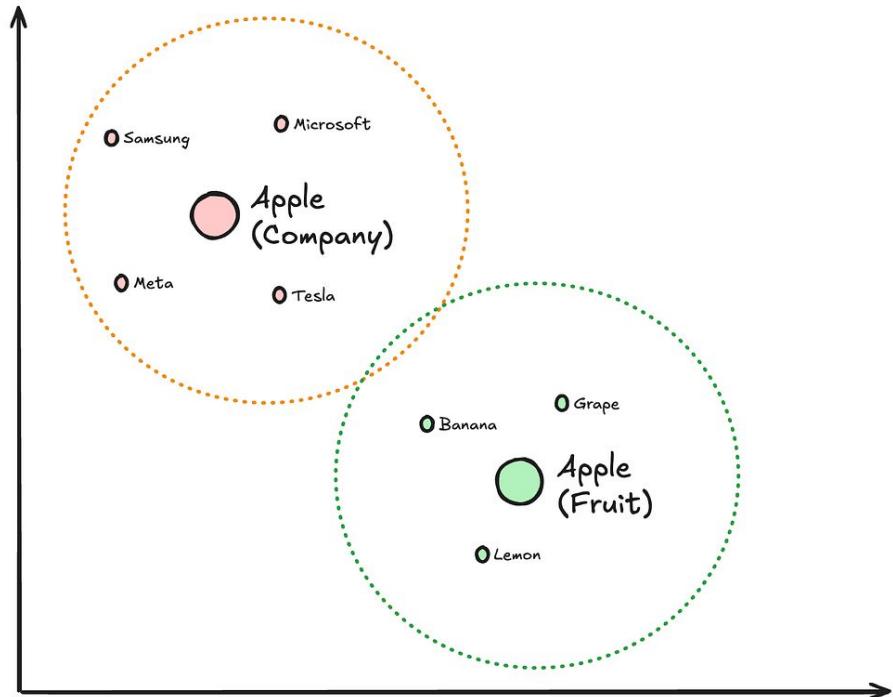
MoE: Mixture of Experts

Contextual tokens

Static Embeddings



Contextual Embeddings





Level 1: input-output



LLM pipelines

RAG



