

Humans are HOOKED

Machines are LEARNING

INTELIGENCIA ARTIFICIAL

# Generative AI

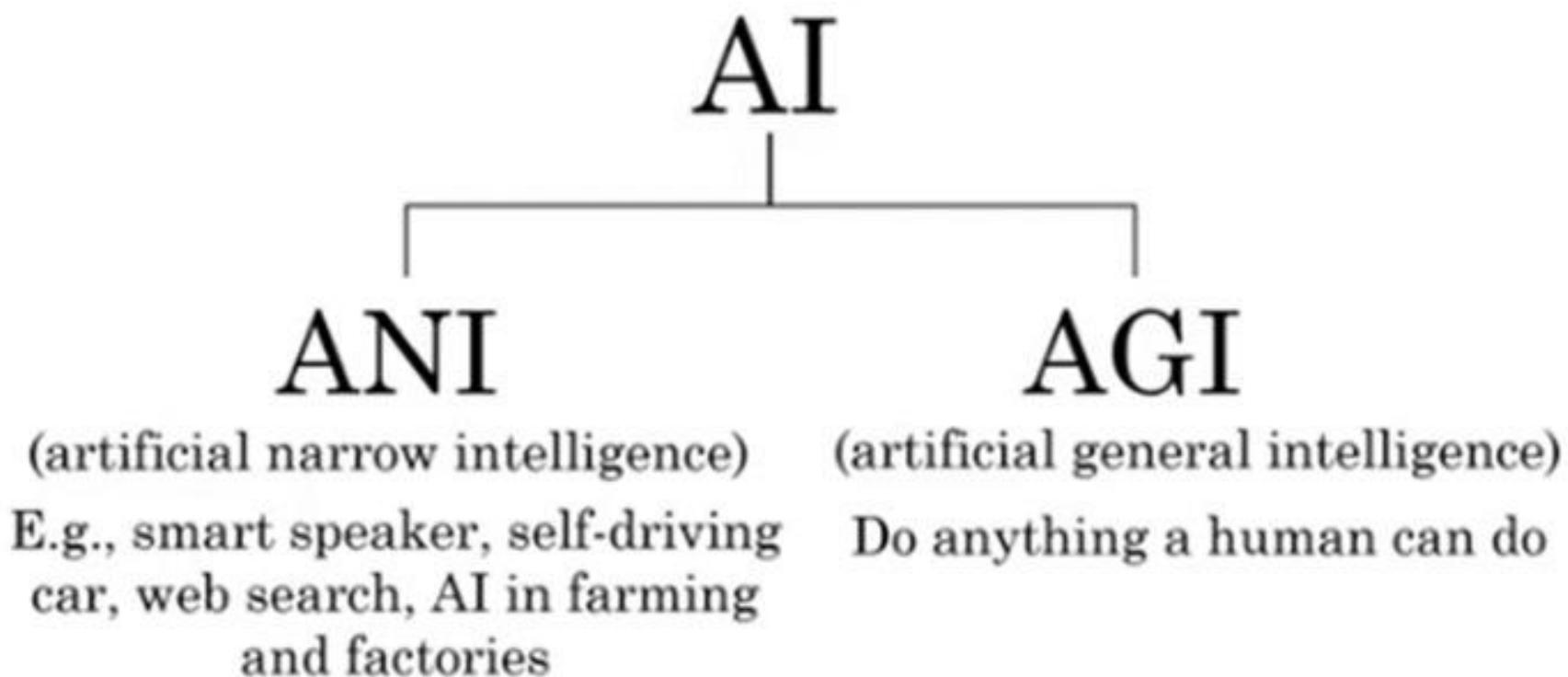


Naveen Kumar Bhansali

Co-Founder BlitzAI | Adjunct Faculty @  
IIM Bangalore



## AGI Vs. ANI



# AGI Levels

## AGI tier list

-- Google DeepMind

### Levels of AGI

Performance (rows) x Generality (columns)	Narrow <i>clearly scoped task or set of tasks</i>	General <i>wide range of non-physical tasks, including metacognitive abilities like learning new skills</i>
<b>Level 0: No AI</b>	Narrow Non-AI calculator software; compiler	General Non-AI human-in-the-loop computing, e.g., Amazon Mechanical Turk
<b>Level 1: Emerging</b> <i>equal to or somewhat better than an unskilled human</i>	Emerging Narrow AI GOFAI <sup>4</sup> ; simple rule-based systems, e.g., SHRDLU ( <a href="#">Winograd, 1971</a> )	Emerging AGI ChatGPT ( <a href="#">OpenAI, 2023</a> ), Bard ( <a href="#">Anil et al., 2023</a> ), Llama 2 ( <a href="#">Touvron et al., 2023</a> )
<b>Level 2: Competent</b> <i>at least 50th percentile of skilled adults</i>	Competent Narrow AI toxicity detectors such as Jigsaw ( <a href="#">Das et al., 2022</a> ); Smart Speakers such as Siri ( <a href="#">Apple</a> ), Alexa ( <a href="#">Amazon</a> ), or Google Assistant ( <a href="#">Google</a> ); VQA systems such as PaLI ( <a href="#">Chen et al., 2023</a> ); Watson ( <a href="#">IBM</a> ); SOTA LLMs for a subset of tasks (e.g., short essay writing, simple coding)	Competent AGI not yet achieved
<b>Level 3: Expert</b> <i>at least 90th percentile of skilled adults</i>	Expert Narrow AI spelling & grammar checkers such as Grammarly ( <a href="#">Grammarly, 2023</a> ); generative image models such as Imagen ( <a href="#">Shaharia et al., 2022</a> ) or Dall-E 2 ( <a href="#">Ramesh et al., 2022</a> )	Expert AGI not yet achieved
<b>Level 4: Virtuoso</b> <i>at least 99th percentile of skilled adults</i>	Virtuoso Narrow AI Deep Blue ( <a href="#">Campbell et al., 2002</a> ), AlphaGo ( <a href="#">Silver et al., 2016, 2017</a> )	Virtuoso AGI not yet achieved
<b>Level 5: Superhuman</b> <i>outperforms 100% of humans</i>	Superhuman Narrow AI AlphaFold ( <a href="#">Jumper et al., 2021; Varadi et al., 2021</a> ), AlphaZero ( <a href="#">Silver et al., 2018</a> ), StockFish ( <a href="#">Stockfish, 2023</a> )	Artificial Superintelligence (ASI) not yet achieved

