

AI/ML Overview (2025)

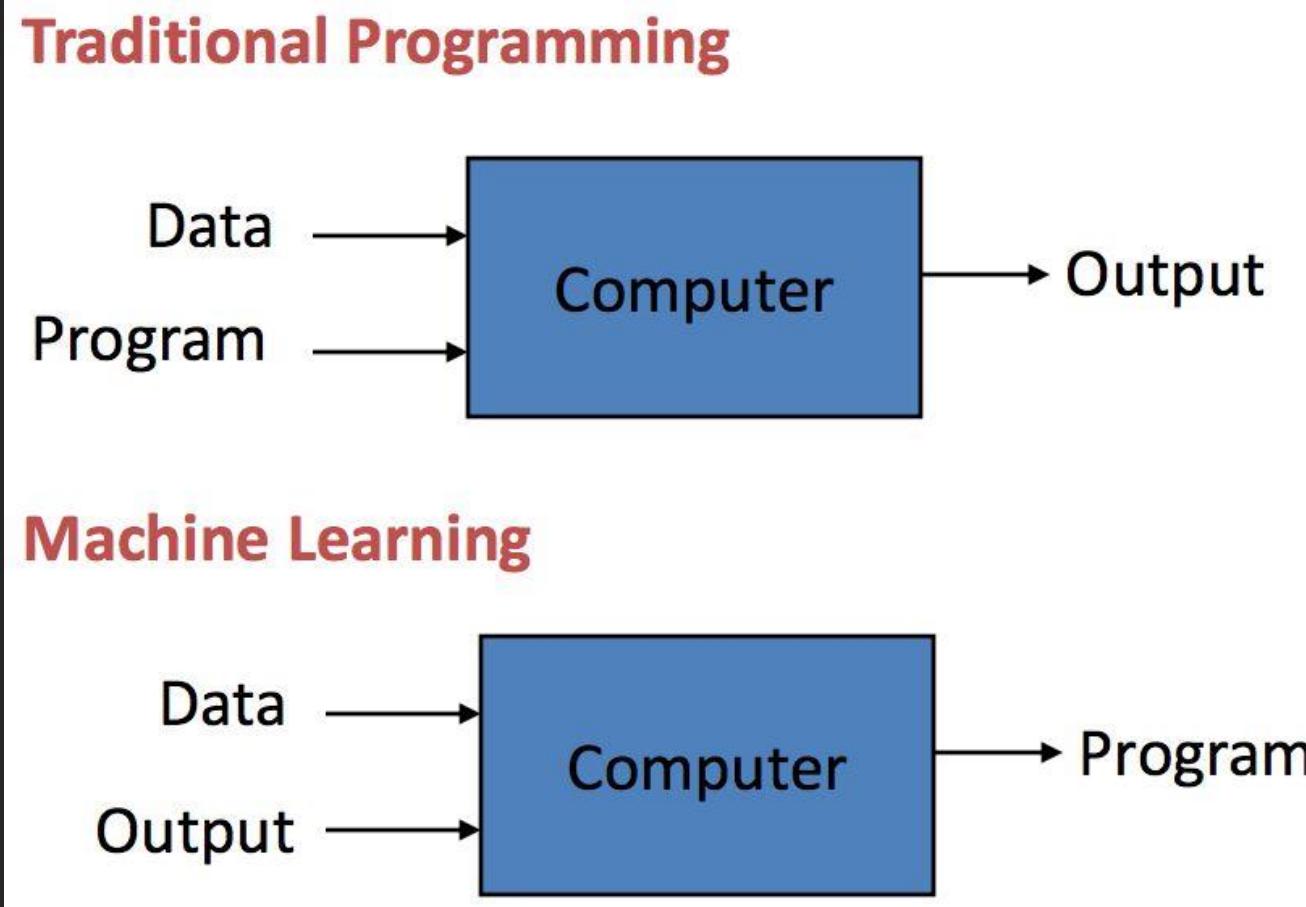
By: Pranpaveen Laykaviriyakul

Agenda

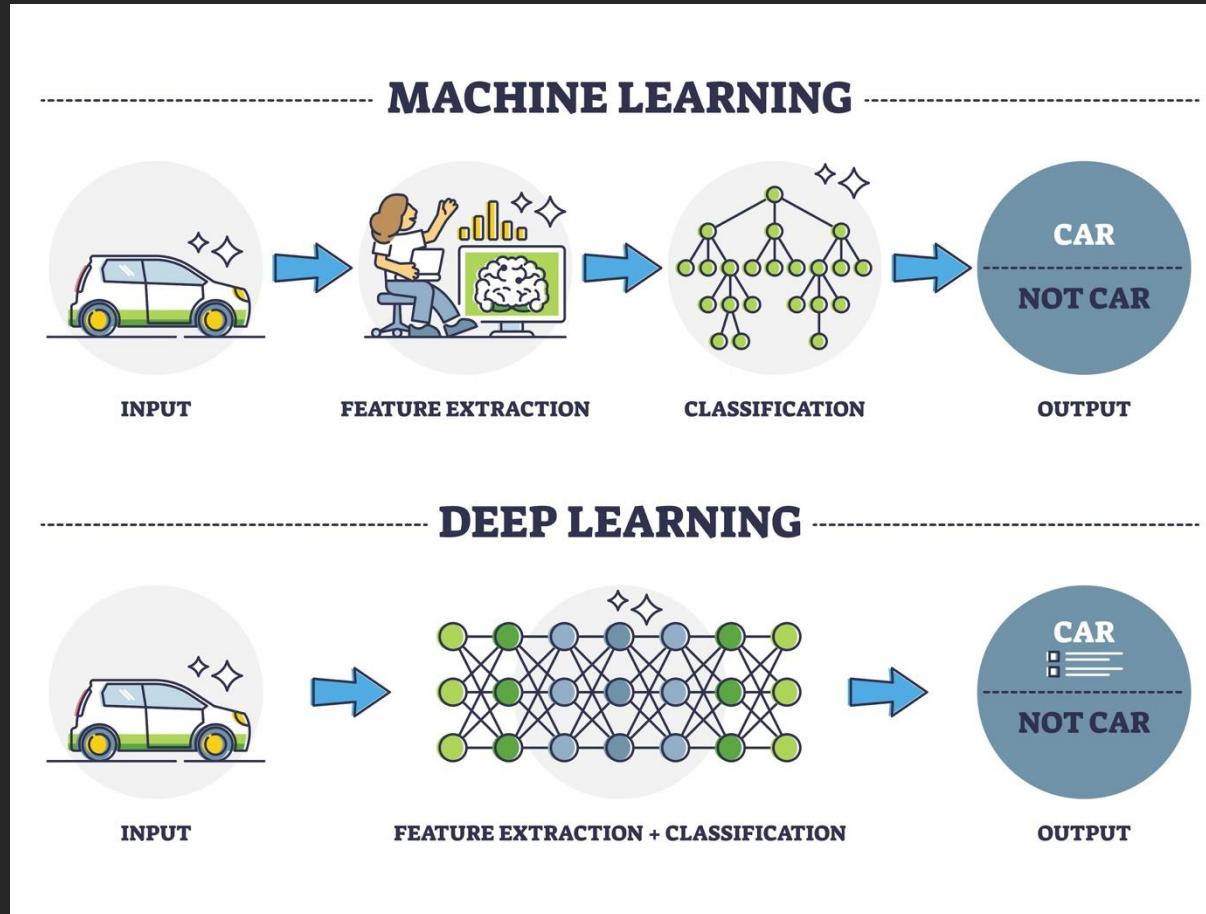
- Introduction to AI/ML
- Deep Dive into Machine Learning
- Tools / Platforms
- Current Trends in AI/ML
- ML related Careers

Introduction to AI/ML

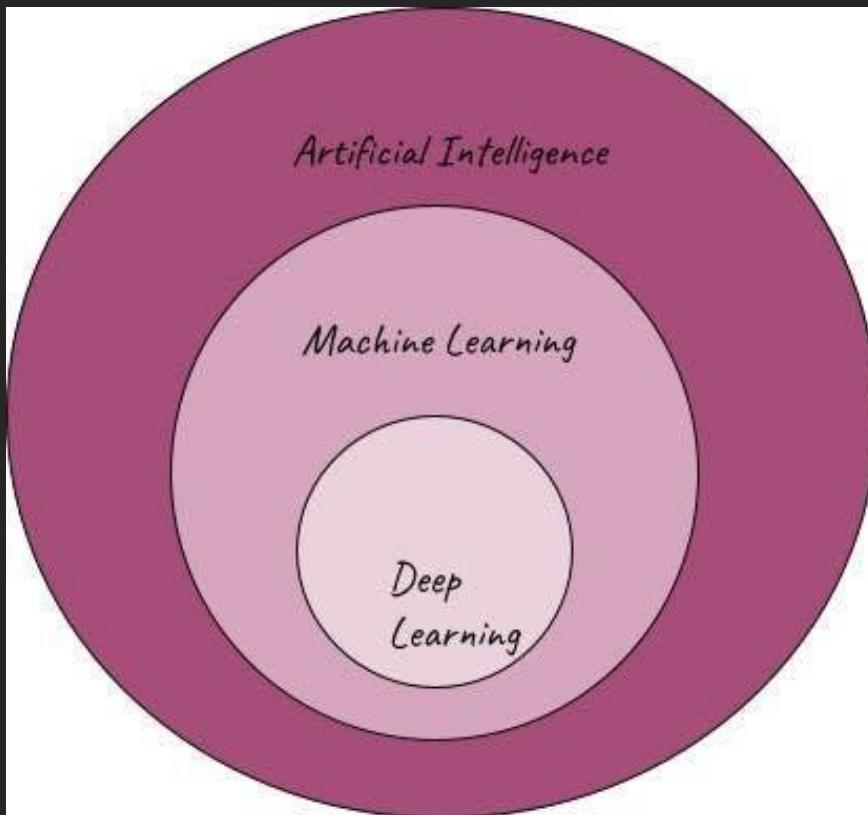
Machine Learning



ML vs DL



AI, ML, and DL



AI: Technique which enables computer to mimic the intelligence or behavior of human.

Machine learning (ML): Subset of AI techniques which enable computer to improve from experiences without being explicitly programmed.

Deep learning (DL) Subset of ML which uses deep artificial neural networks as a model and automatically build a hierarchy of data representation .

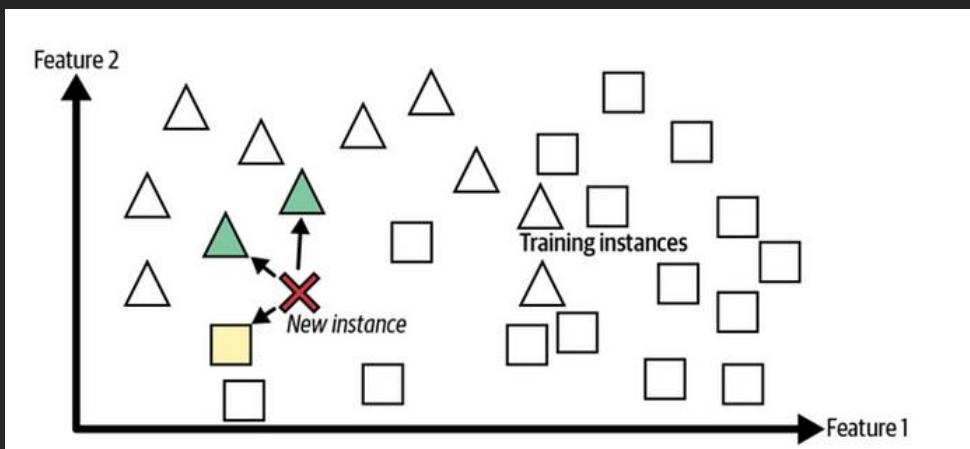
ML/AI capabilities

SUMMARY OF ML/AI CAPABILITIES				
CAPABILITIES	USE CASES			
	PERCEPTION (interpreting the world)	VISION understanding images	AUDIO audio recognition	SPEECH <ul style="list-style-type: none">text-to-speechSpeech-to-text conversions NATURAL LANGUAGE understanding & generating text
	COGNITION (reasoning on top of data)	REGRESSION <ul style="list-style-type: none">predicting a numerical value CLASSIFICATION <ul style="list-style-type: none">predicting a category for a data point PATTERN RECOGNITION <ul style="list-style-type: none">identifying relevant insights on data	PLANNING <ul style="list-style-type: none">determining the best sequence of steps for a goal OPTIMISATION <ul style="list-style-type: none">identifying the most optimal parameters. RECOMMENDATION <ul style="list-style-type: none">predicting user's preferences	
LEARNING (types of ML/AI)	SUPERVISED <ul style="list-style-type: none">learning on labelled data pairs: (input, output) UNSUPERVISED <ul style="list-style-type: none">inferring hidden structures in an unlabelled data REINFORCEMENT LEARNING <ul style="list-style-type: none">learning by experimentingmaximizing reward			

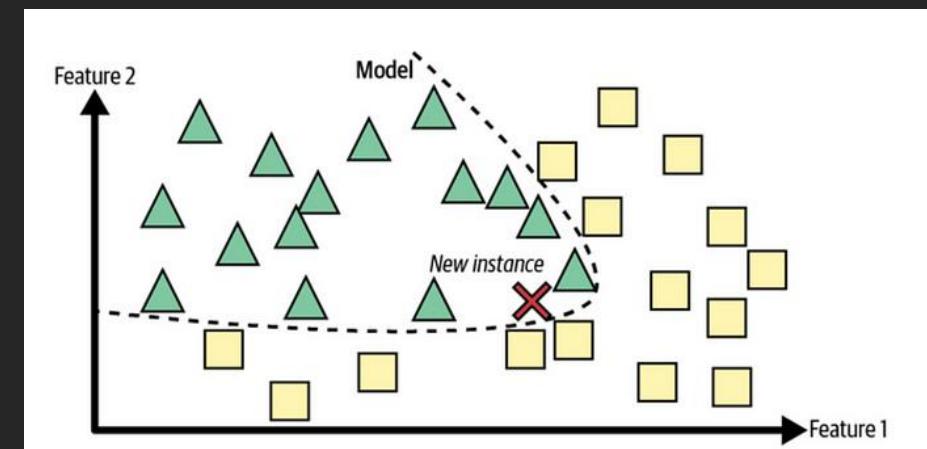
Types of Machine Learning

Instance-based vs Model-based learning

Instance-based

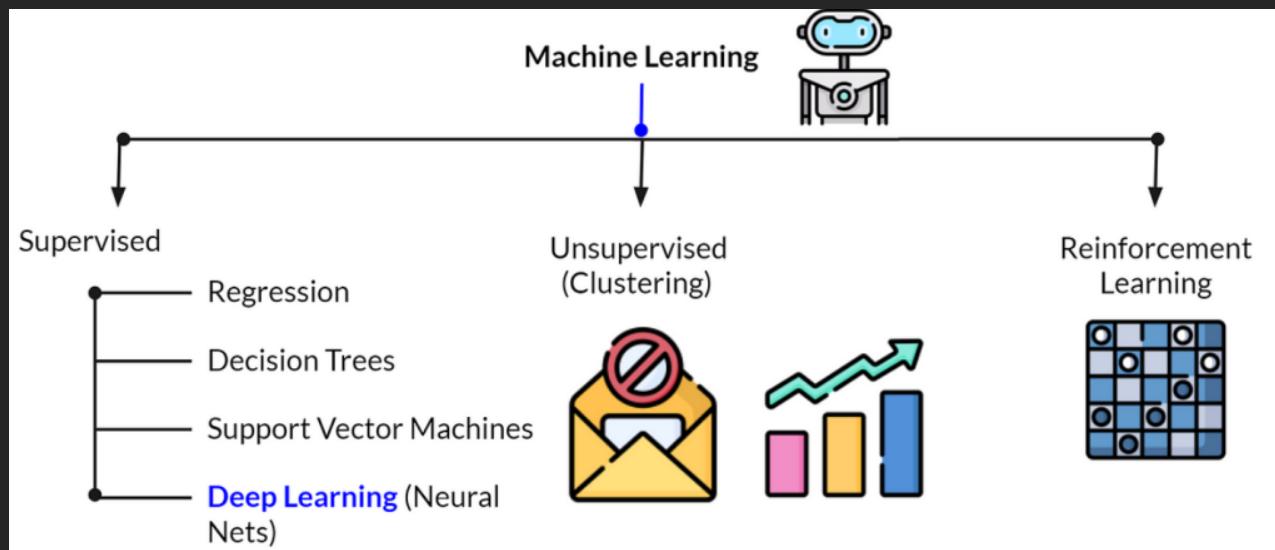
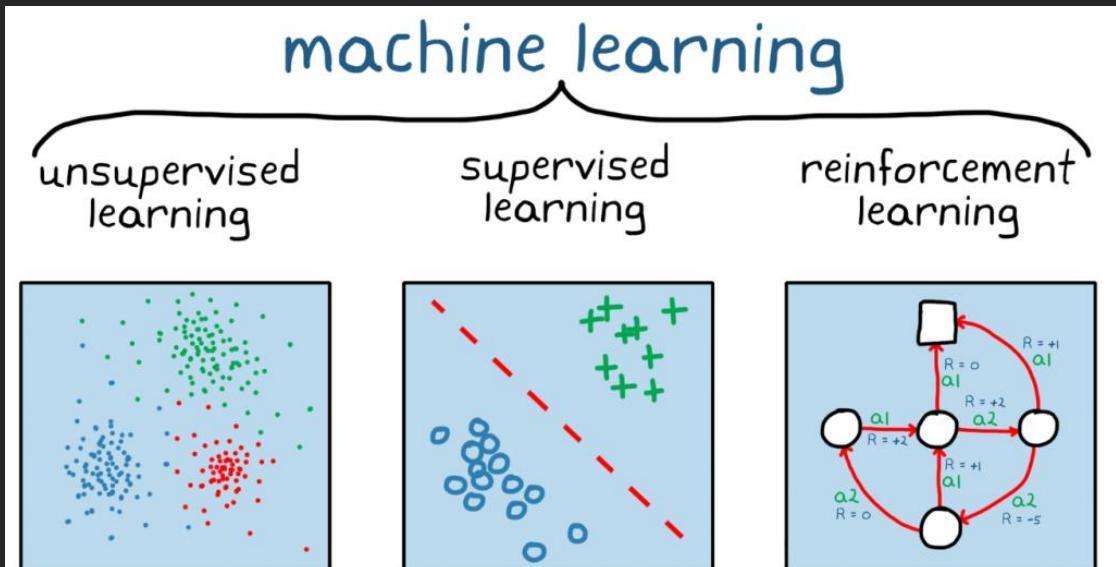