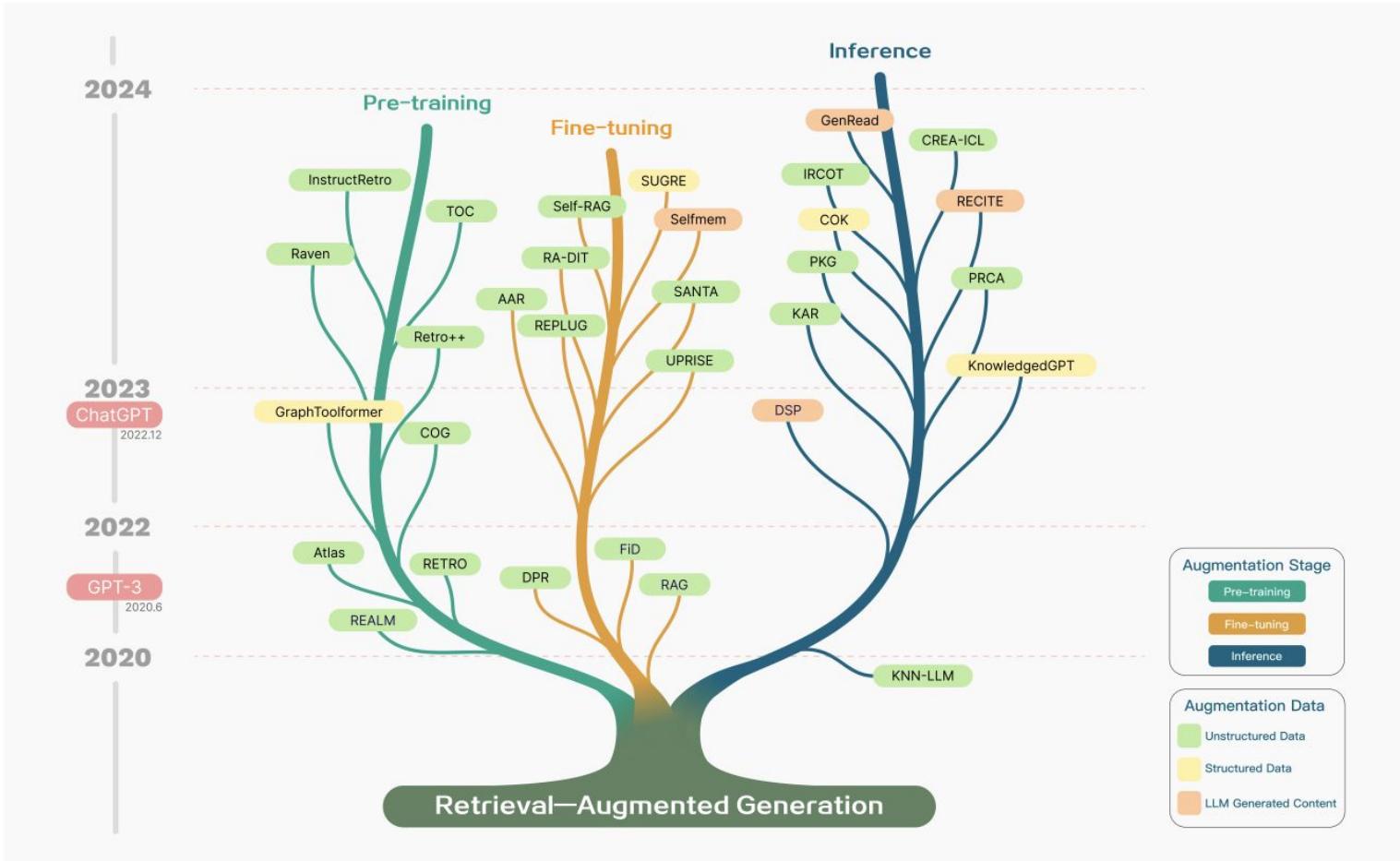
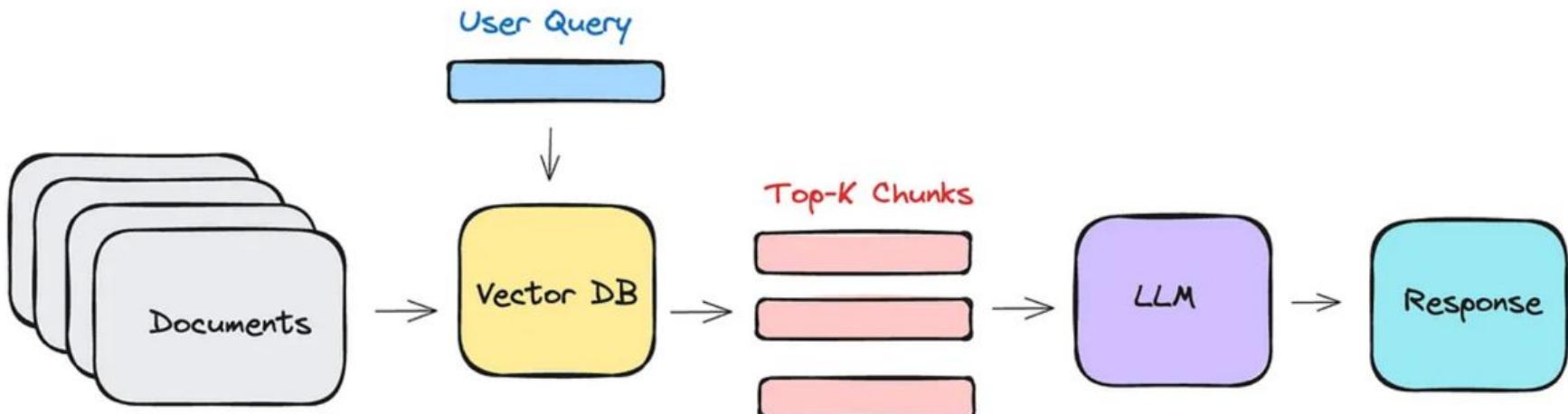


Figure 1: A timeline of existing RAG research. The timeline was established mainly according to the release date.