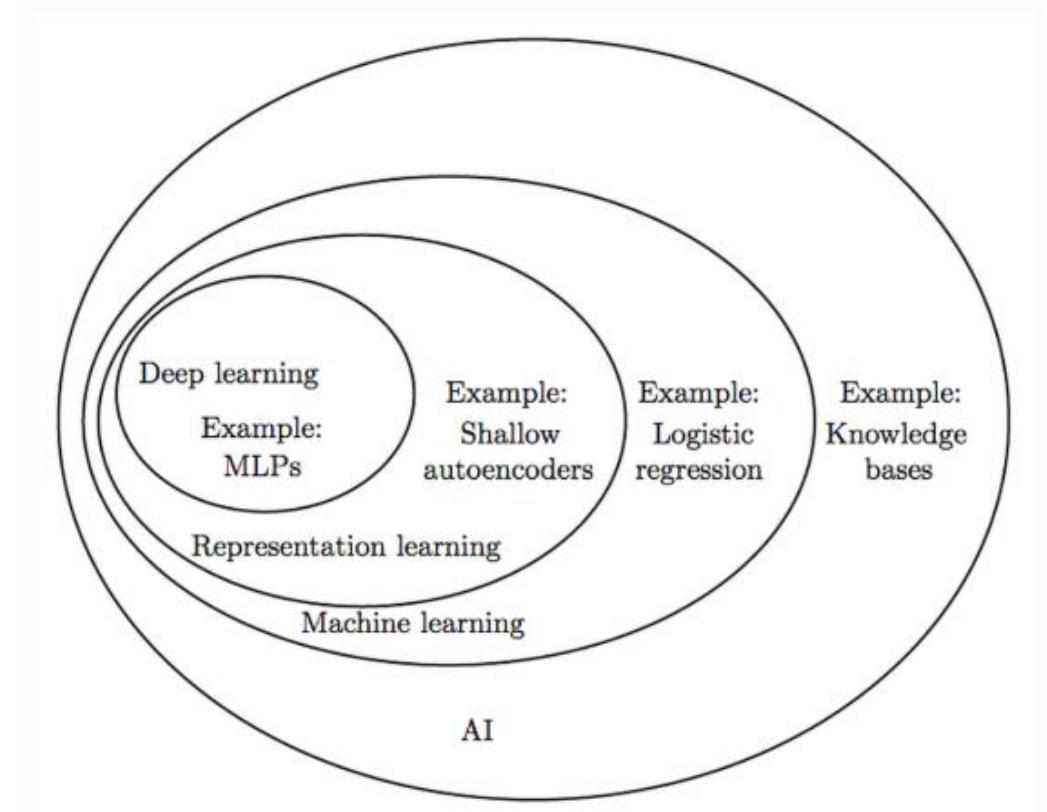


# Agenda

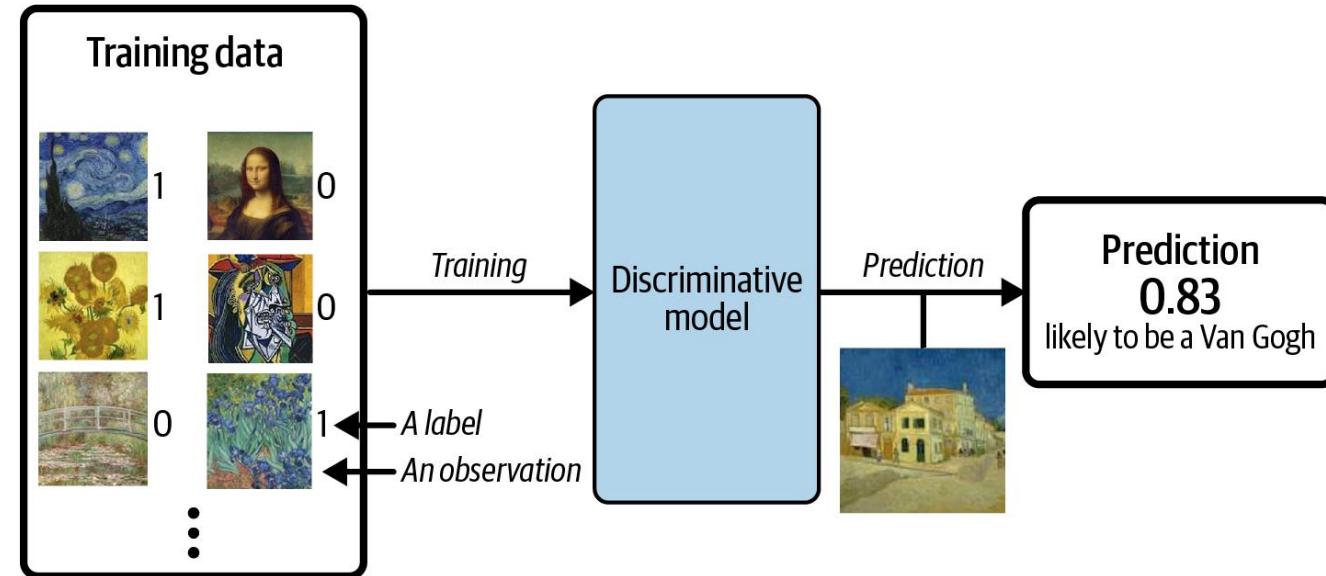
- Generative AI - Large Language Models / Large Multimodal Models
- Applications
- Case Studies