Model-based



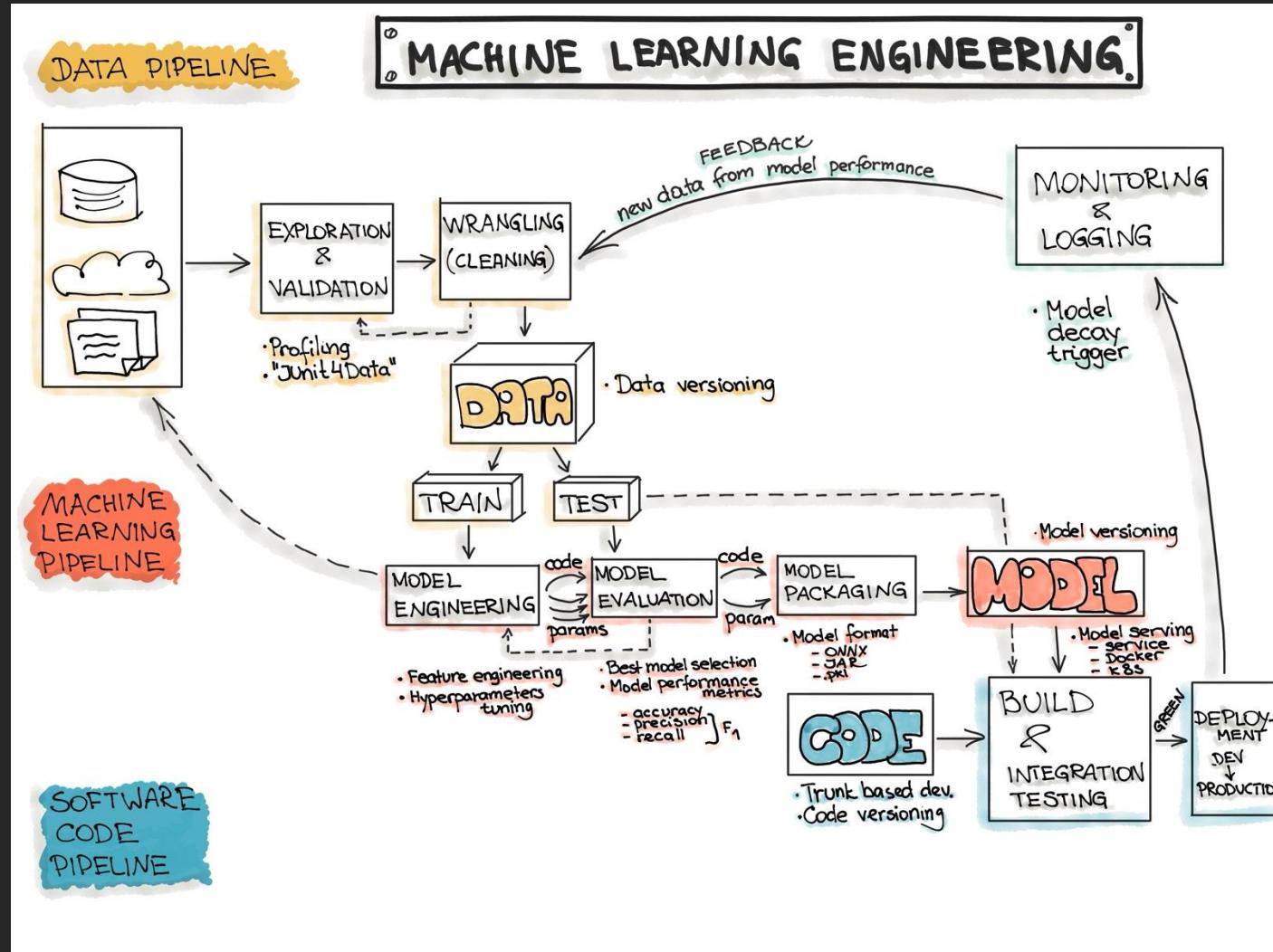
Types of Machine Learning

Learn style



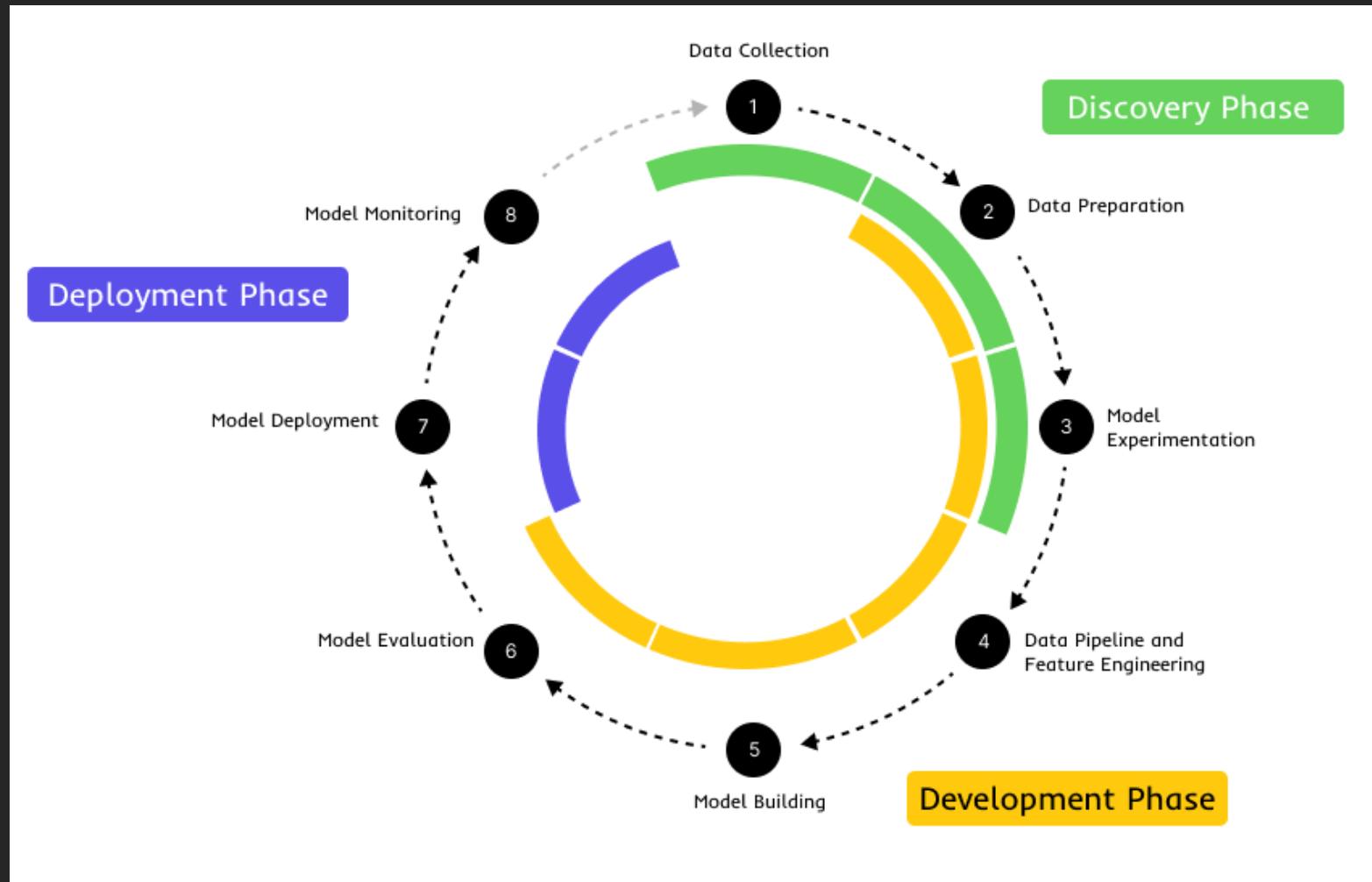
Deep Dive into Machine Learning

Machine Learning Workflow

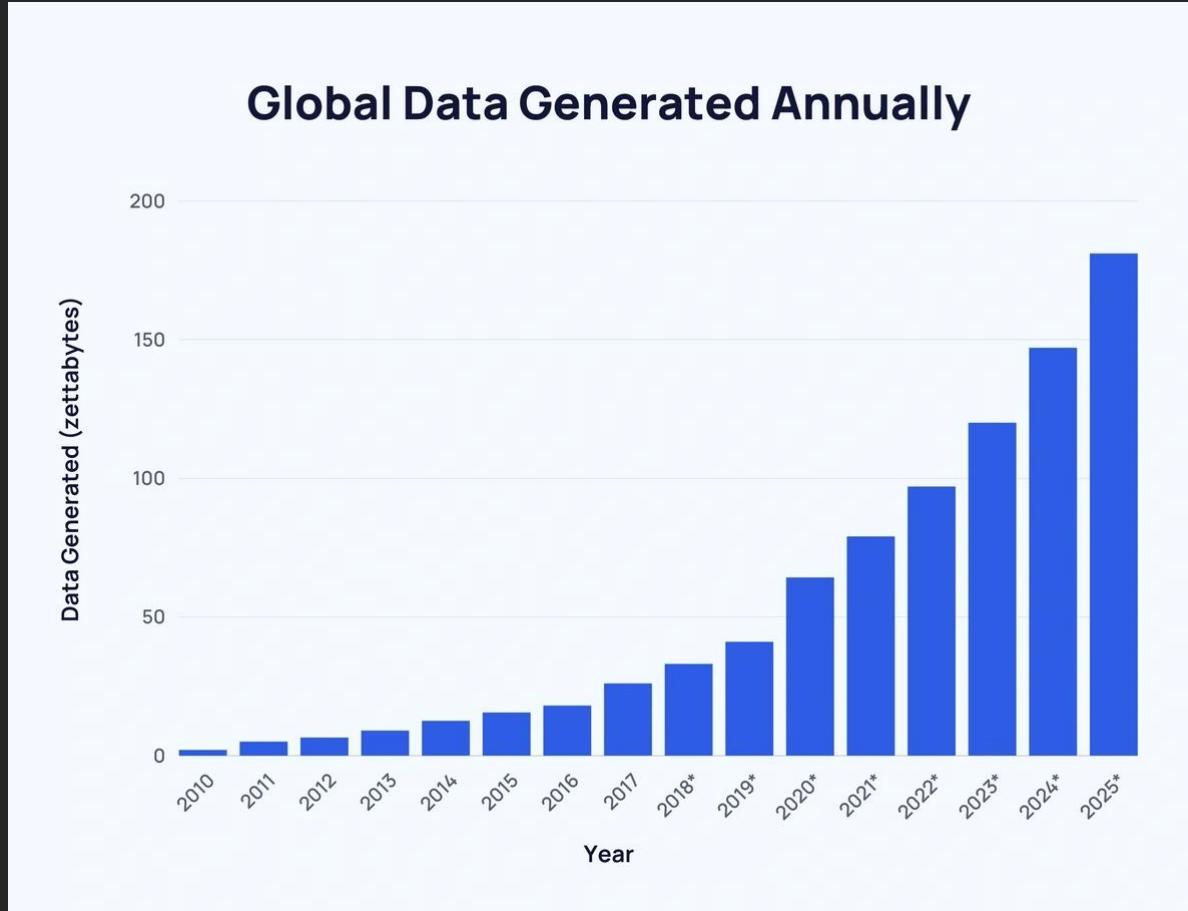


- Data Engineering
 - data acquisition & data preparation
- ML Model Engineering
 - ML model training & serving, and
- Code Engineering
 - integrating ML model into the final product.