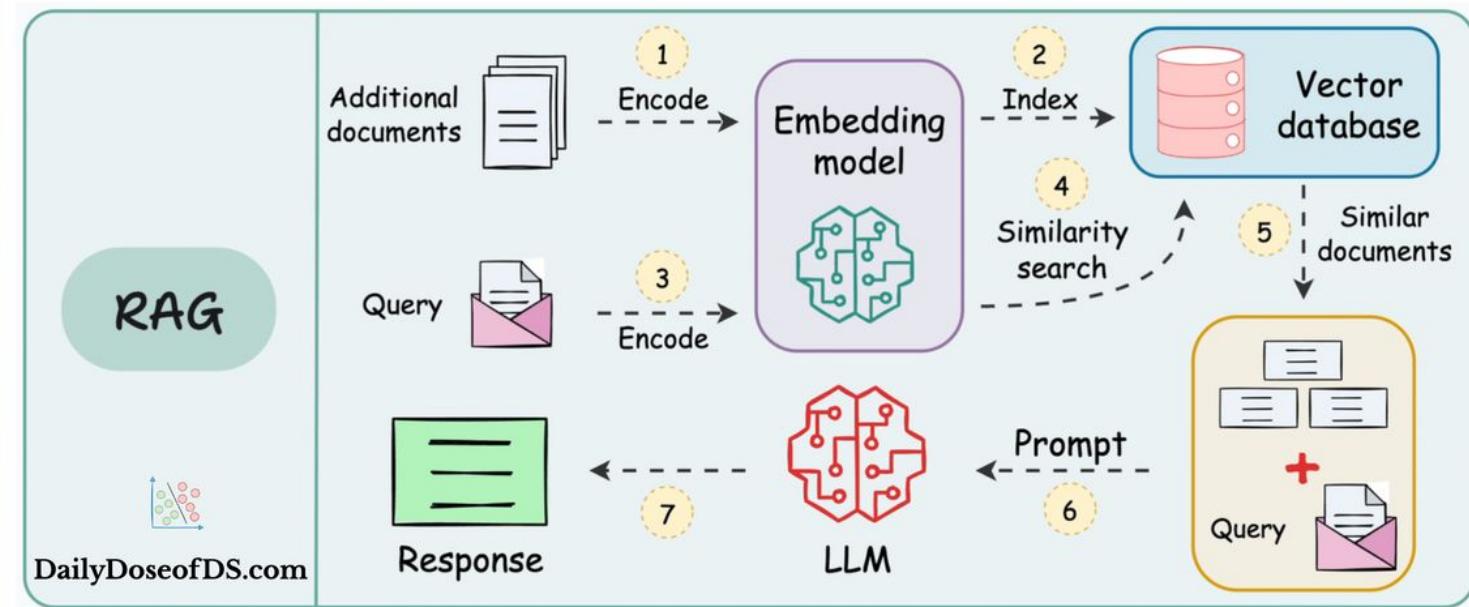
Basic RAG Pipeline



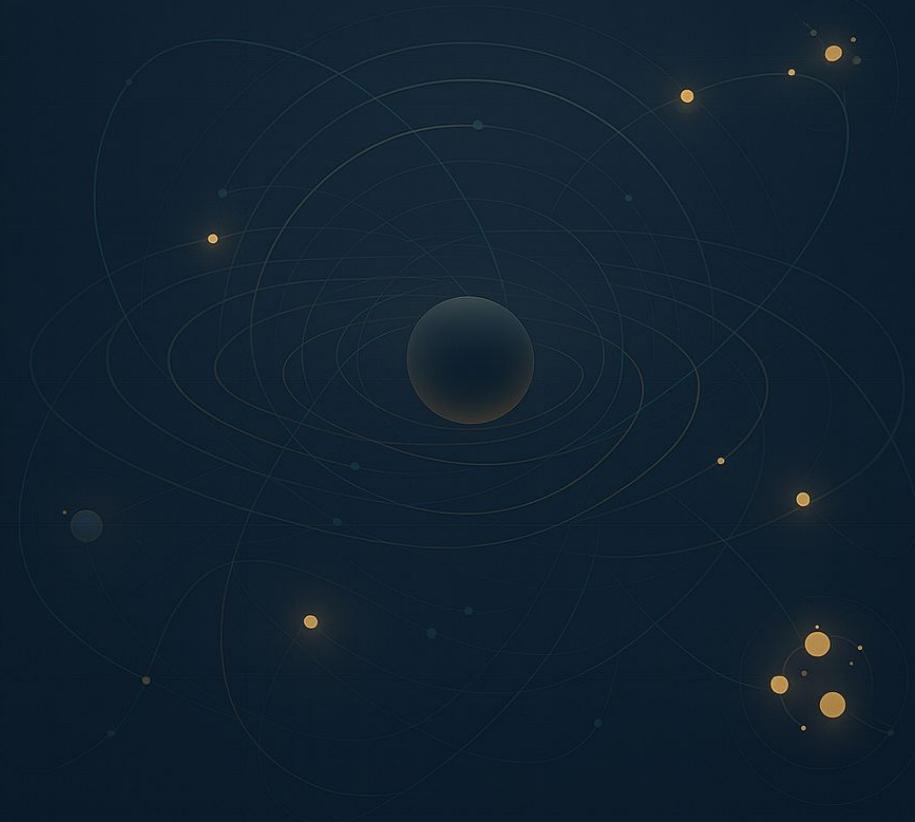
Step 1: Data Indexing

Step 2: Data Retrieval & Generation

16 techniques to supercharge RAG systems



Agentic



Raw LLM

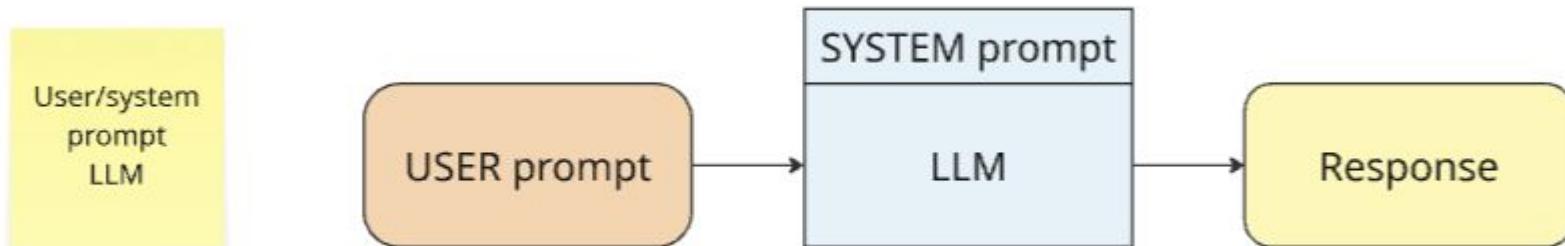