- Generative AI project lifecycle
  - Cheat Sheet
  - Prompt Engineering
  - LLM Powered Applications

# AI / ML / DL / GenAI

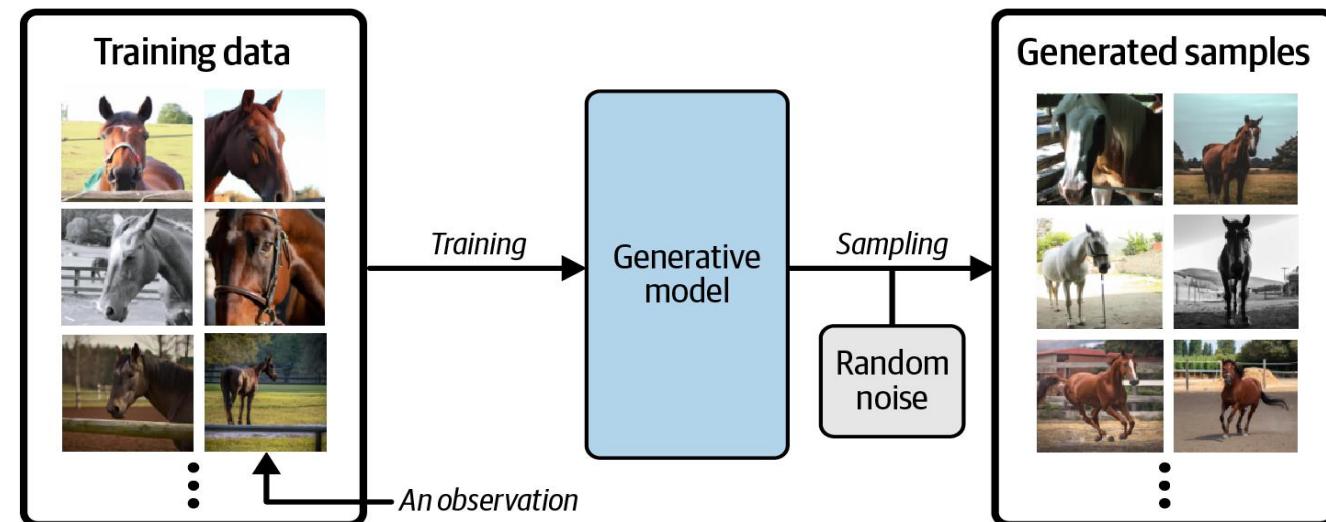
## – How are they related?