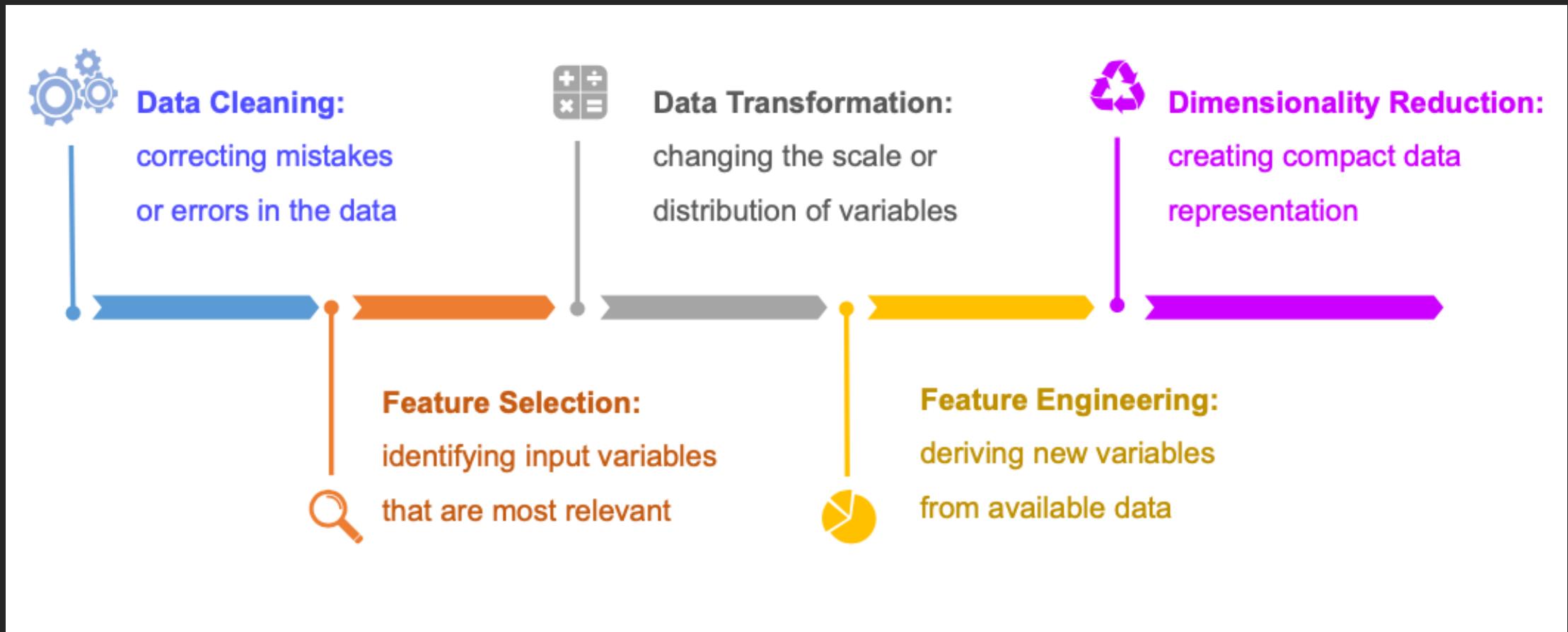
Machine Learning Lifecycle



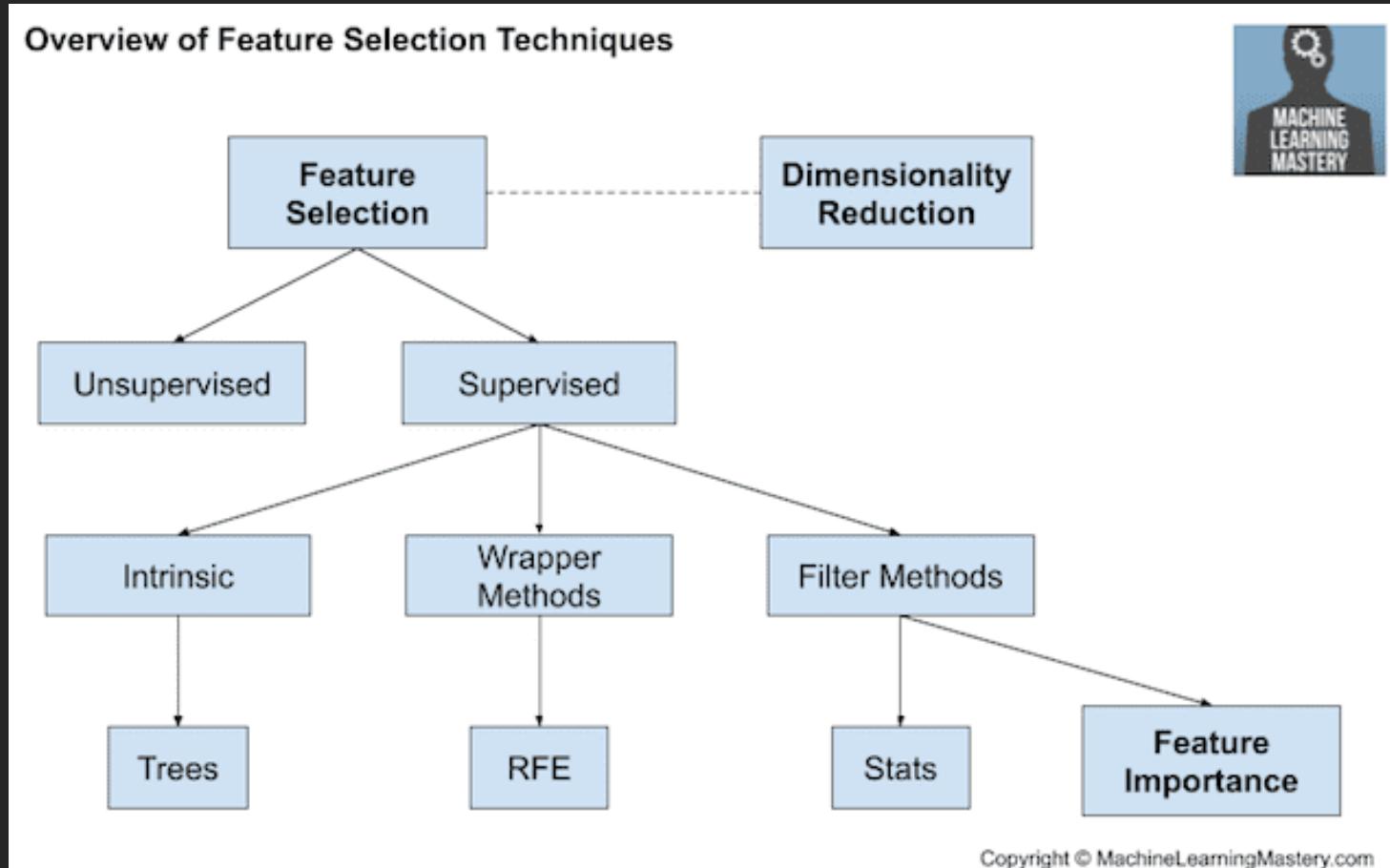
Data Is the New Gold



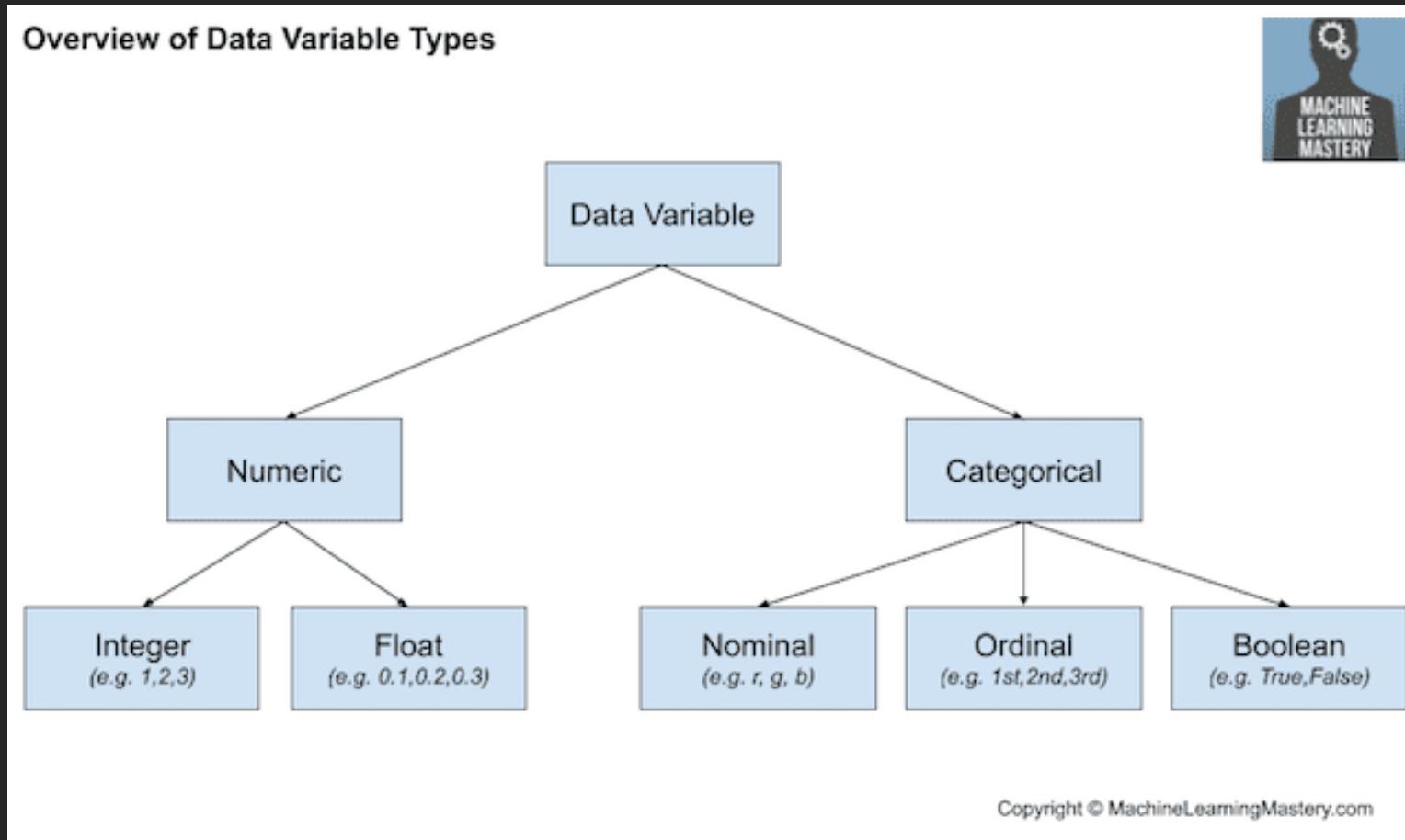
Data preparation



Feature Selection



Data Transformation



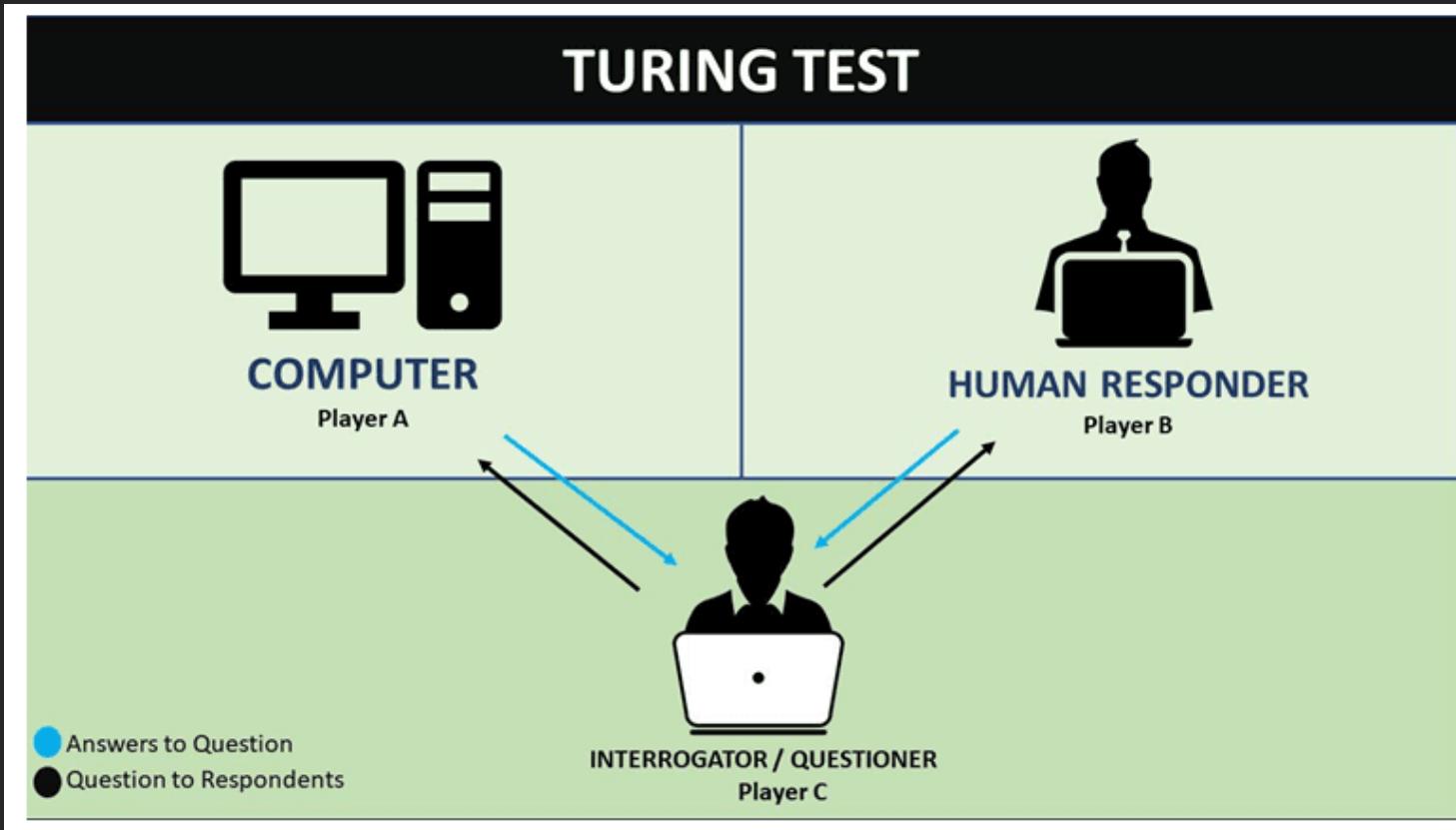
<https://machinelearningmastery.com/feature-selection-with-real-and-categorical-data/>

Features Engineering

We should clarify how should the input data be represented.

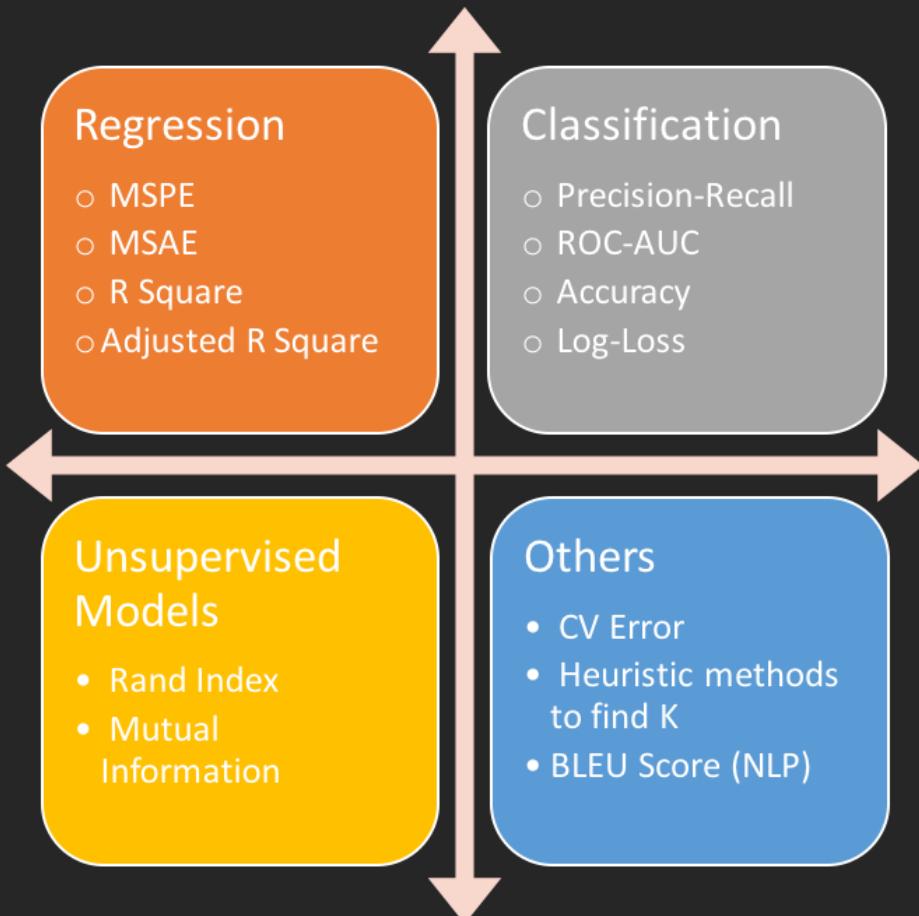
- How do we extract features from raw sources?
- Consider to include domain experts to specify what data aspects are most important for the particular ML task.

Evaluation

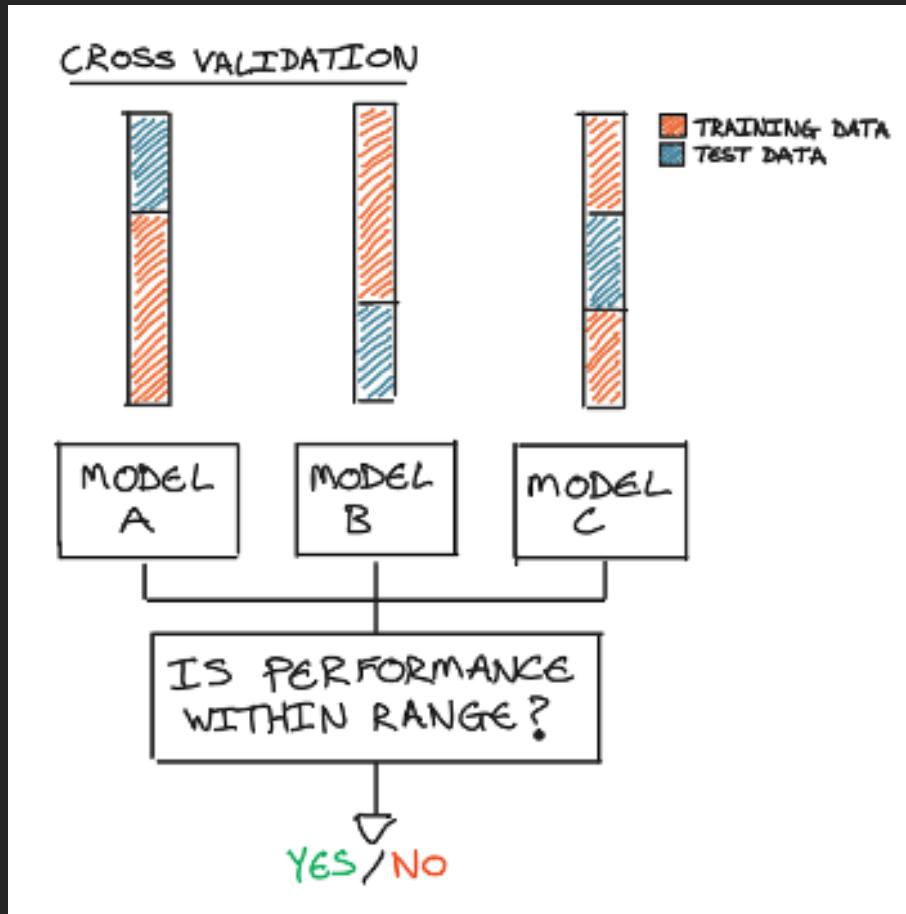


- The Turing Test, proposed by Alan Turing in 1950
- Assesses a machine's ability to exhibit human-like intelligence
- If a human evaluator cannot distinguish between the machine and a human based on their responses, the machine passes the test