LLM (raw)

→ get an answer based on the trained model



LLM (system prompt)

- get an answer based on the trained model
- Separates request (USER prompt) and model behaviour (POST-training SYSTEM prompt)



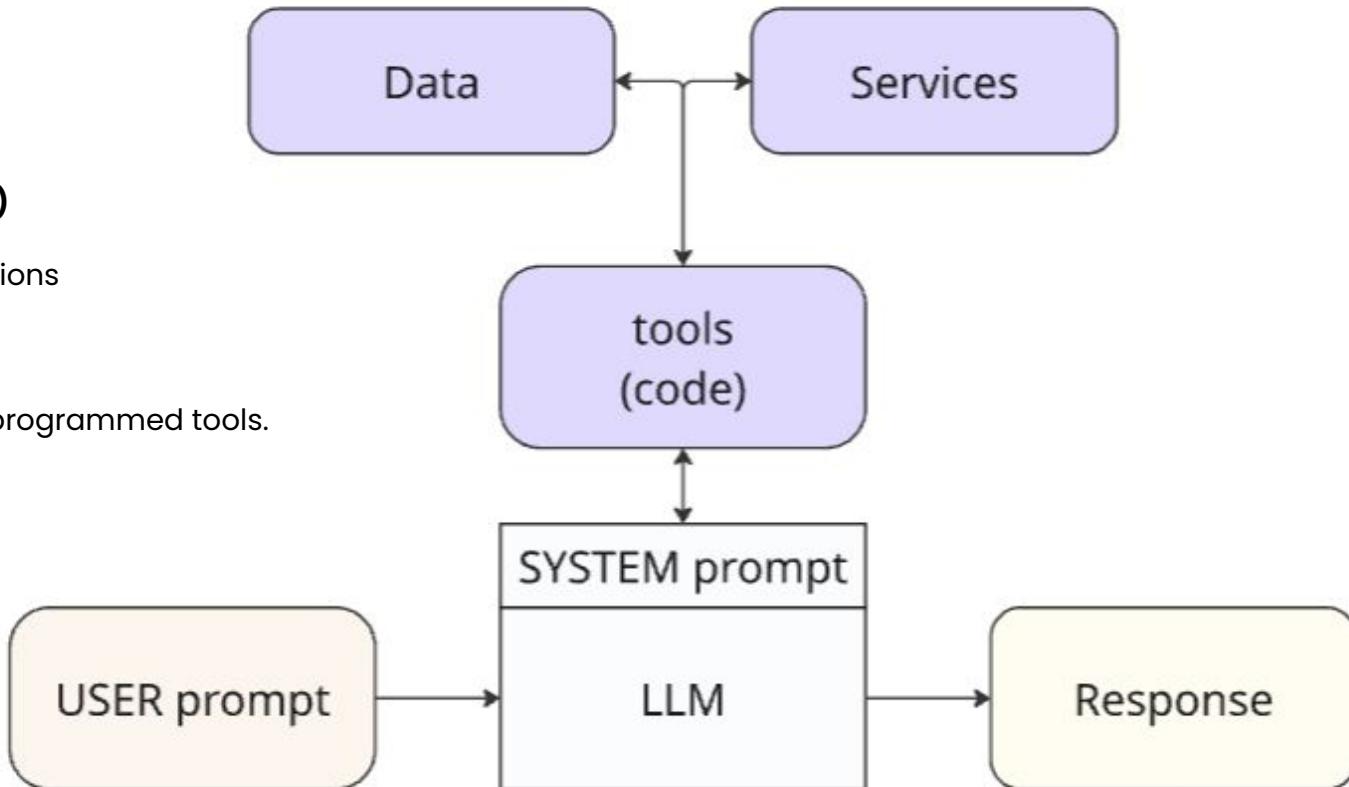
Tools (functions for LLMS)