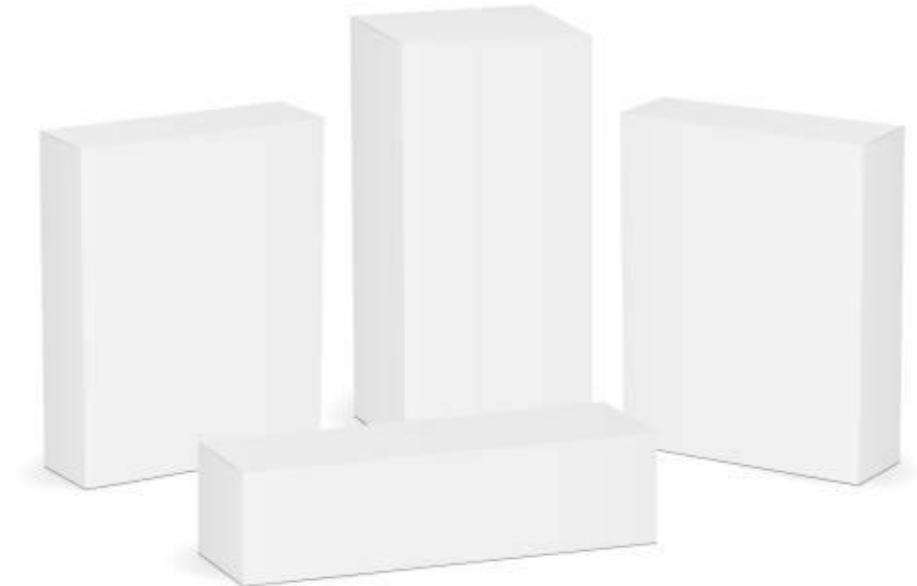


## Discriminative AI



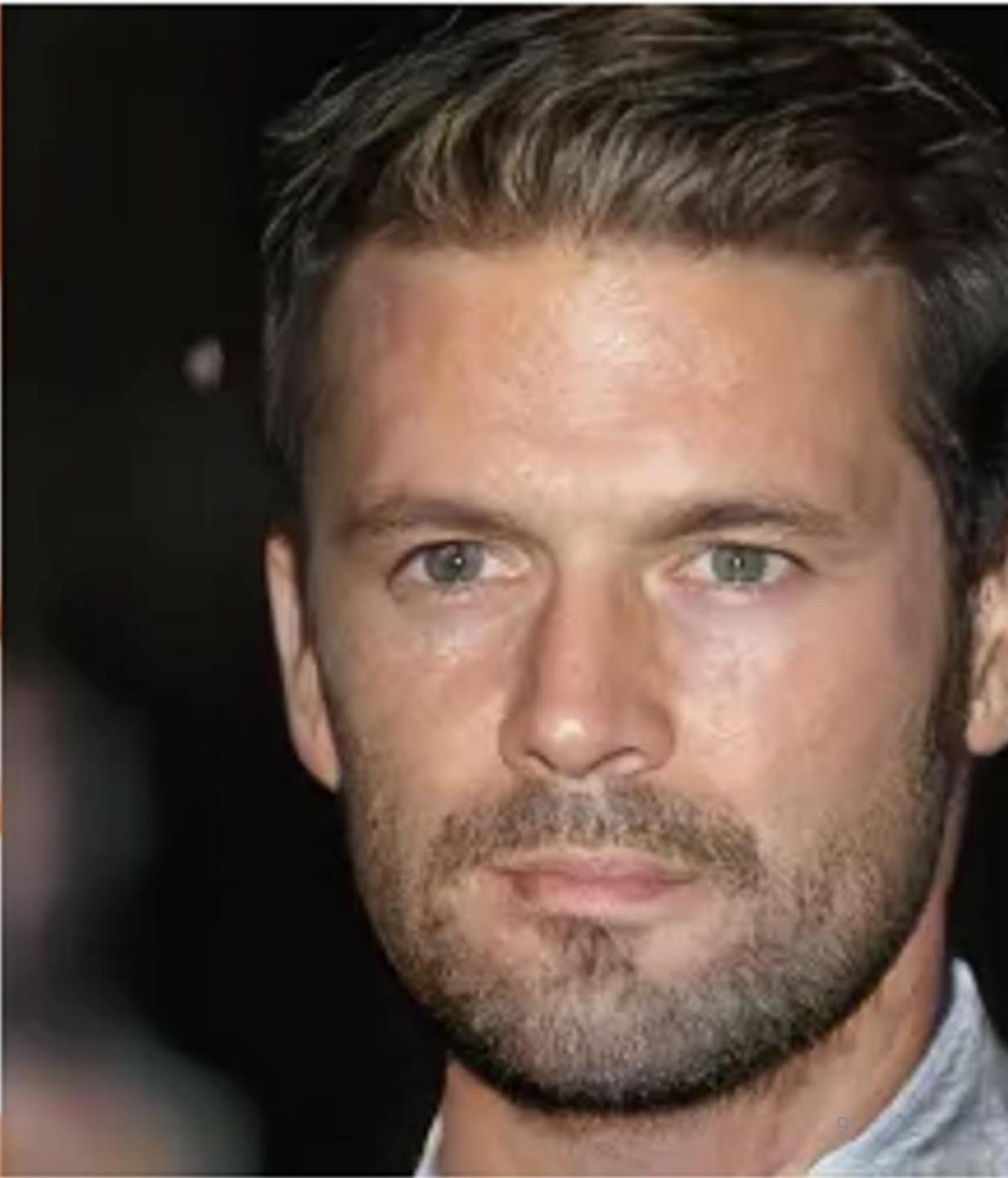
## Generative AI

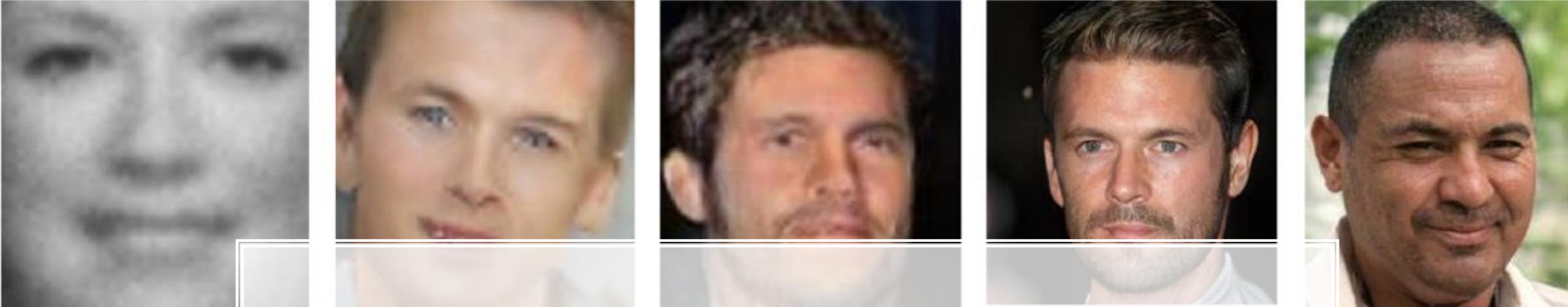