Evaluation - Metrics



Cross-validation



Experiments Tracking - Test

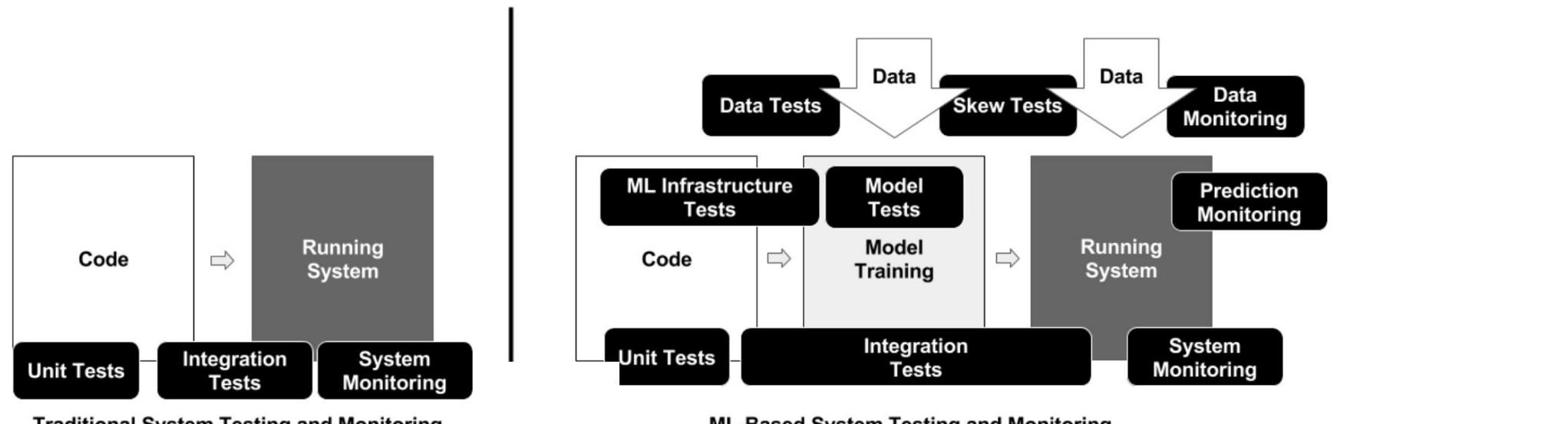
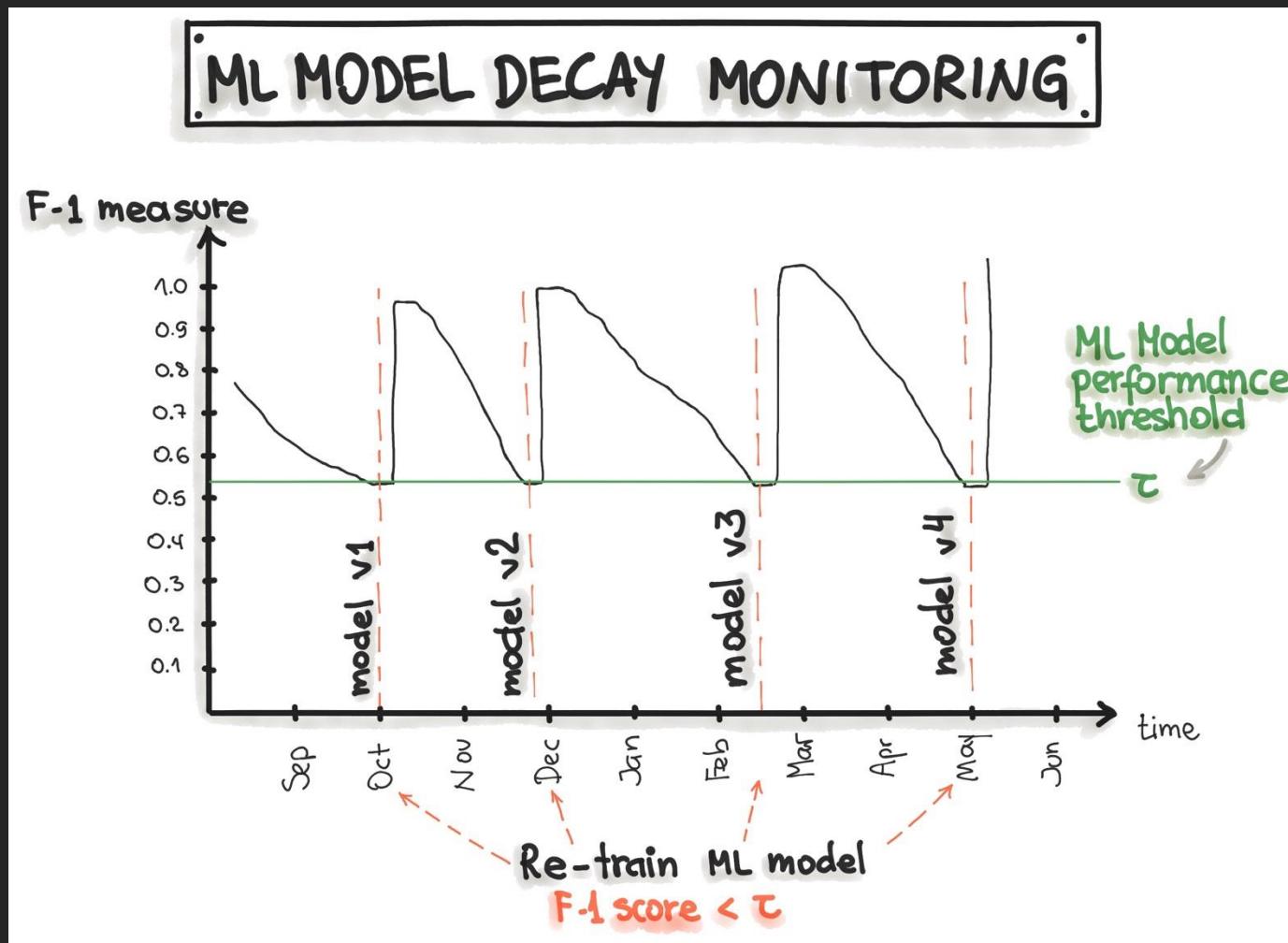
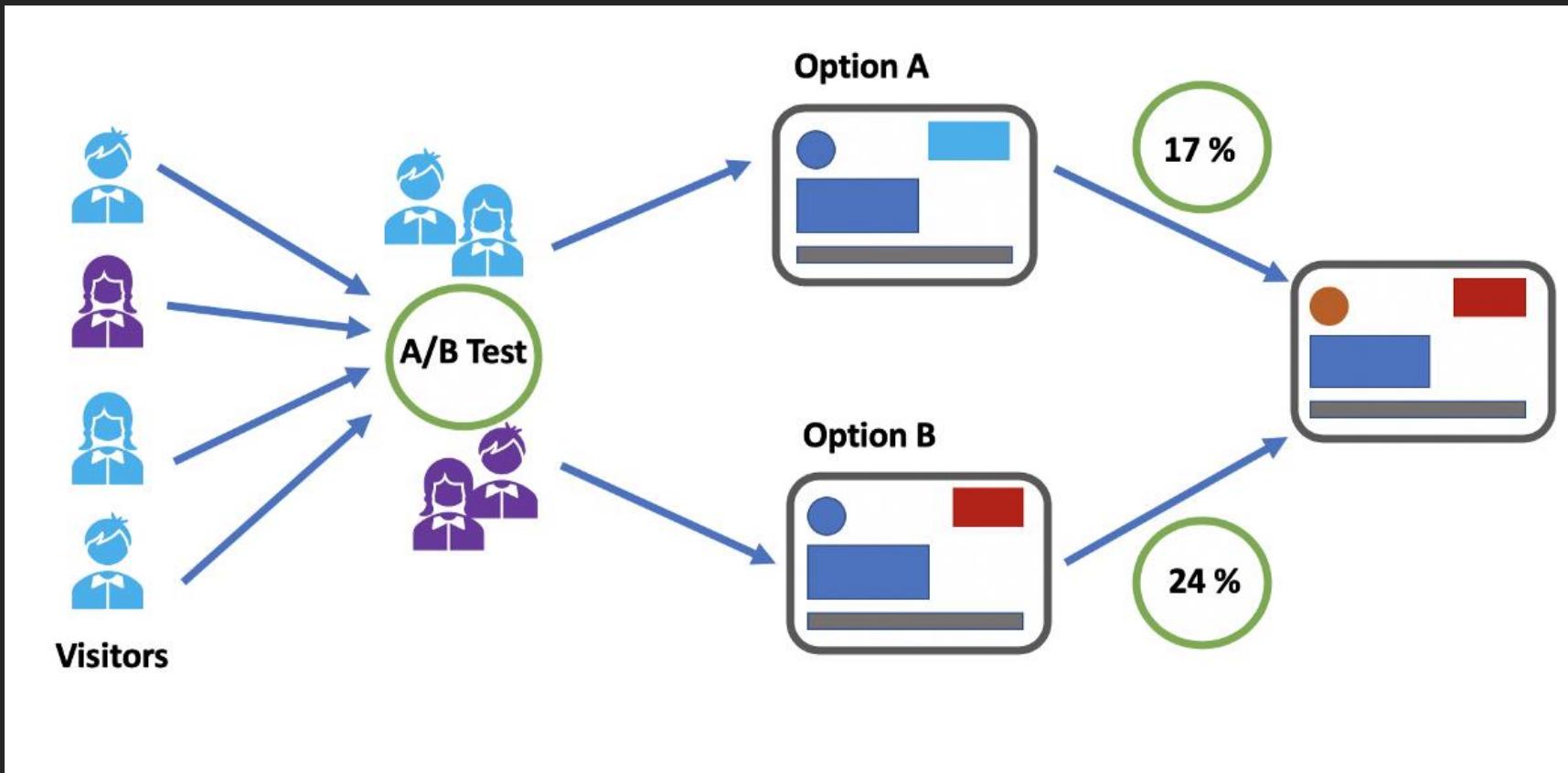


Figure 1. **ML Systems Require Extensive Testing and Monitoring.** The key consideration is that unlike a manually coded system (left), ML-based system behavior is not easily specified in advance. This behavior depends on dynamic qualities of the data, and on various model configuration choices.

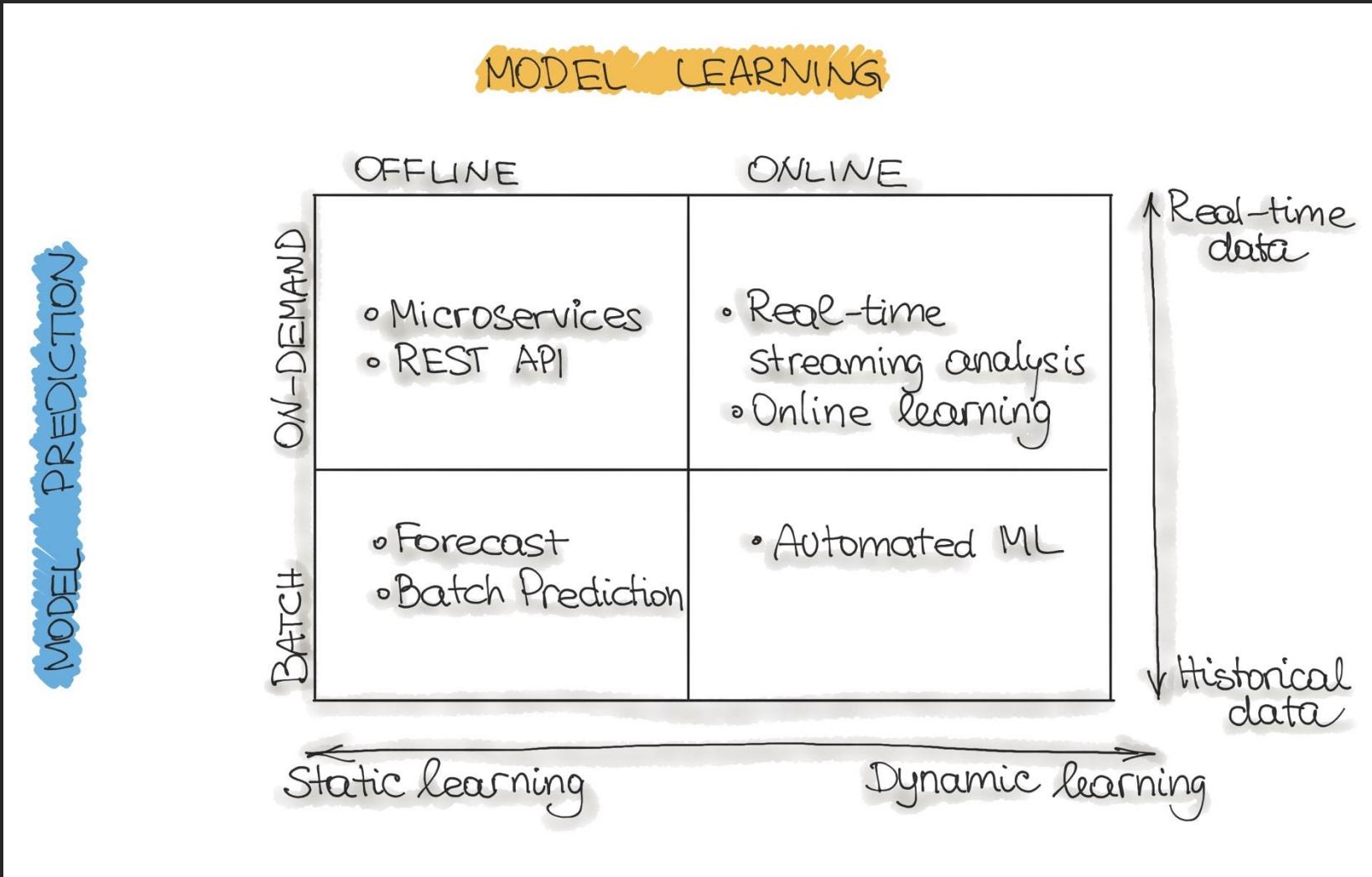
Experiments Tracking - Monitoring



Experiments Tracking - A/B Testing



Model Serving Patterns



Tools / Platforms

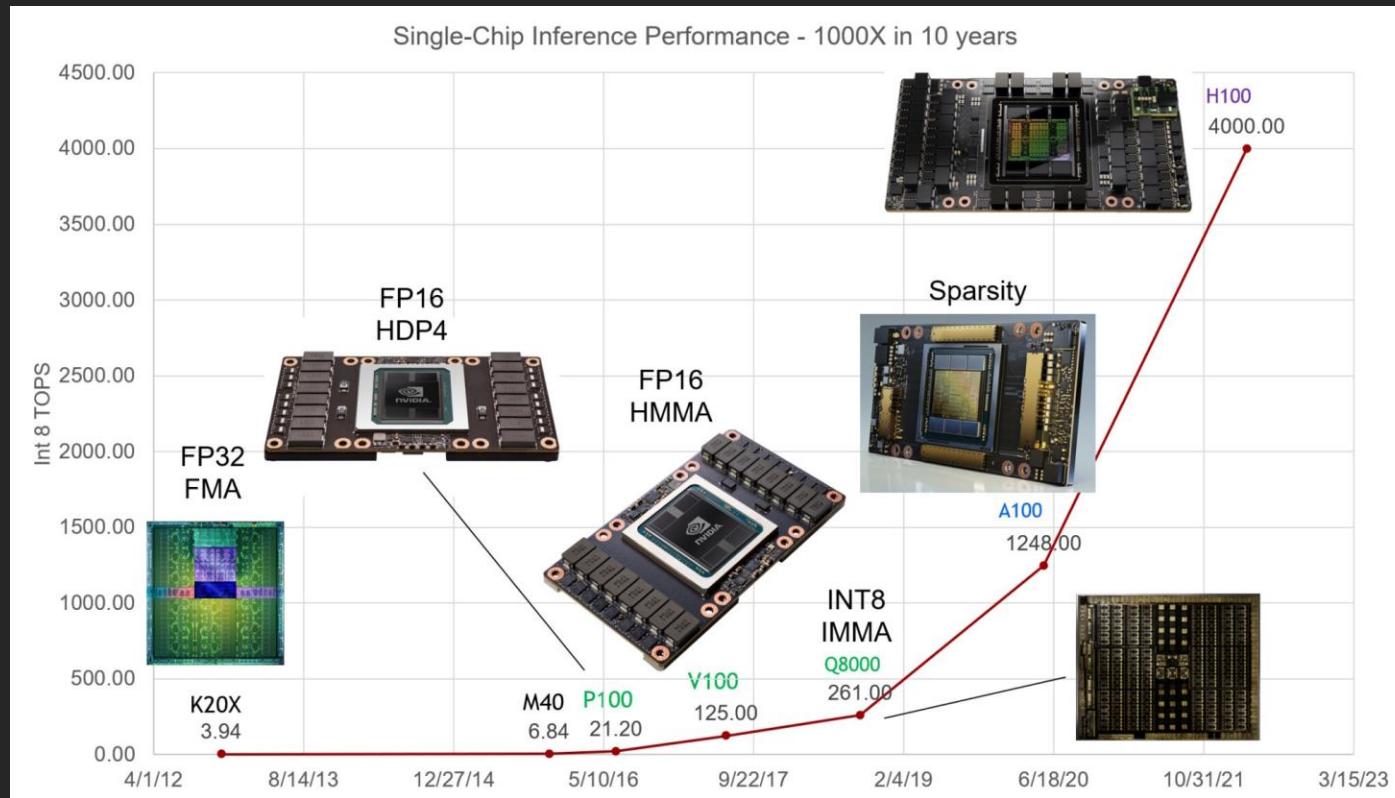
GPU - What is a GPU?

- A Graphics Processing Unit (GPU)
- Specialized processor originally designed to accelerate graphics rendering
- Unlike CPUs (Central Processing Units), which handle general-purpose tasks sequentially
- GPUs are built to process many operations in parallel, making them ideal for certain compute-intensive tasks.

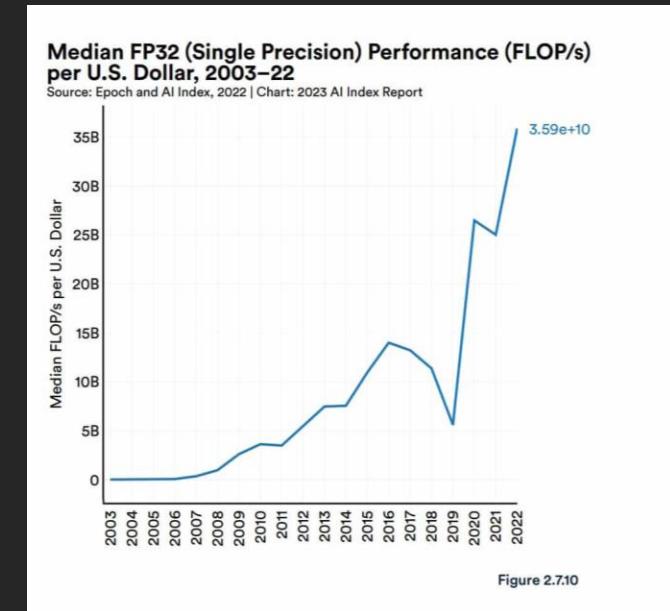
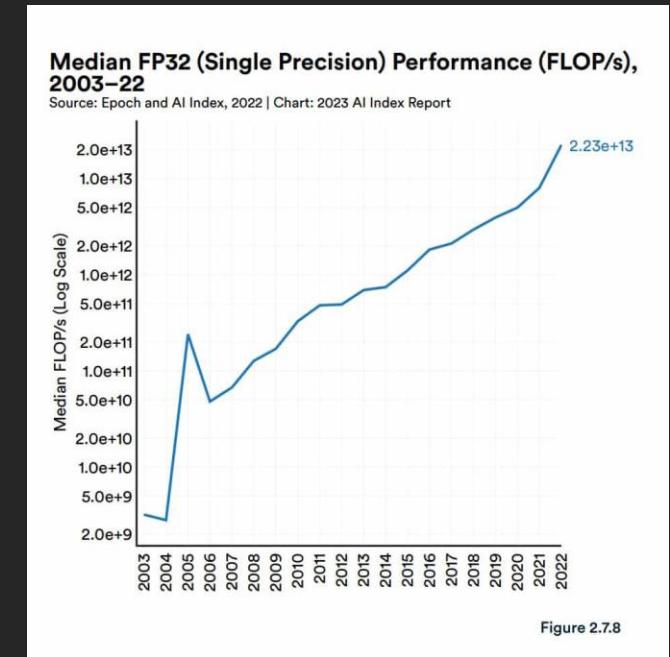
GPU

Models Grow, Systems Expand

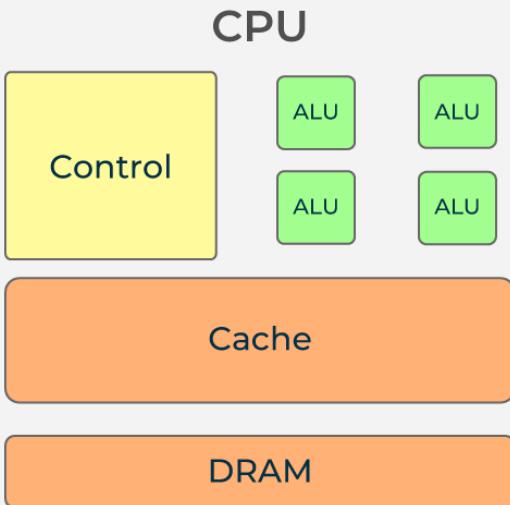
The complexity of AI models is expanding a whopping 10x a year.



https://blogs.nvidia.com/blog/why-gpus-are-great-for-ai/?utm_source=chatgpt.com



GPU vs CPU



- **CPU (Central Processing Unit)**

- Executes general-purpose tasks (arithmetic, logic).
- Handles complex instructions and task switching.
- CPU has a lower core count (4-8 cores) than others, focusing on delivering higher performance for tasks that rely on a single core.

GPU



- **GPU (Graphics Processing Unit)**

- Used for parallel processing tasks (graphics, AI).
- Has many cores (100s or 1000s of cores) for handling multiple operations simultaneously.
- GPUs use a specialized type of DRAM known as VRAM (Video RAM), specifically designed to handle graphical data and textures.