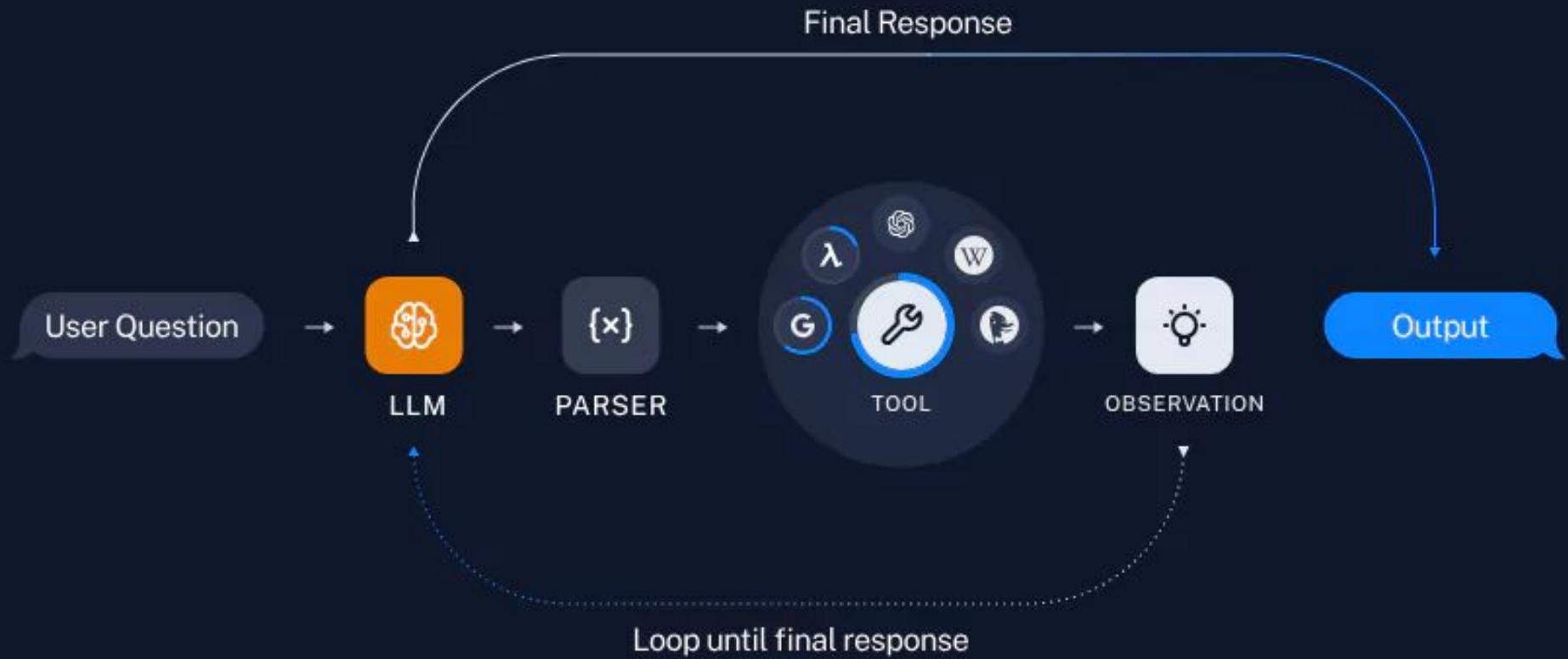
→ Let LLMs run code functions

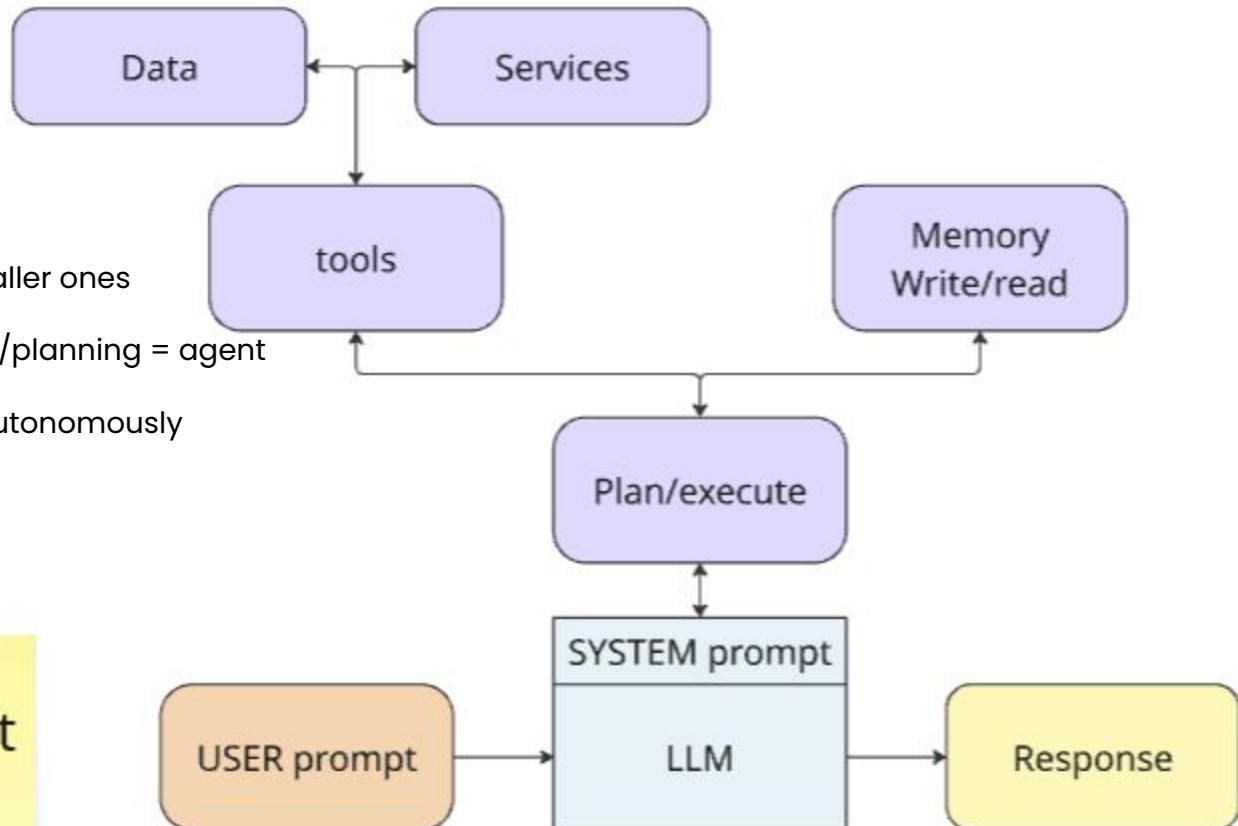
→ tool-augmented LLM:

Gives LLM access to pre-programmed tools.

Tools







Agents (LLMs with tools & planning)

- Can break down big tasks into smaller ones
- LLM + tools + control loop/memory/planning = agent
- Able to complete requests more autonomously

planner &
executor pattern

Agent

Raw LLM



User/system
prompt
LLM

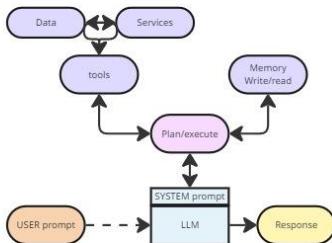
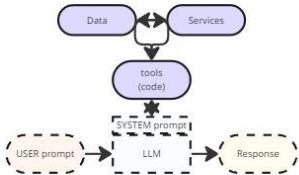
prompt
engineering