# Core Principle – Representation Learning

# Applications





2014

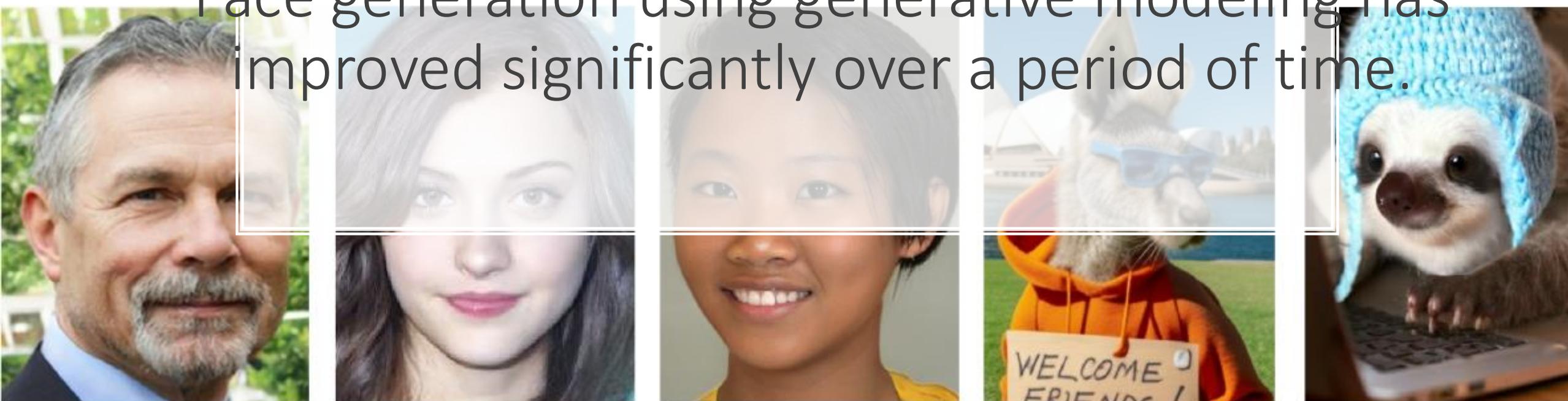
2015

2016

2017

2018

Face generation using generative modeling has improved significantly over a period of time.