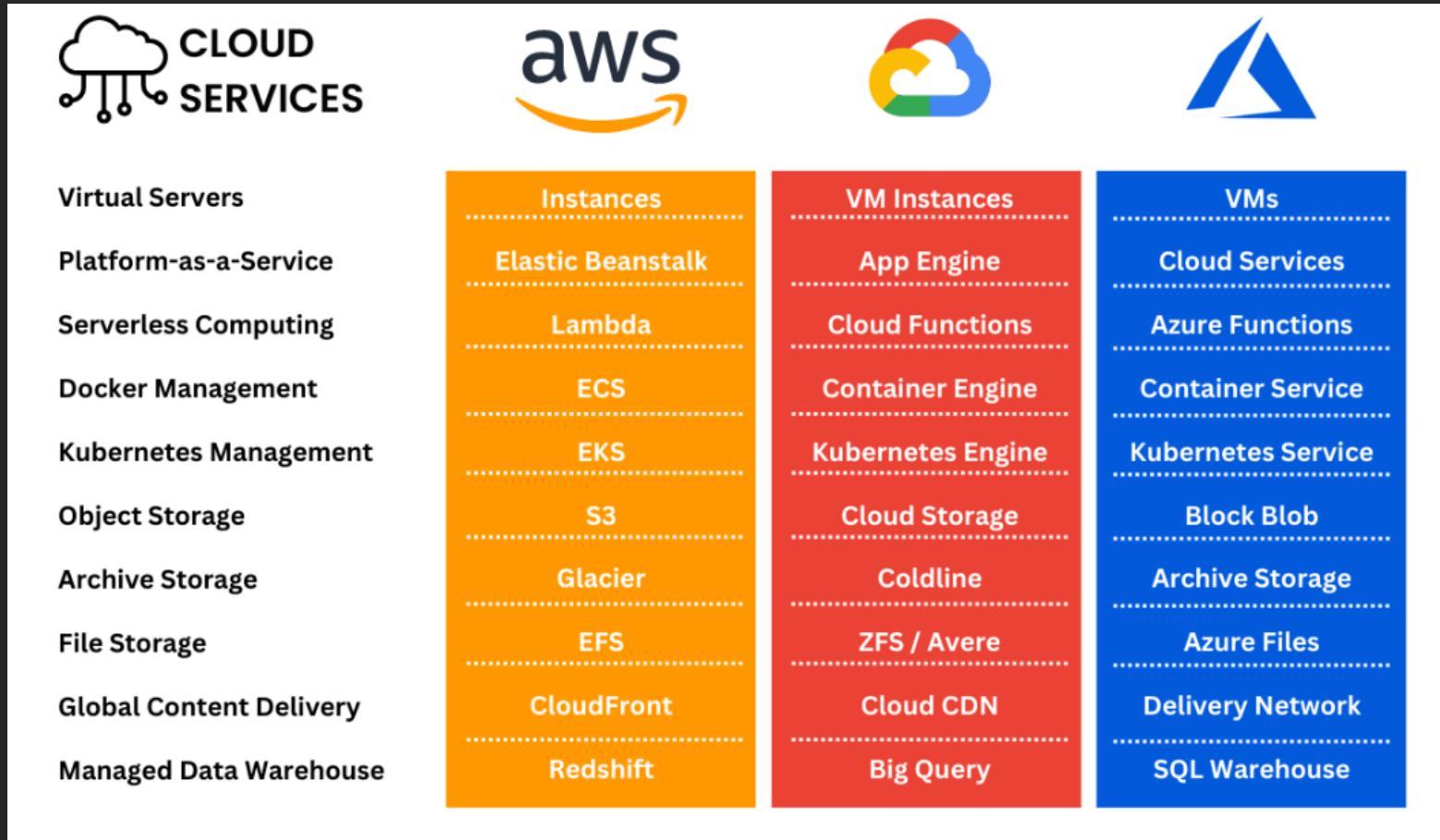
GPU vs CPU

CPU	GPU
Used for General-purpose computation.	Used for Specialized computation for graphics and parallel tasks.
Handles single-threaded, complex tasks.	Handles highly parallel tasks (e.g., graphics rendering).
Optimized for sequential processing.	Optimized for parallel processing.
Smaller cache memory (L1, L2, L3).	Larger memory (VRAM) optimized for high-speed data transfer.
More energy-efficient for general tasks.	Consumes more power due to parallel processing needs.
CPU emphasis on low <u>latency</u> .	While GPU emphasis on high <u>throughput</u> .
Runs operating system, applications, and tasks.	Handles graphics rendering, AI, machine learning.
Generally less expensive.	More expensive due to specialized hardware.

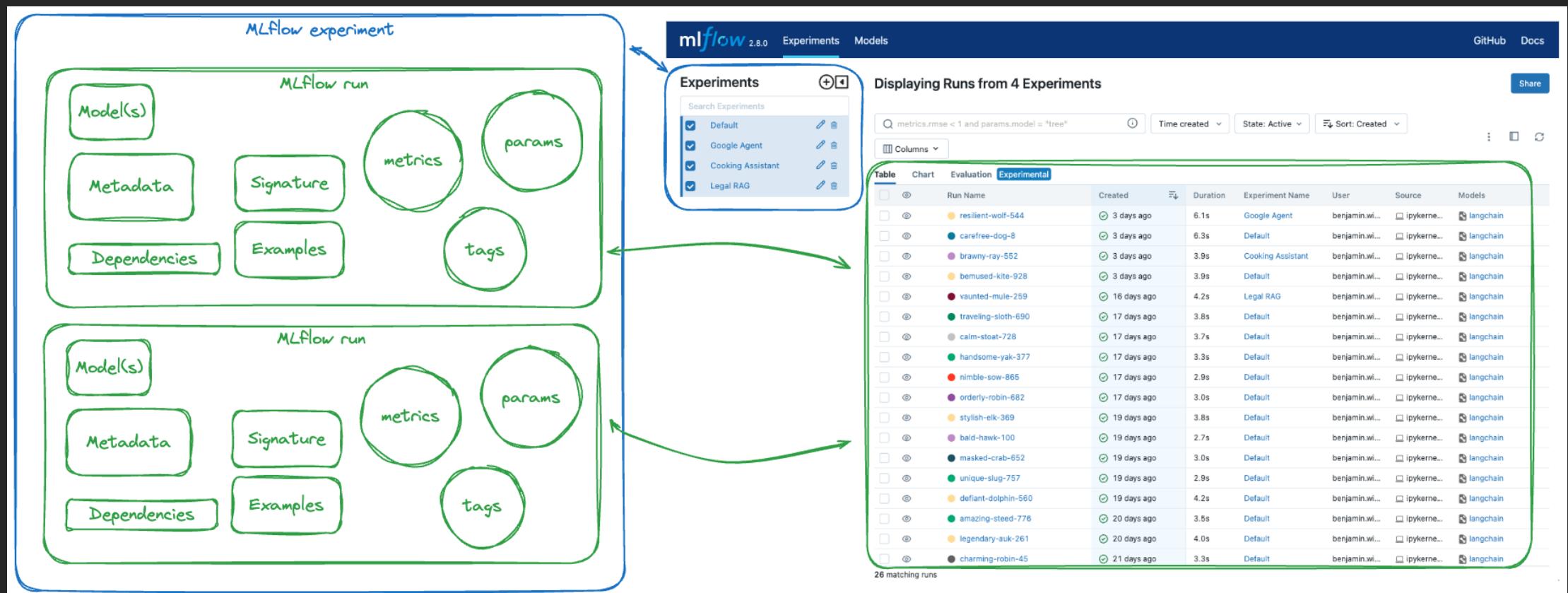
Cloud

- Massive Compute Power
- Scalability and Elasticity: Cloud infrastructure can elastically scale resources up or down based on AI workload demands.
- Cost Efficiency: Cloud's pay-as-you-go model reduces capital expenses and makes advanced AI capabilities accessible to both large enterprises and smaller organizations. (???)
- Collaboration and Accessibility
- Automation and Integration
- Hybrid and Multi-Cloud Flexibility

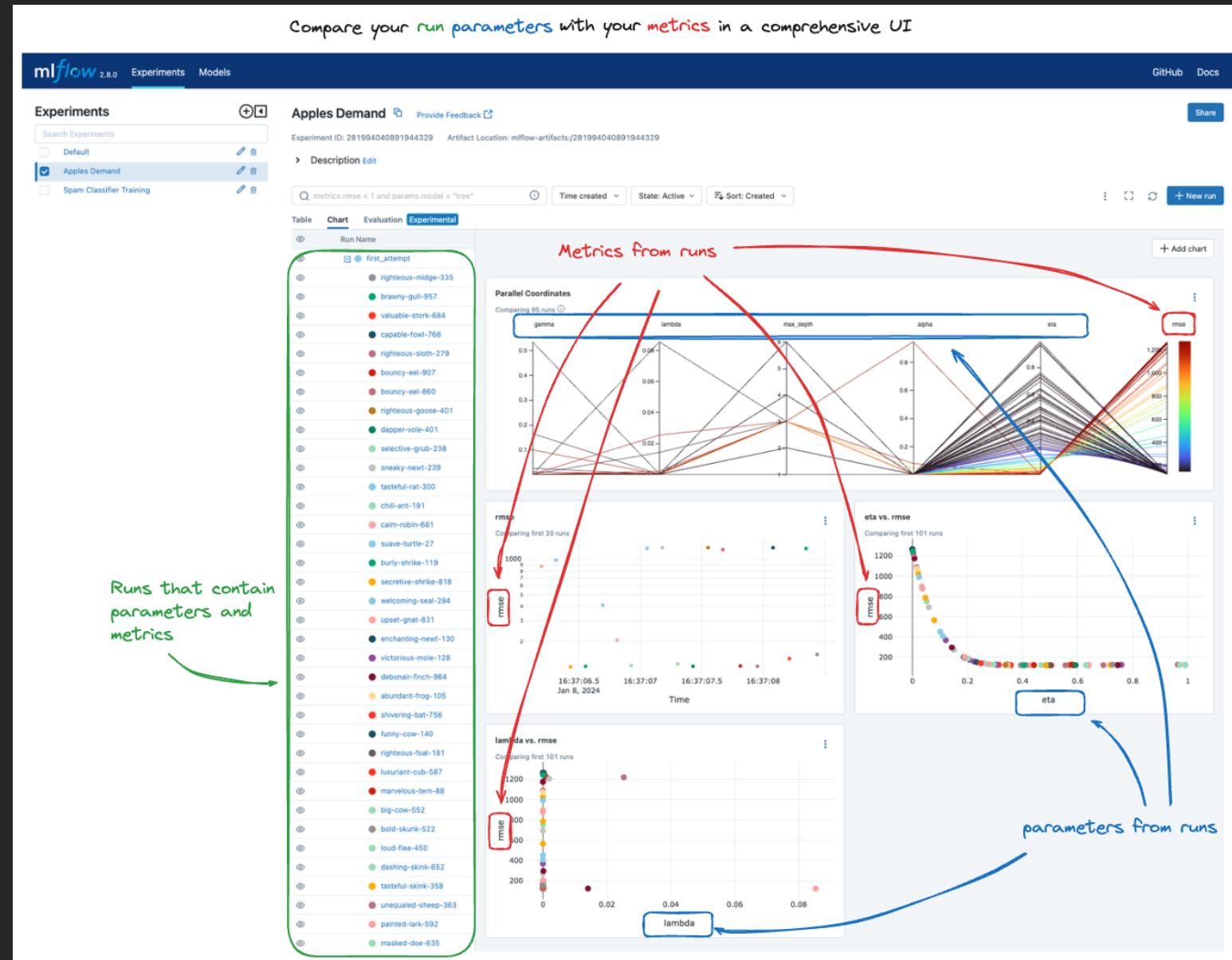
Cloud



MLflow

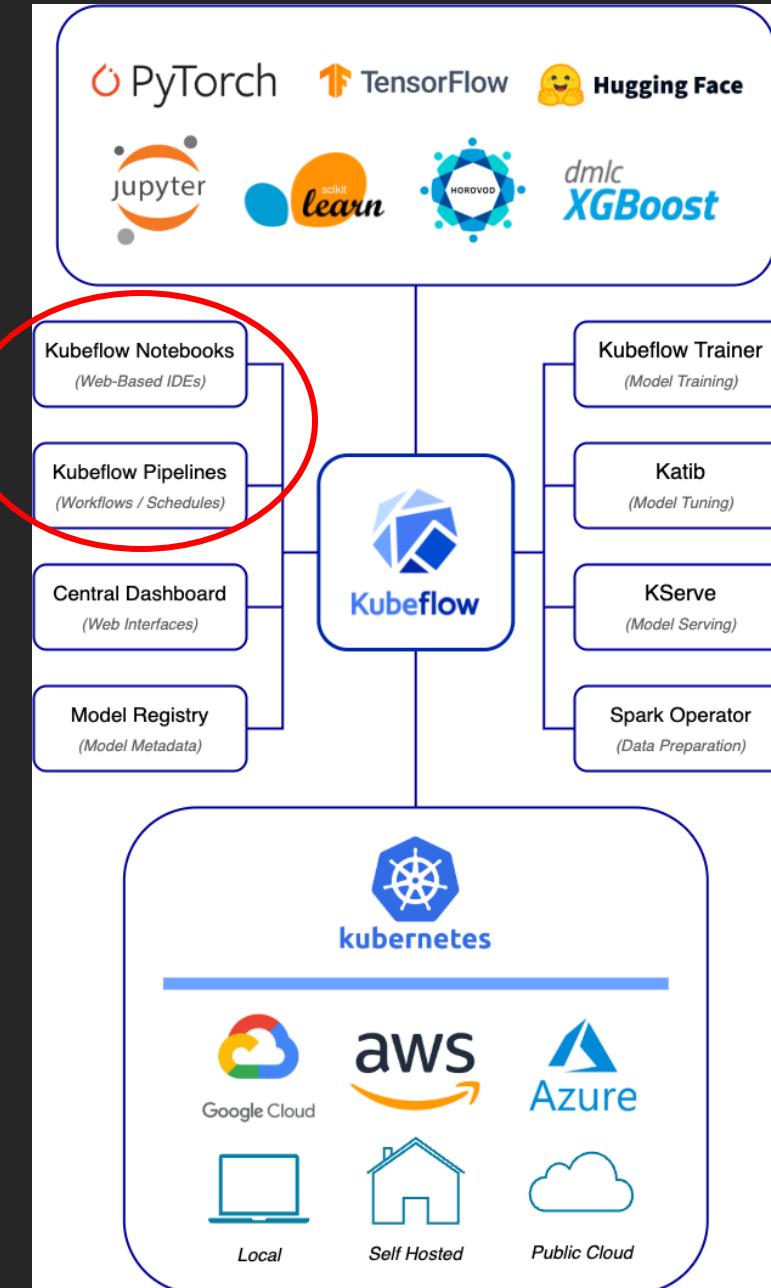


MLflow

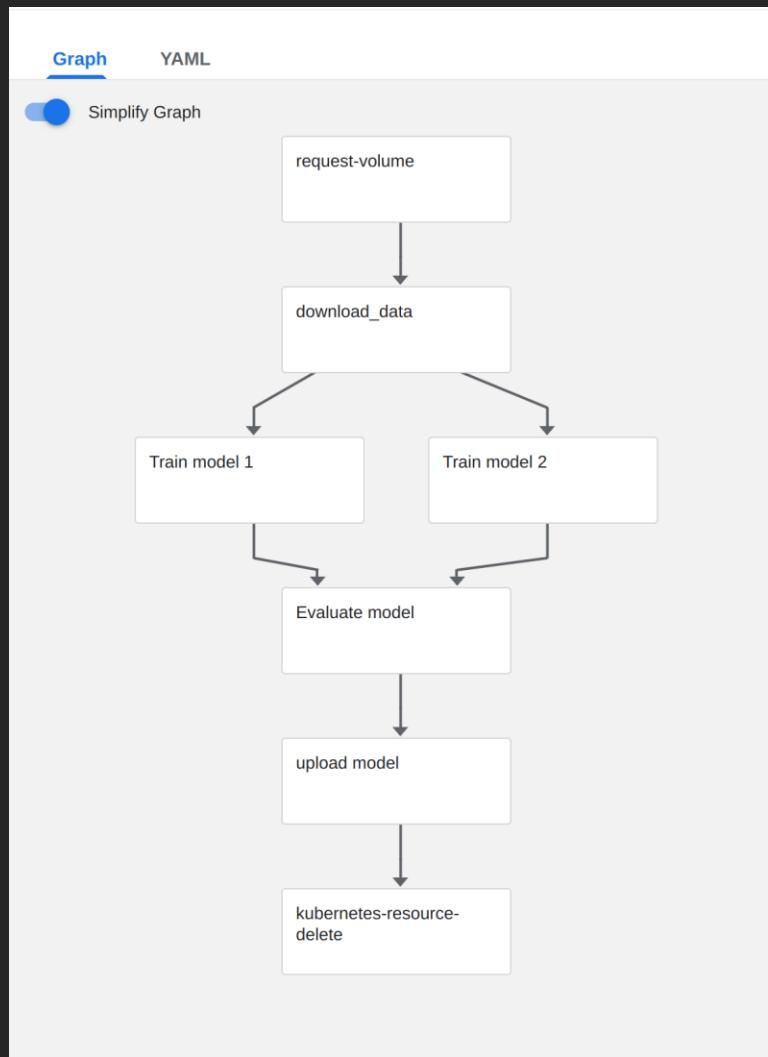


Kubeflow

- Kubeflow is a community and ecosystem of open-source projects to address each stage in the machine learning (ML) lifecycle with support for best-in-class open source tools and frameworks. Kubeflow makes AI/ML on Kubernetes simple, portable, and scalable.