LLM
+ Tools

prompt
engineering

Tools



Agent

Context
engineering

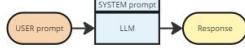
Tools
Subagents

reasoning

MCP

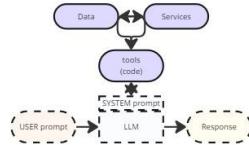


Raw LLM



User/system
prompt
LLM

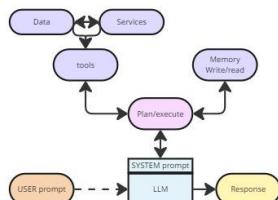
prompt
engineering



LLM
+ Tools

prompt
engineering

Tools



Agent

Context
engineering

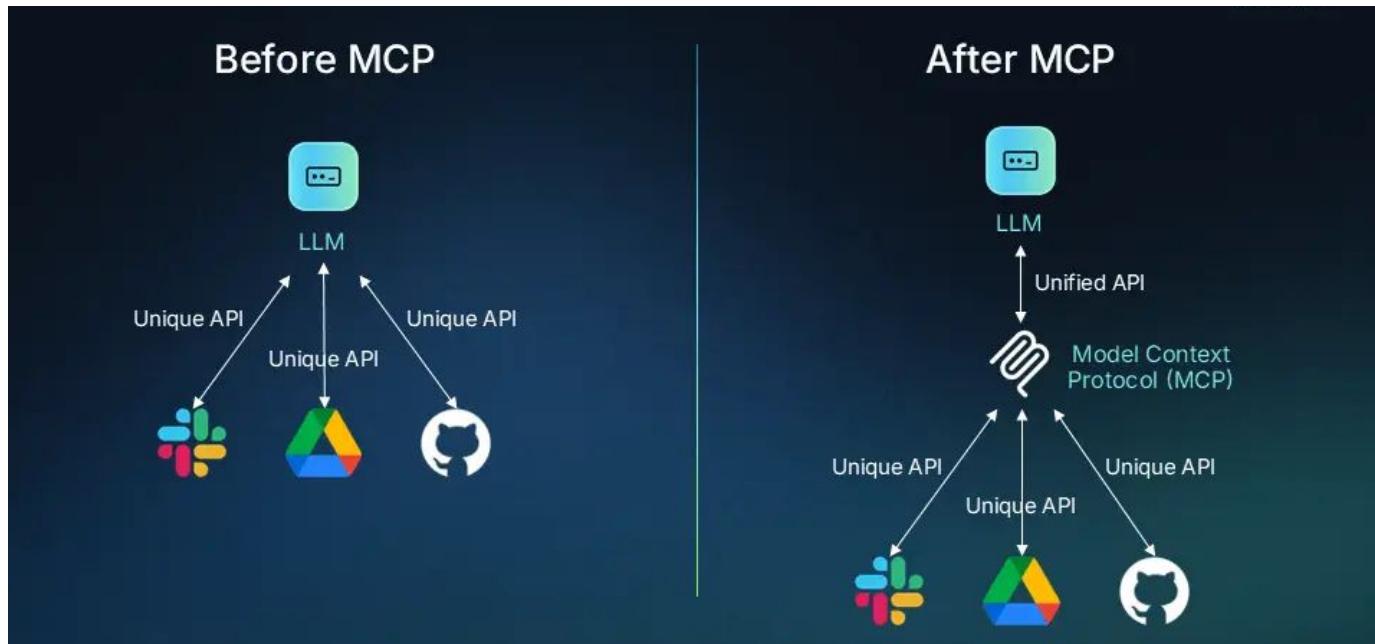
Tools
Subagents

reasoning

MCP

MCP (Model Context Protocol)

- standardized interface for how models access tools, prompts, and resources.
- equip local LLM agent with discoverable external tools.



Server Key Components

MODEL CONTROLLED

Tools

Functions invoked by the model

Retrieve / search

Send a message

Update DB records

APPLICATION CONTROLLED

Resources

Data exposed to the application

Files

Database Records

API Responses

USER CONTROLLED

Prompts

Pre-defined templates for AI interactions

Document Q&A

Transcript Summary

Workflow Automation

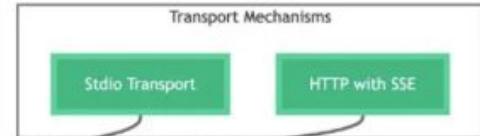
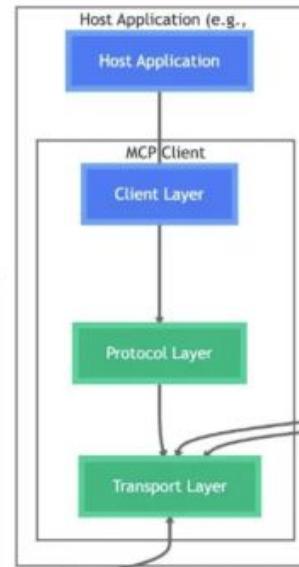
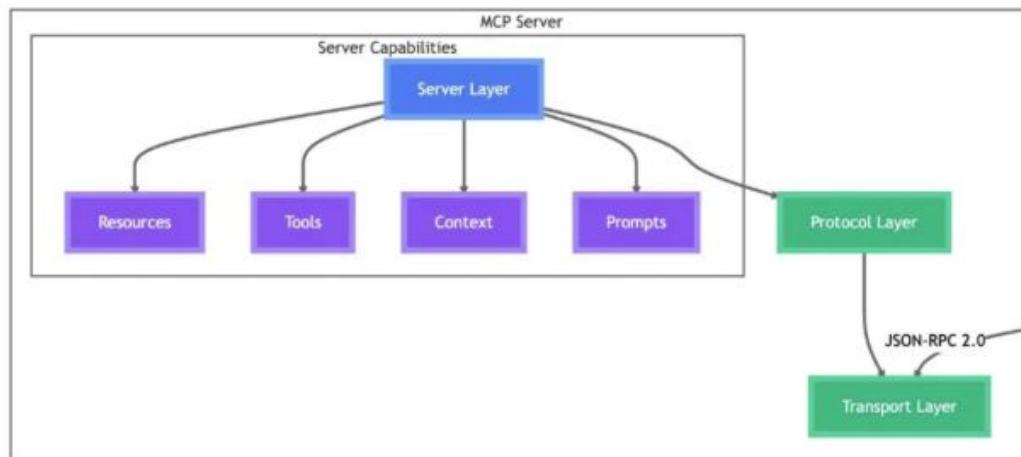
Protocol Layer