2019

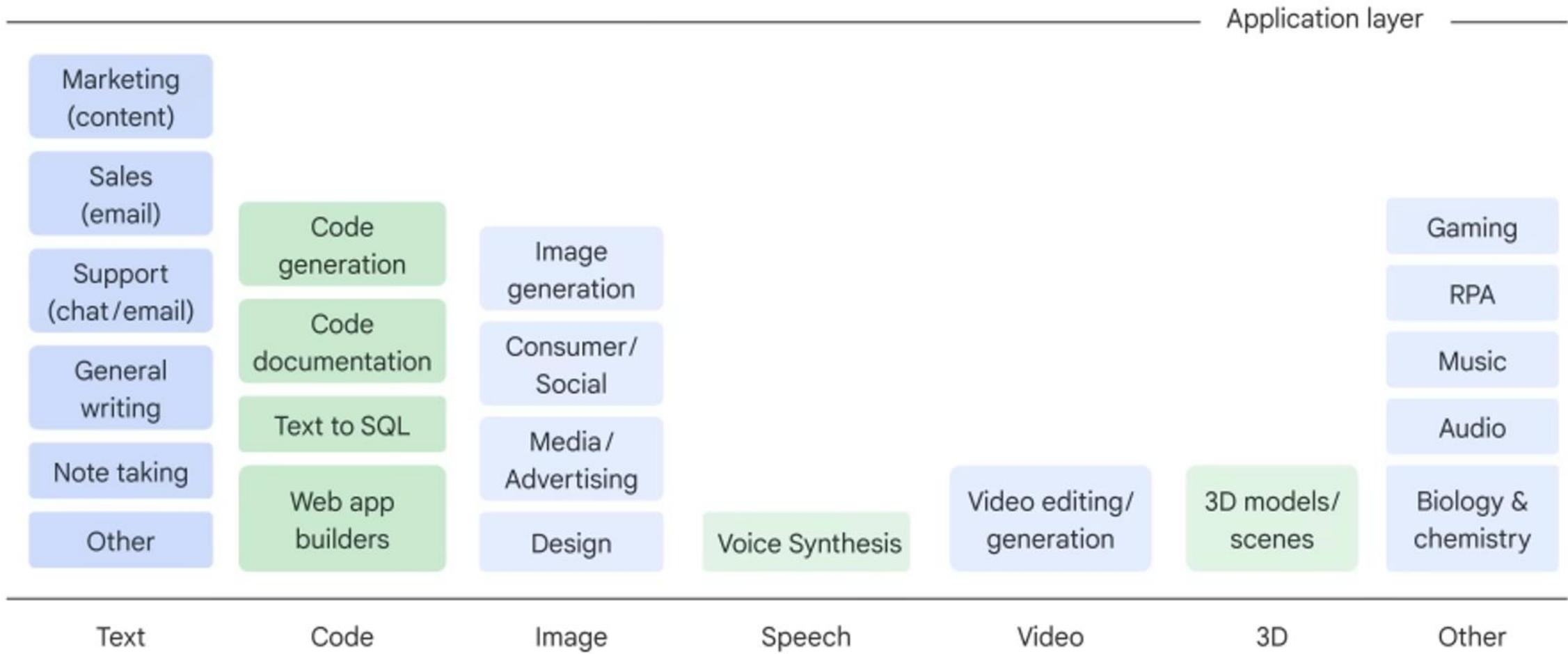
2020

2021

2022

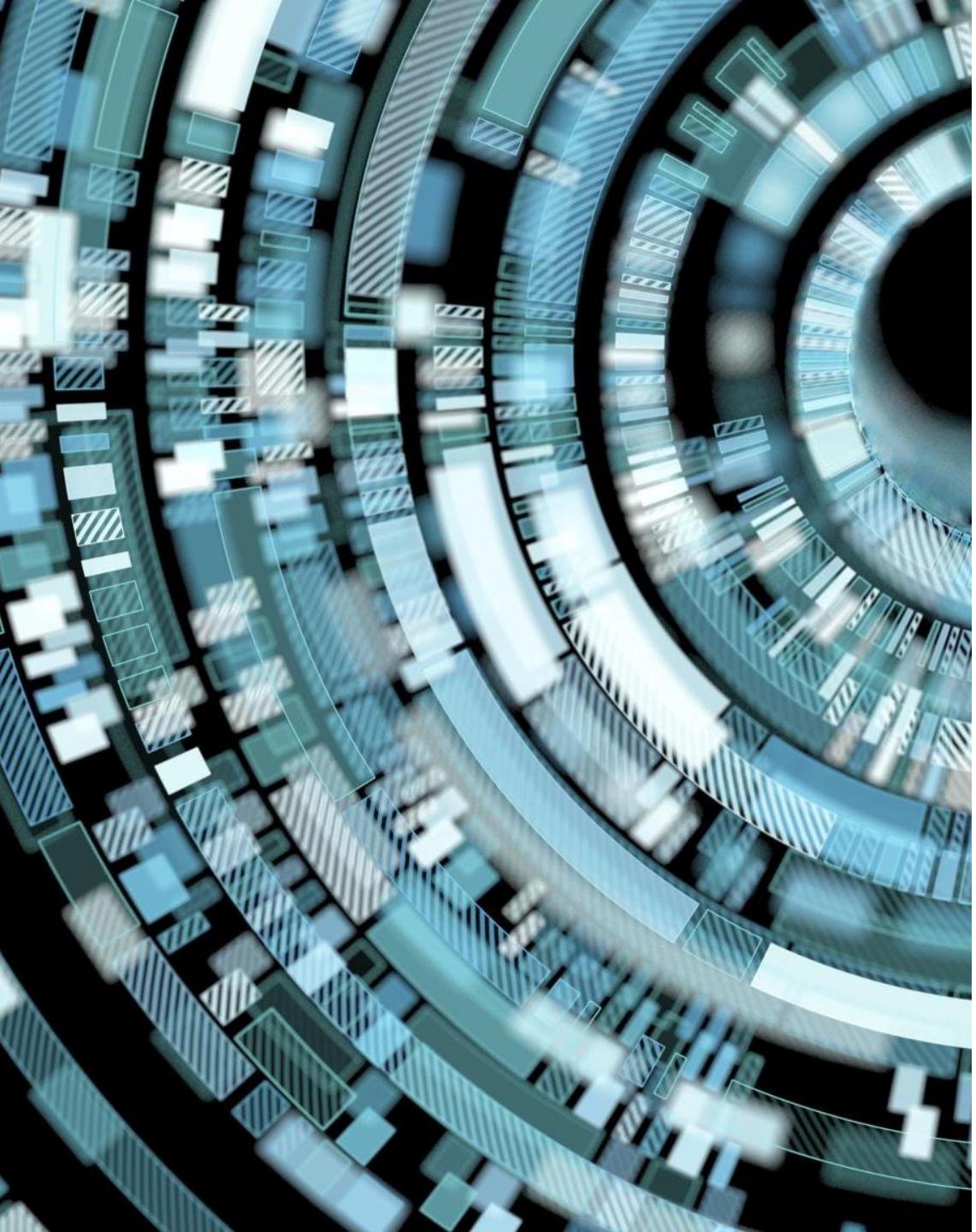
2023

# Generative AI Application Landscape



A large, abstract graphic on the left side of the slide features a complex pattern of overlapping triangles in shades of red, pink, and light purple, creating a polygonal, crystalline effect.