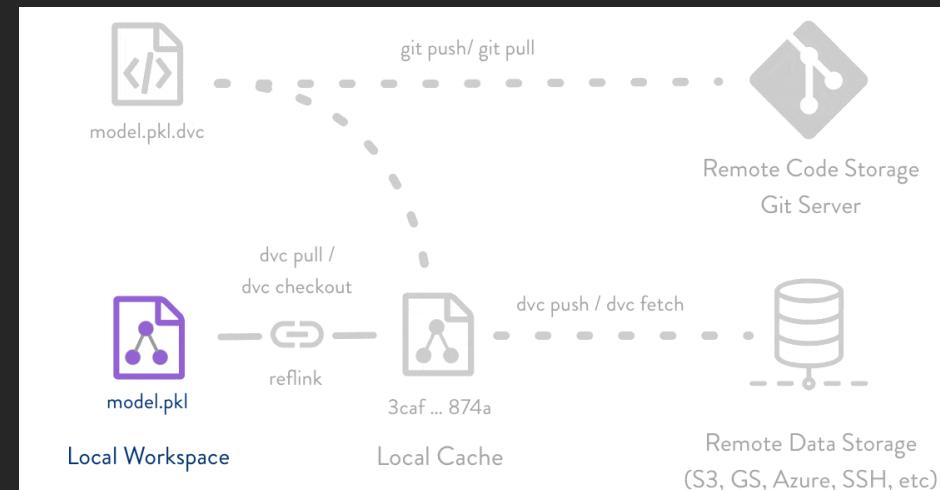


Kubeflow - Pipeline



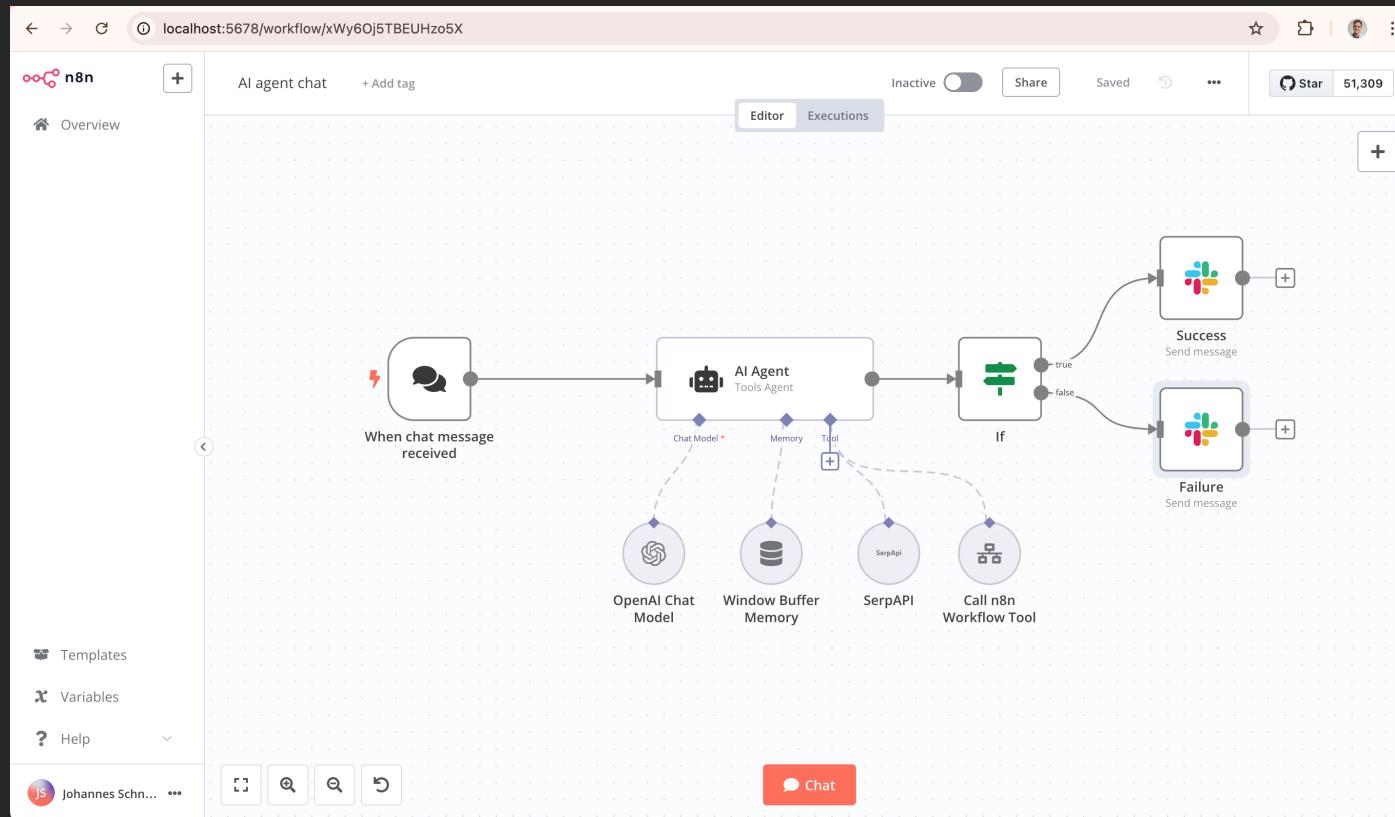
DVC (Data Version Control)

- Version your data and models. Store them in your cloud storage but keep their version info in your Git repo.
- Iterate fast with lightweight pipelines. When you make changes, only run the steps impacted by those changes.
- Track experiments in your local Git repo (no servers needed).
- Compare any data, code, parameters, model, or performance plots.
- Share experiments and automatically reproduce anyone's experiment.

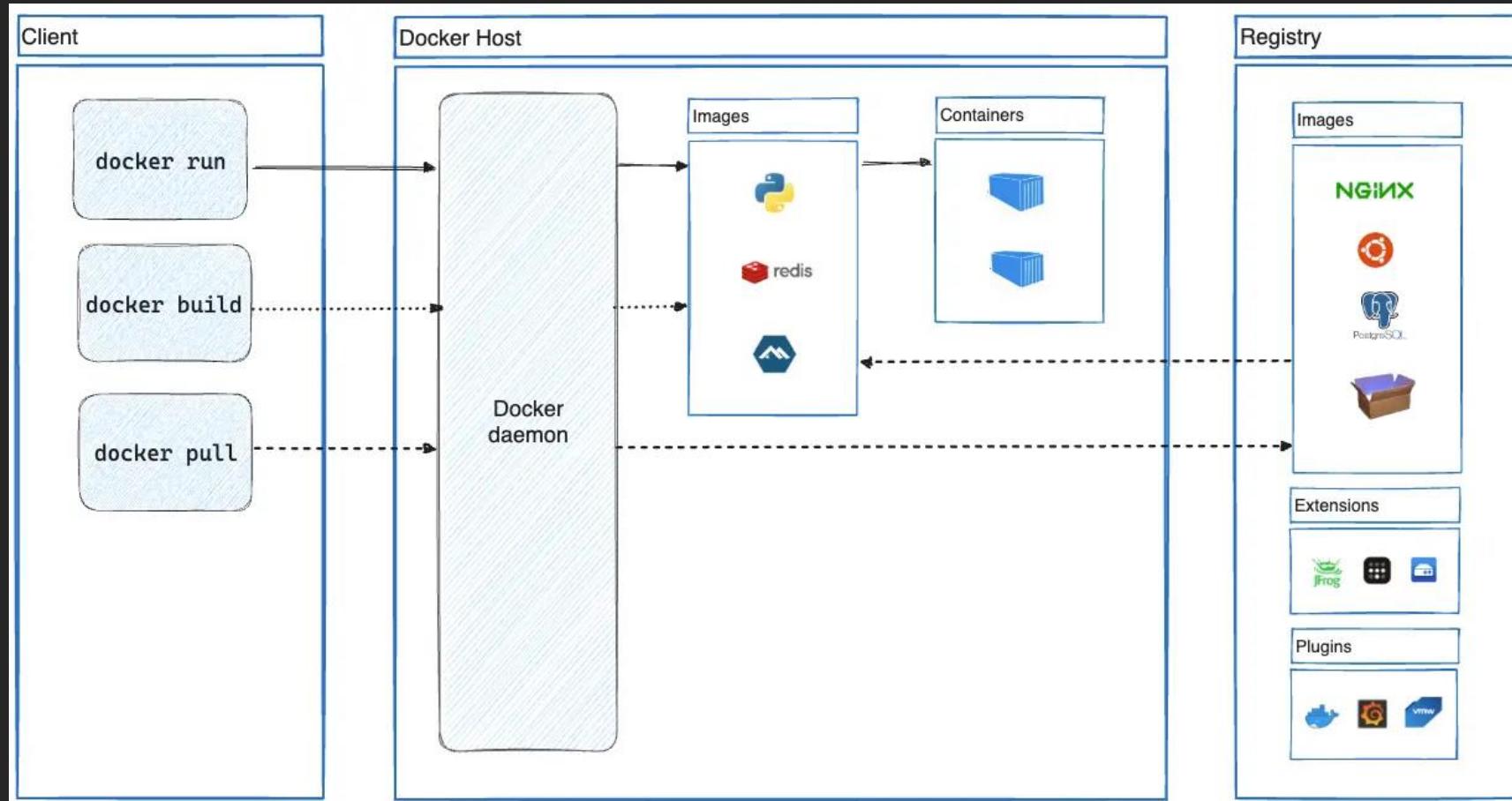


n8n

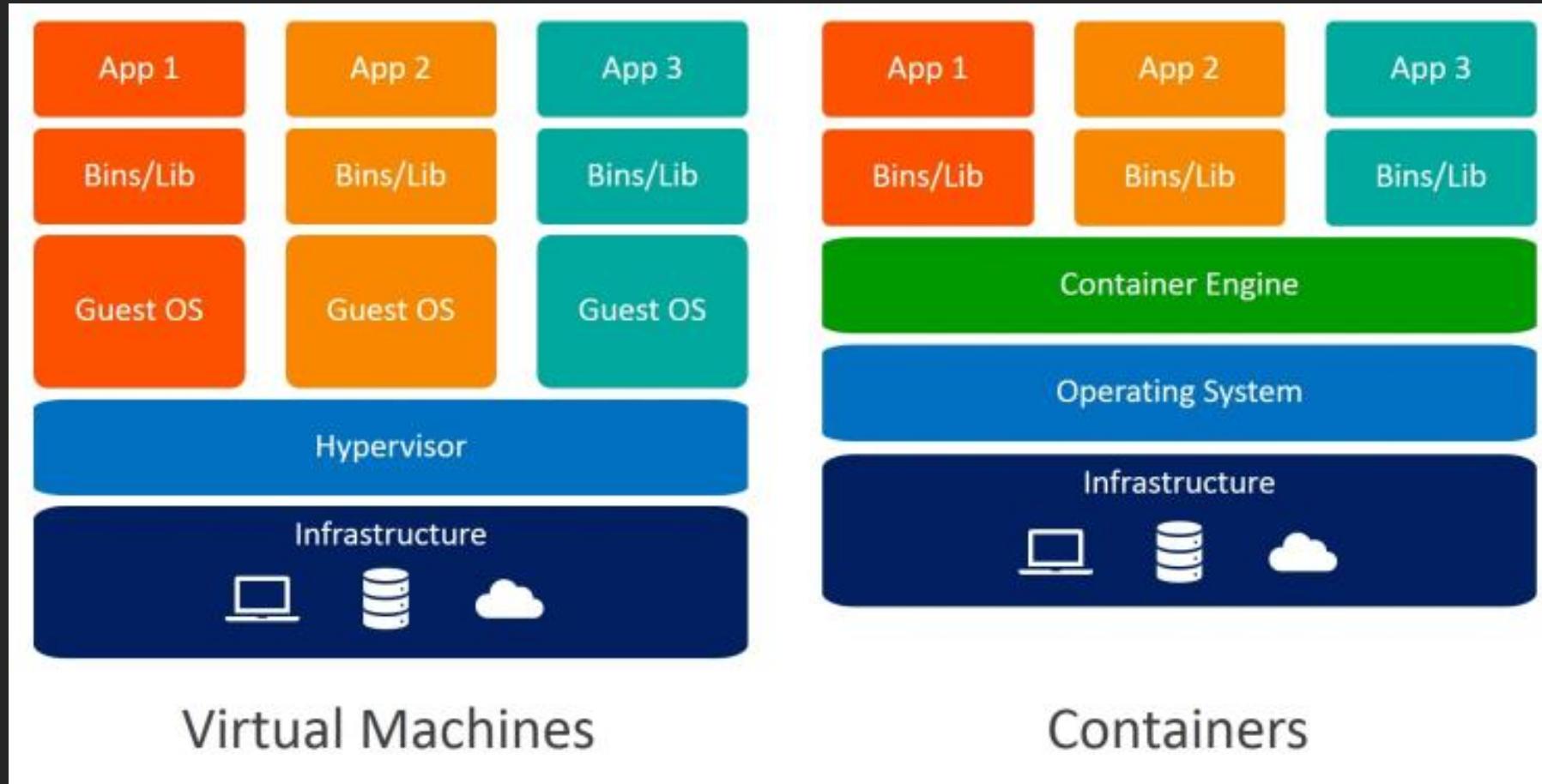
- AI agent workflow automation



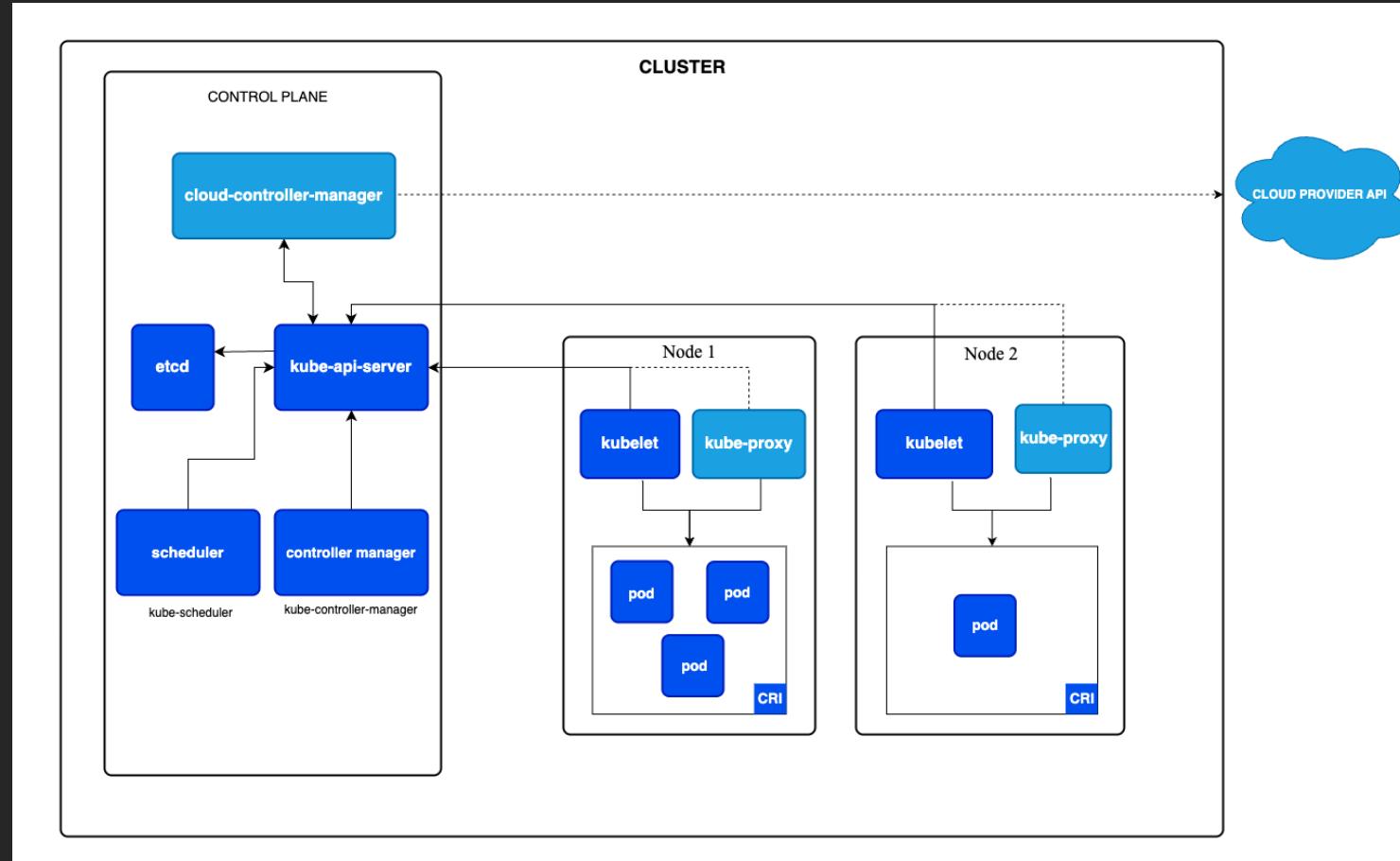
Docker (Container)



Docker (Container)