Message Framing & Routing
Request/Response Management
Communication Patterns
Versioned Protocol Schemas

Transport Layer

JSON-RPC 2.0 Wire Format
Protocol Message Transmission
Multiple Transport Methods
Extensible Channel System

Core Architecture



Pre-LLM

LLM

RAG

Agentic

Symbolic AI

Attention mechanism

System Prompting

Autonomous Workflow

Transformer Architecture

Hybrid Retrieval

Feature Engineering

Vector Databases & Embeddings

Goal Decomposition

Word Embeddings

Knowledge Distillation

LLM to SQL

Multi-Agent systems

Instruction-Tuning

Semantic search

Tool Use & Function Calling

Meta-Prompting

Tokenization

MoE

Chunking & Indexing

MCP

Sequence Models (RNN/LSTM)

Context Window

Context Compression / Engineering

RLHF

Human-centered AI

Human-Computer Interaction (HCI)

NLP (Natural Language Processing)

Prompt engineering

AGI

Whats next?

Agentic economy



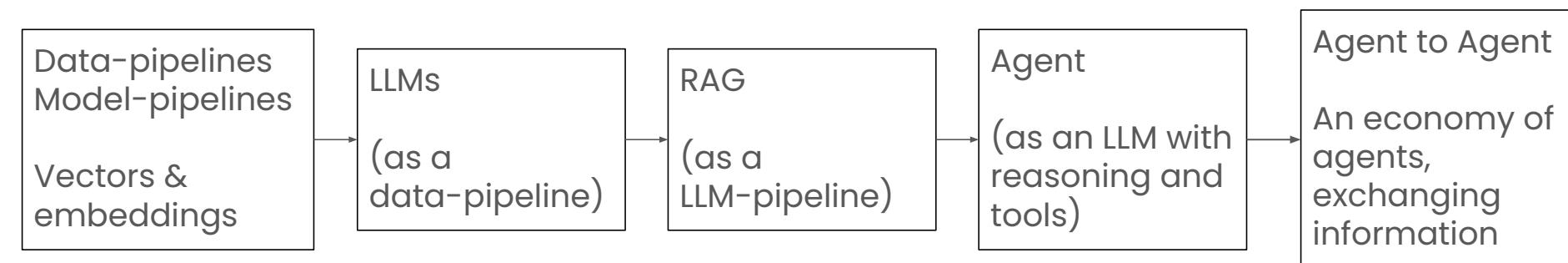
Pre-LLM

LLM

RAG

Agentic

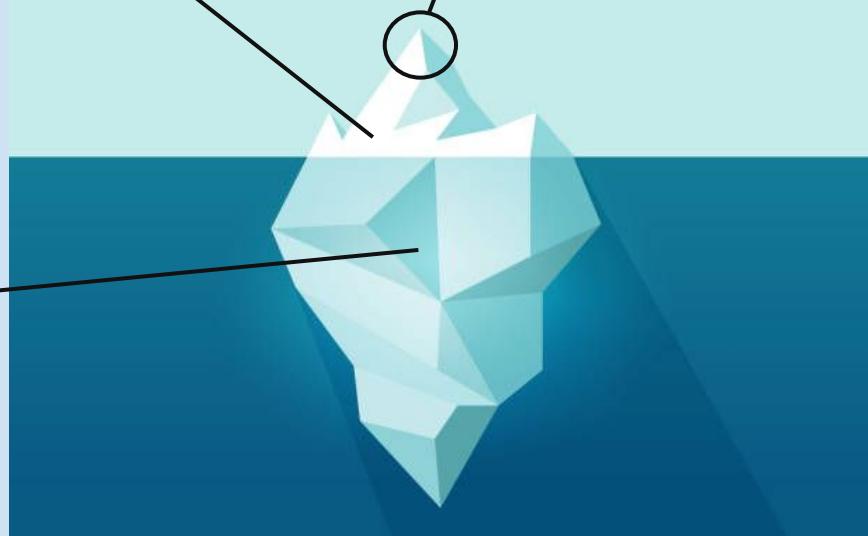
Agentic economy



Techical

Societal,
mental/cognitive,
political,
Economic,
Cultural,
Academic,
Environmental,
Ethical,
Educational

This presentation



"Crash course in LLMs"

End / Q&a

Prompt Engineering Guide - <https://www.promptingguide.ai/>

Daily Dose of Data Science - <https://www.dailydoseofds.com/>

Model Context Protocol - <https://modelcontextprotocol.io/>

Costs

Effects

Costs

Effects

Costs Effects Costs Effects Costs

Costs Effects Costs Effects Costs