Kubernetes



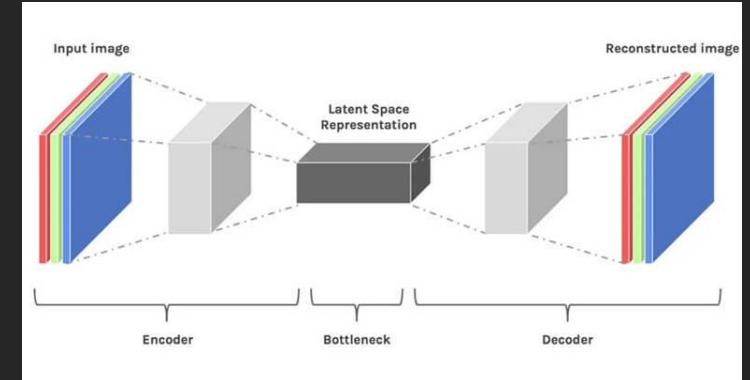
Current Trends in AI/ML

Generative AI

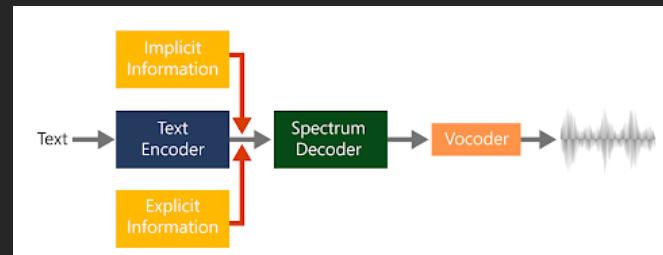
Text: LLM, ChatGPT



Image: Stable diffusion, DALL-E, Deepfake

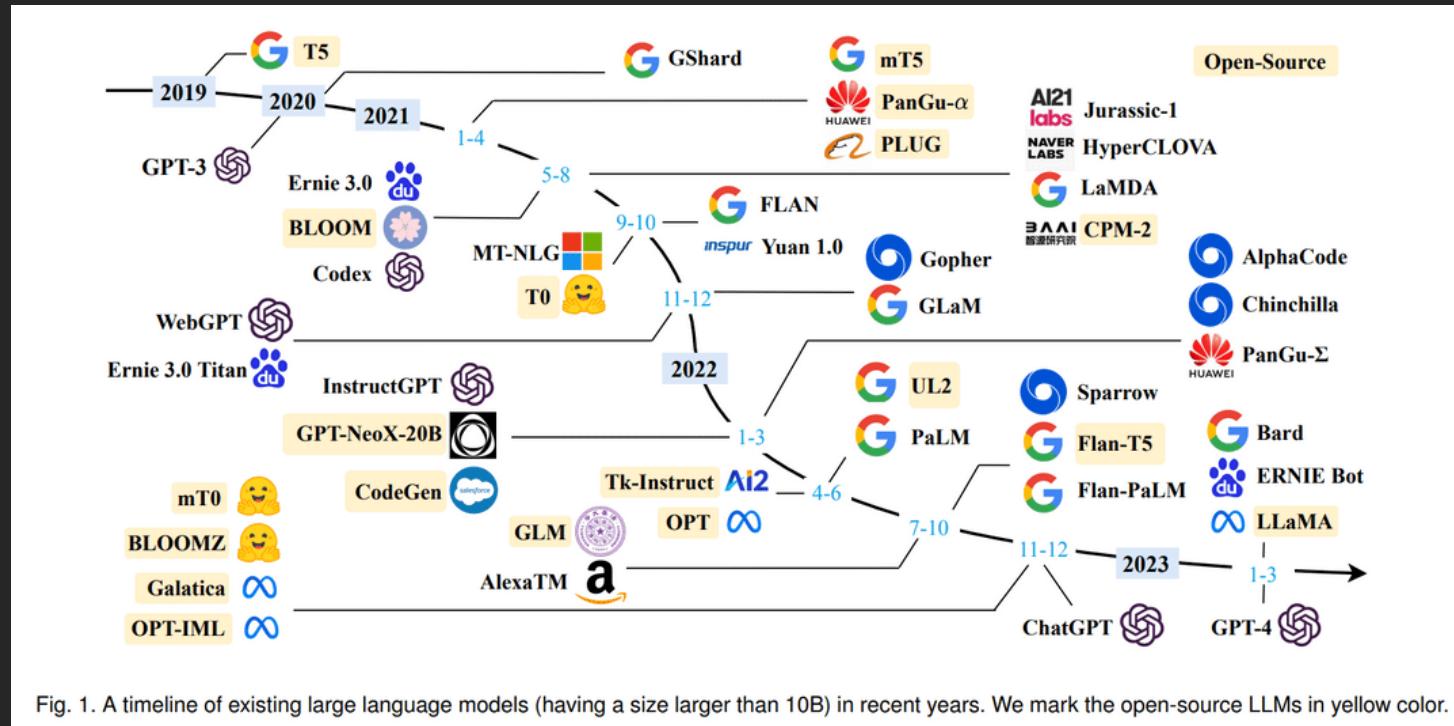


Voice: Text-to-Speech, Retrieval-based Voice Conversion



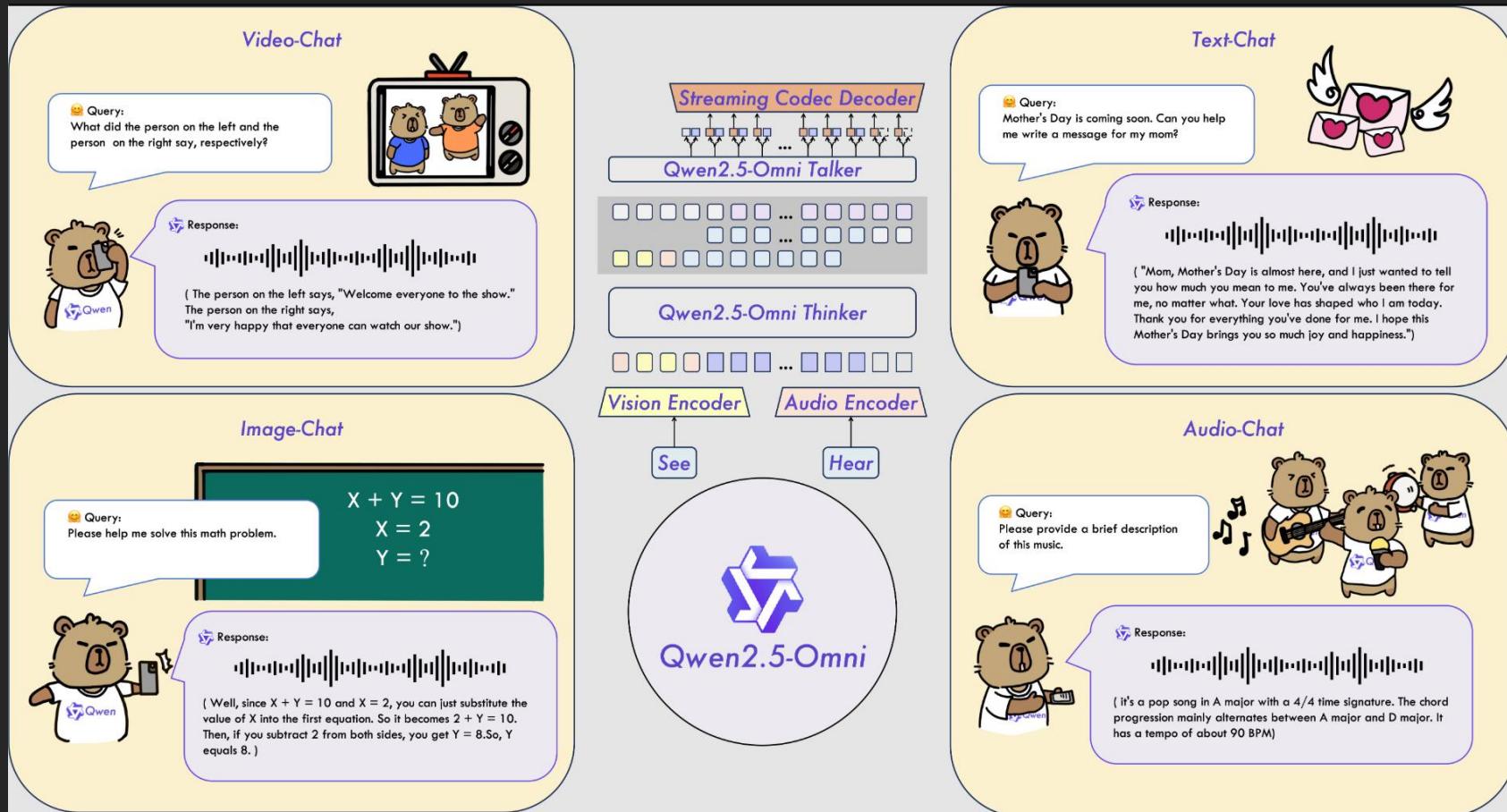
LLM (Large Language Model)

- A large language model is an AI system trained on vast amounts of text data to understand and generate human-like text for various language tasks, though it can sometimes produce inaccurate or biased outputs.

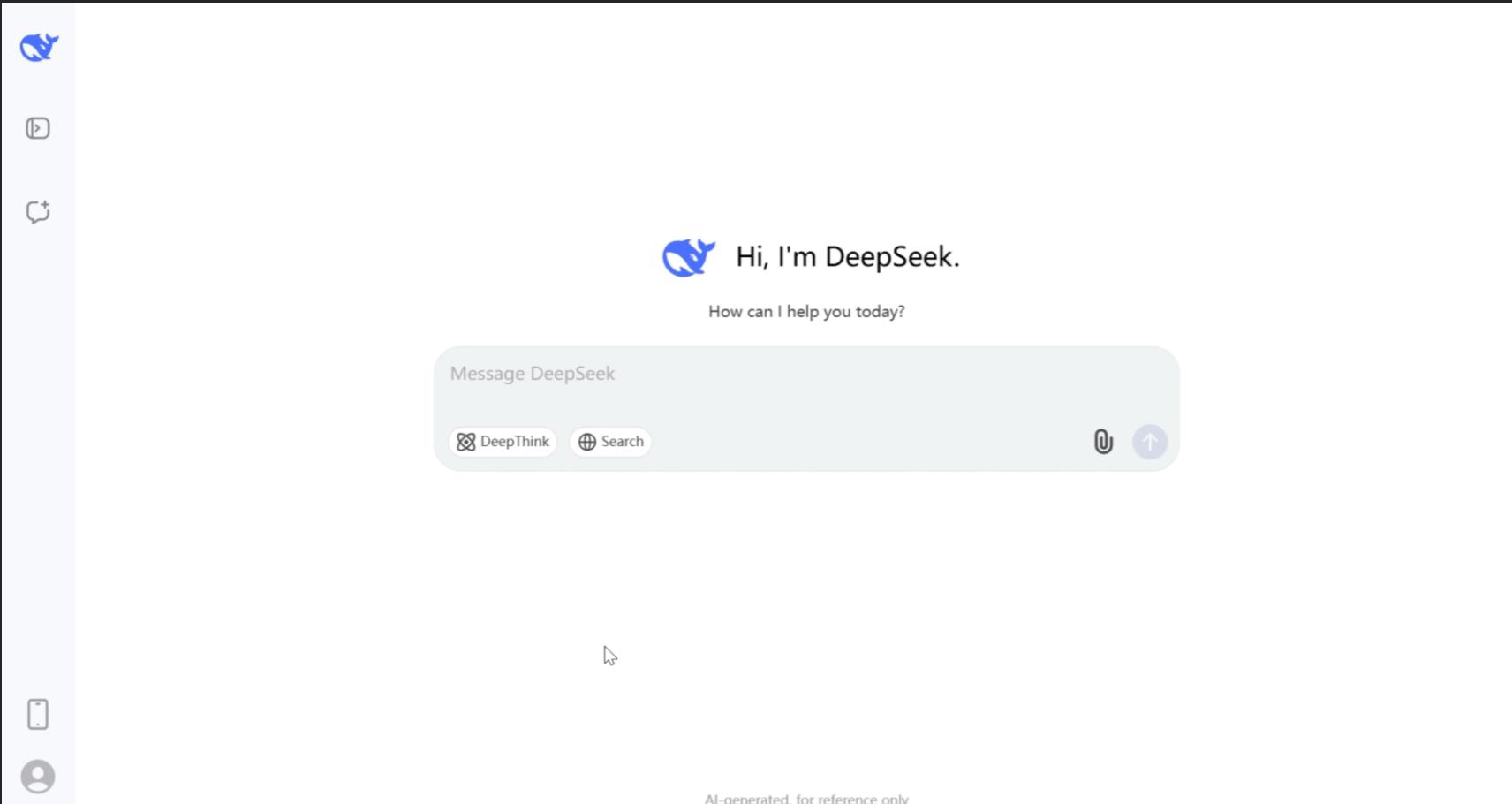


https://www.reddit.com/r/dataisugly/comments/12a8heb/a_timeline_of_large_language_models/?rdt=40625

LLM - OMNI (Multimodal)



LLM - Reasoning

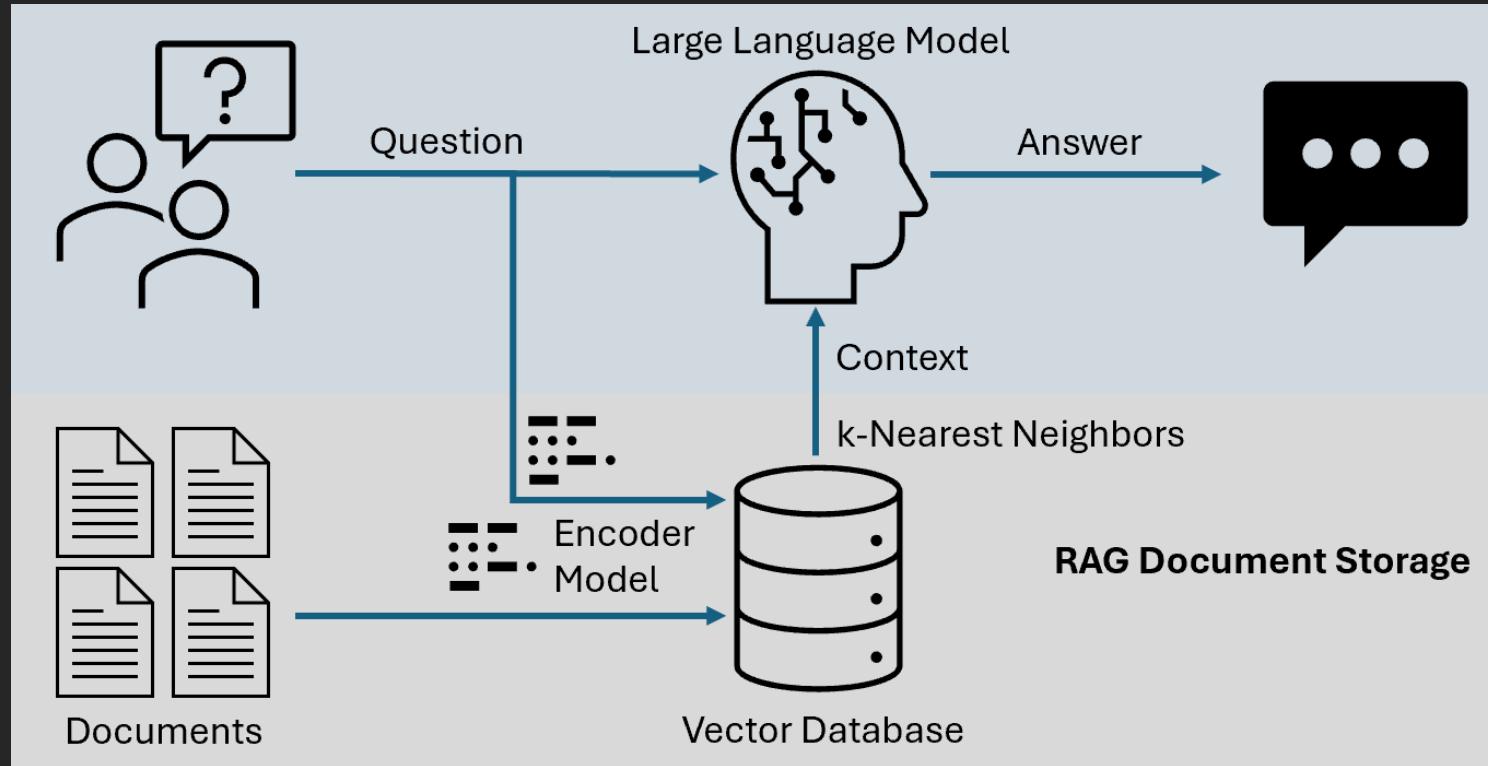


Prompt Engineering

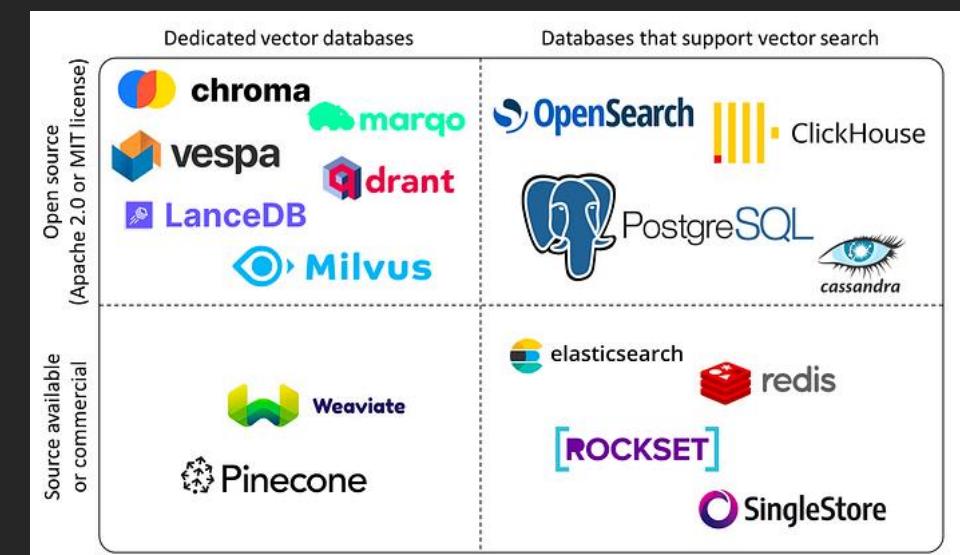
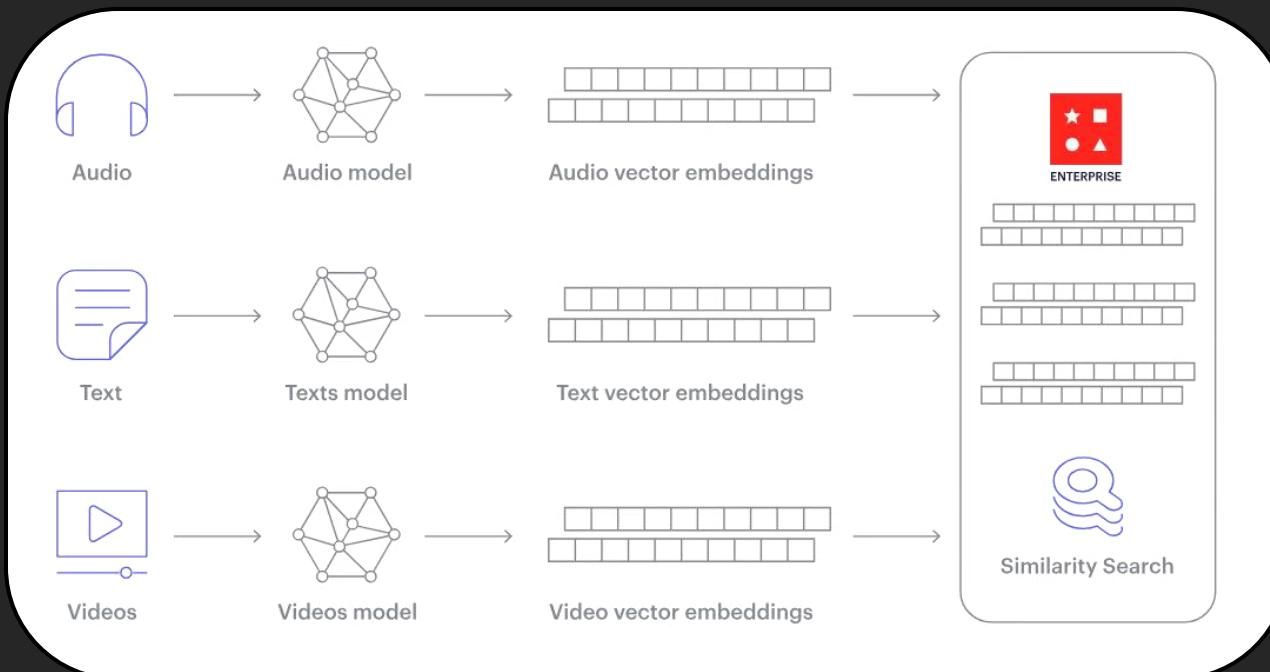
Crafting precise and effective prompts to communicate with AI models

1. **Clarity and Specificity:** Ensuring prompts are clear and specific to reduce ambiguity.
2. **Contextual Information:** Providing background context to guide accurate responses.
3. **Iterative Refinement:** Testing and refining prompts to improve results.
4. **Model Understanding:** Knowing how AI models process text to inform better prompts.
5. **Ethical Considerations:** Being aware of biases and ethical implications in prompt design.
6. **Applications:** Used in areas like customer support, content creation, and coding.

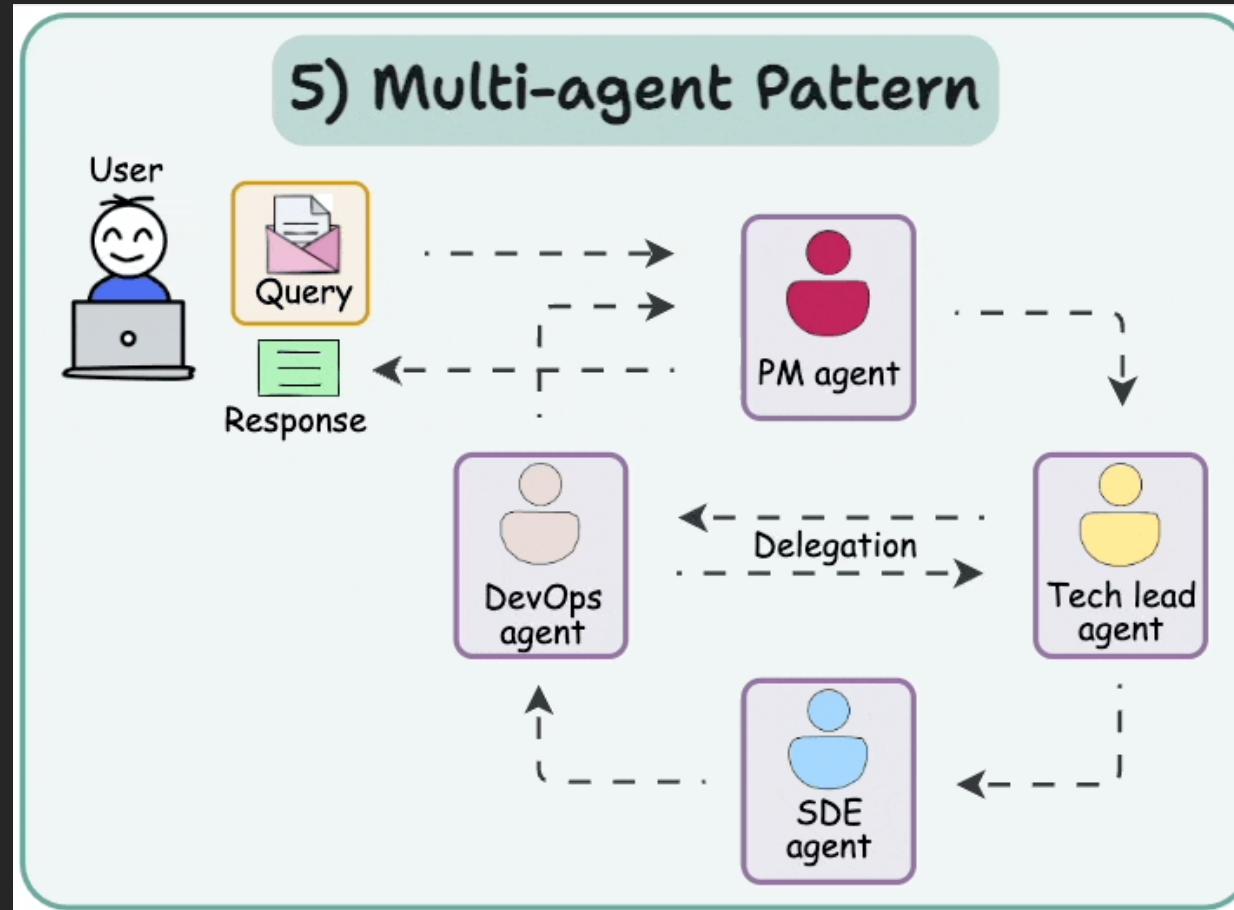
RAG - Retrieval-Augmented Generation



Vector Database

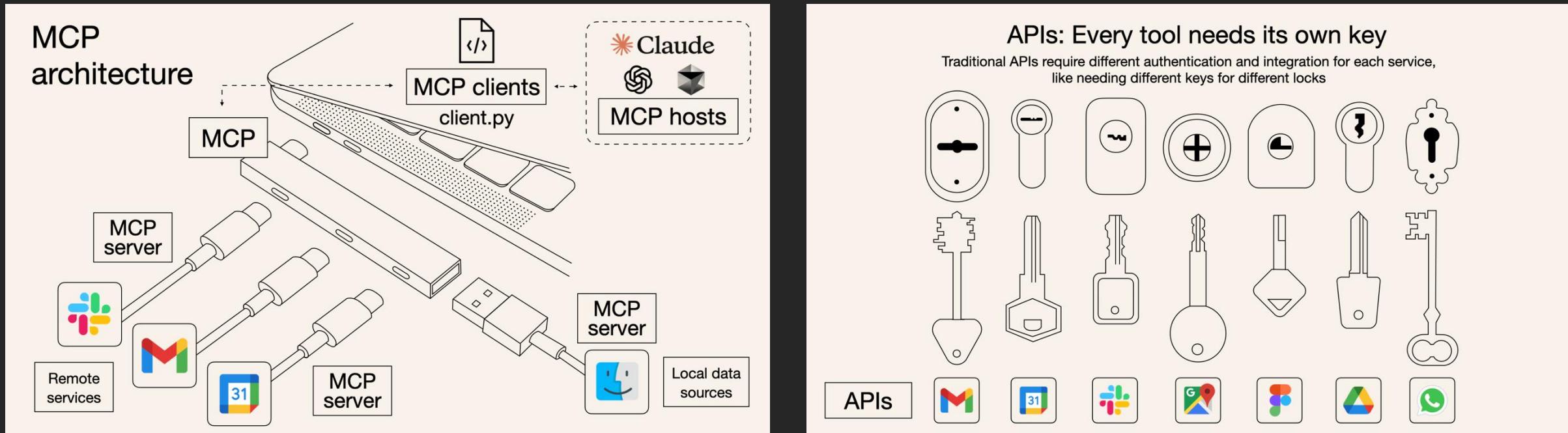


Agentic AI (Multi Agents)



Model Context Protocol (MCP)

- MCP is an open protocol that standardizes how applications provide context to LLMs. Think of MCP like a USB-C port for AI applications.



Wrap-up

Generative AI in a Nutshell

Henrik Kniberg
Jan 2024

Prompt Engineering/Design

- Bad prompt**: You give me an agenda for a workshop.
- Good prompt**: You give me an agenda for a workshop with context: I'm inviting a leadership team at an aerospace consulting firm. The goal of the workshop is figure out how they can use AI. They are new to it. We have 8 people for 4 hours.
- Good prompt**: You give me an agenda for a workshop with questions: Feel free to ask me any clarifying questions first!

Iterate!

Autonomous Agents with tools

Side effect: Better communication skills overall

Computers have gotten smarter

Execute Instructions

Learn + Think

Einstein in your basement

Biggest limitation is You

How it works

Neural Network

Dogs are animals

Dogs are animals → that

Dogs are animals that are known for... (walk...)

Models, models everywhere

Speed

Capability

Cost

Downloadable

Open Source

Easy to use

Specialized

Integrated

Training

- ① Unsupervised, Generative Pretraining (lots + lots of text...)
- ② Reinforcement Learning with Human Feedback (R+HF)

Model Types

Multimodal Models

Text → Text

Text → Image

Image → Image

Image → Text

Text → Audio

Text → Video

Taking AI for a walk

Emergent Capabilities

Small model: "John fell" → "down"

Large model: etc. teach, coach, legal/medical advice

AI

Artificial Intelligence

Generative AI

Machine Learning

Computer vision

Product: ChatGPT, Llama, etc

Model: GPT, Llama, etc

API

Product: ChatGPT, Llama, etc

Model: GPT, Llama, etc

API

you

Your Product

your users

your data

The role of Humans

Is human role X still needed?

doctor, developer, lawyer, ceo, teacher, etc...

Evaluate results: beware of hallucinations!

Decide what to ask and how

Provide context

Legal compliance

Data security

etc. etc...

Human + **AI** = **Heart**

Mindset

Positive: I will be insanely productive!

Panic: AI is going to steal my job!

Denial: Non AI can't do my job!

AI might not take your job, but people/companies using AI will

The Age of AI

Fast Revolution

Slow Revolutions

Time: ~80 years

Human Intelligence → **Artificial Intelligence**

Roleplay, etc. teach, coach, legal/medical advice

Poetry

Code

Strategy

What will happen when I use the scissors?

ChatGPT: What's the best way to hide something valuable in this room? Here are a few suggestions:

- Behind the map: You could tape something flat, like documents or cash, behind the map.
- Inside the guitar: If the guitar has a case, you could hide something inside it. It's not a good idea to damage the guitar.
- Under the couch cushion: A common hiding spot, but effective for small and flat items.
- Behind the wood stove: There is any removable panel or space, although this is likely due to the heat when it's on.
- Under the lamp: If the bottom of the lamp opens or there is a space where the lampshade is attached.

 It would prevent small items.

Fire, **Agriculture**, **Printing press**, **Steam power**, **Telegraph**

Careers

PS. my personal opinion

Data/ML/AI Related Careers

Data Scientist

Data Engineer

Data Analyst

ML Engineer

AI Engineer

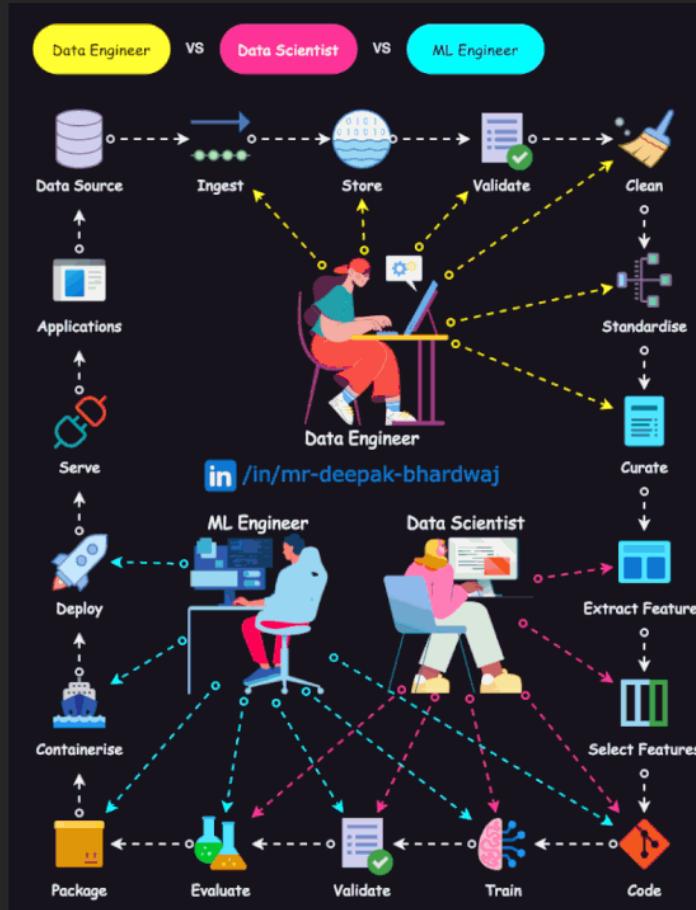
AI Researcher

Computer
Vision Engineer

NLP Engineer

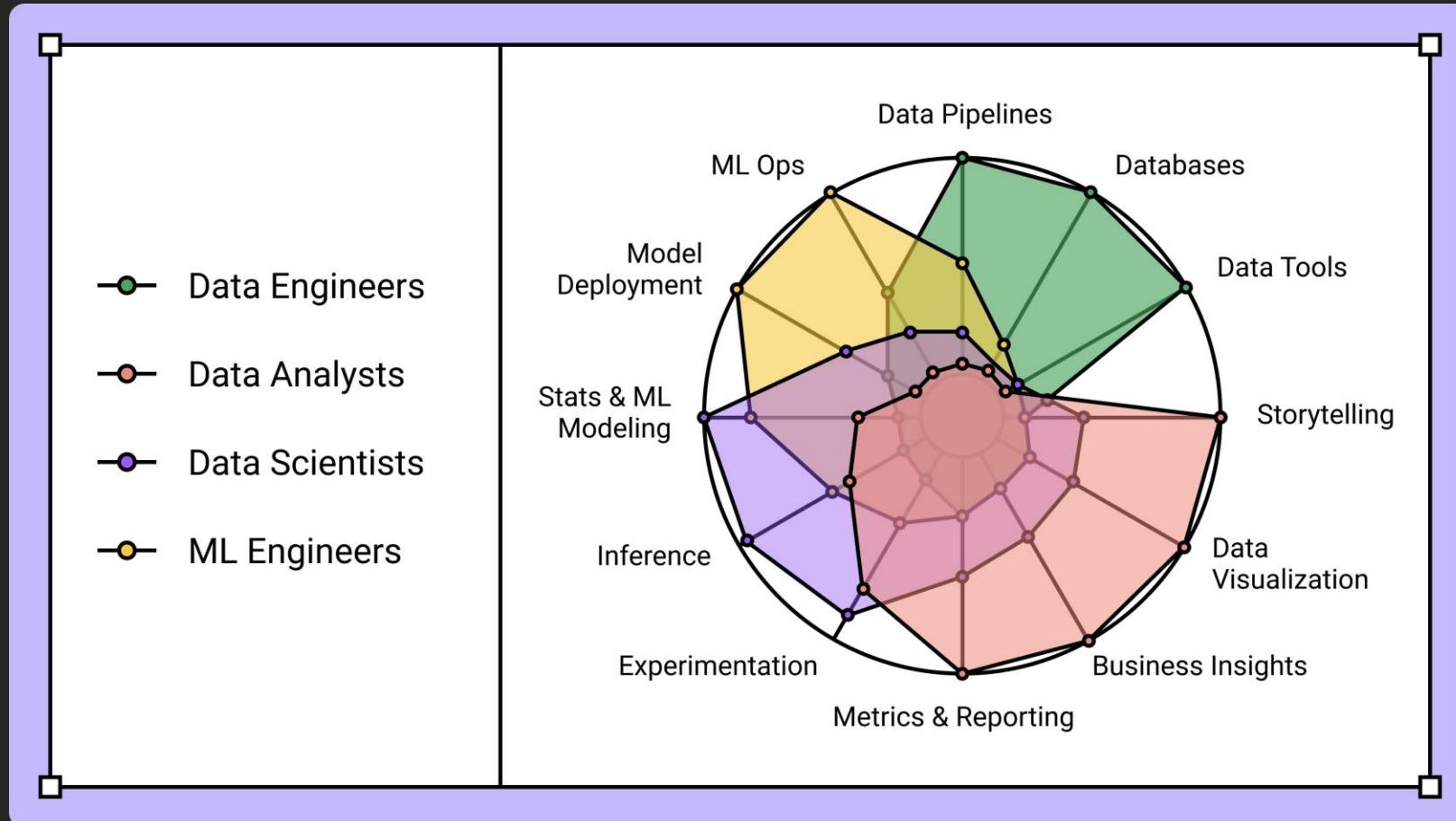
MLOps
Engineer

Data/ML/AI Related Careers



<https://www.linkedin.com/pulse/data-engineer-vs-scientist-ml-heba-al-haddad>

Data/ML/AI Related Careers



Ex. Data scientist

Job Description:

- Apply statistical and machine learning methods to large, complex data sets to draw insights and provide actionable recommendations.
- Solve complex problems on both technical and business sides using advanced analytical methods.
- Work with Engineering teams to implement end-to-end process from model development to testing, validation, and deployment
- Research and develop new quantitative models and frameworks to enhance the company's data science capability

Responsibilities:

- Design and implement predictive models, machine learning algorithms, and advanced statistical analyses to address business needs.
- Analyze large and complex datasets to extract insights and identify trends, patterns, and opportunities.
- Collaborate with cross-functional teams to define business challenges, collect requirements, and ensure alignment with data solutions.
- Build and deploy end-to-end machine learning solutions, from feature engineering to model deployment.
- Research and stay updated on the latest AI/ML trends and techniques to apply innovative solutions to business challenges.
- Support the implementation of data-driven strategies to improve operational efficiency and enhance customer experience.

DUTIES AND RESPONSIBILITIES:

- Identify valuable data sources and automate data collection for AI and chatbot development.
- Preprocess structured/unstructured data for training ML and Generative AI models.
- Analyze large datasets to uncover trends and enhance learning experiences.
- Build and fine-tune predictive models and chatbot assistants using Generative AI.
- Deploy AI models and chatbots to web-based production environments with scalability in mind.
- Visualize insights and model performance via clear, actionable dashboards.
- Collaborate with engineering and product teams to integrate AI into the platform.
- Implement AI assistants with focus on NLU, personalization, and ongoing learning.
Maintain and investigate ETL pipelines that feed into the data warehouse.
- Create proof-of-concept (POC) solutions and present ideas clearly to management and stakeholders.
- Stay up-to-date with Generative AI trends to improve interactive learning tools.

READ JD !!!

Q/A

AI Index Report 2025

ref: <https://hai.stanford.edu/ai-index/2025-ai-index-report>