# Large Language Models (LLMs)



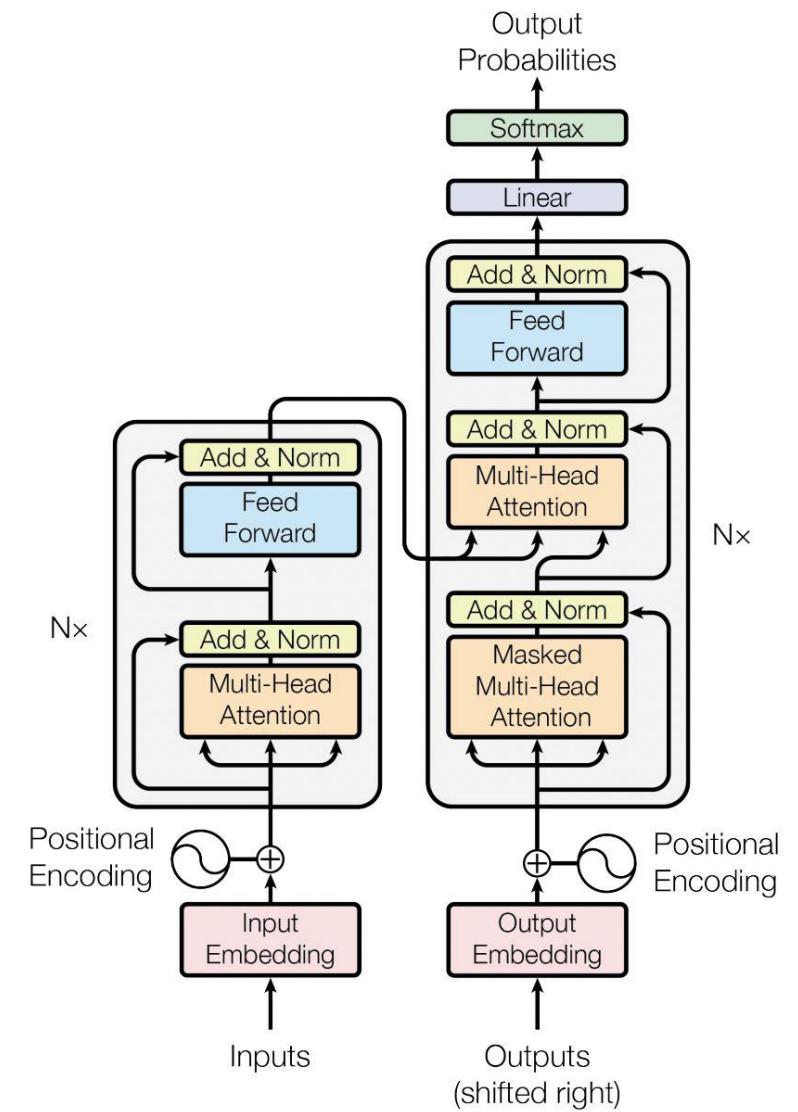
# LLMs

- A Language Model
- Consisting of Neural network with many parameters
  - (typically, billions to trillions of weights or more)
- Trained on large quantities (trillions of words) of unlabelled text.
  - Wikipedia, GitHub, Common Crawl, The Pile etc.
- Using self-supervised learning

# LLMs – Based on?

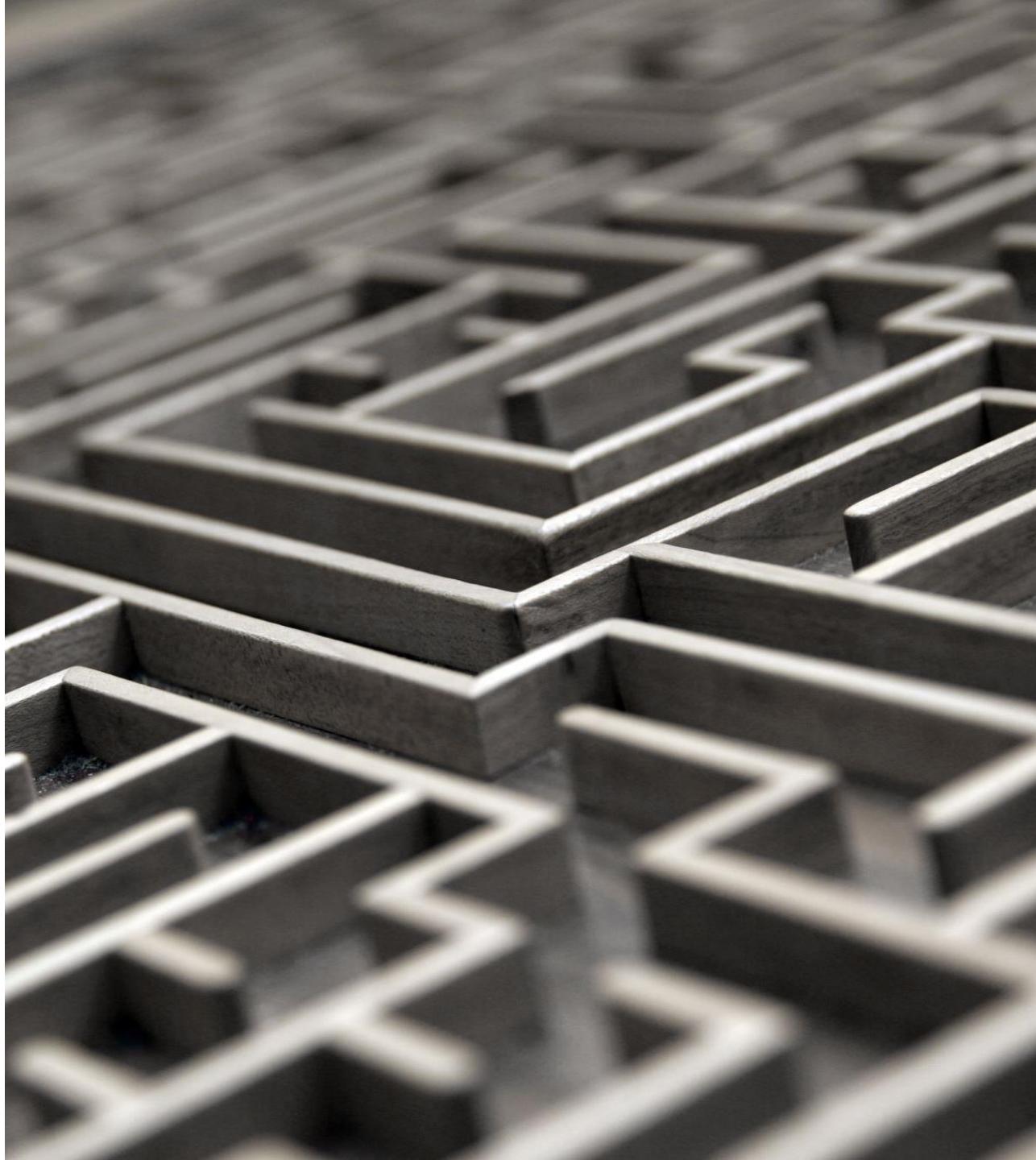
**Attention is all you need**

<https://ig.ft.com/generative-ai/>

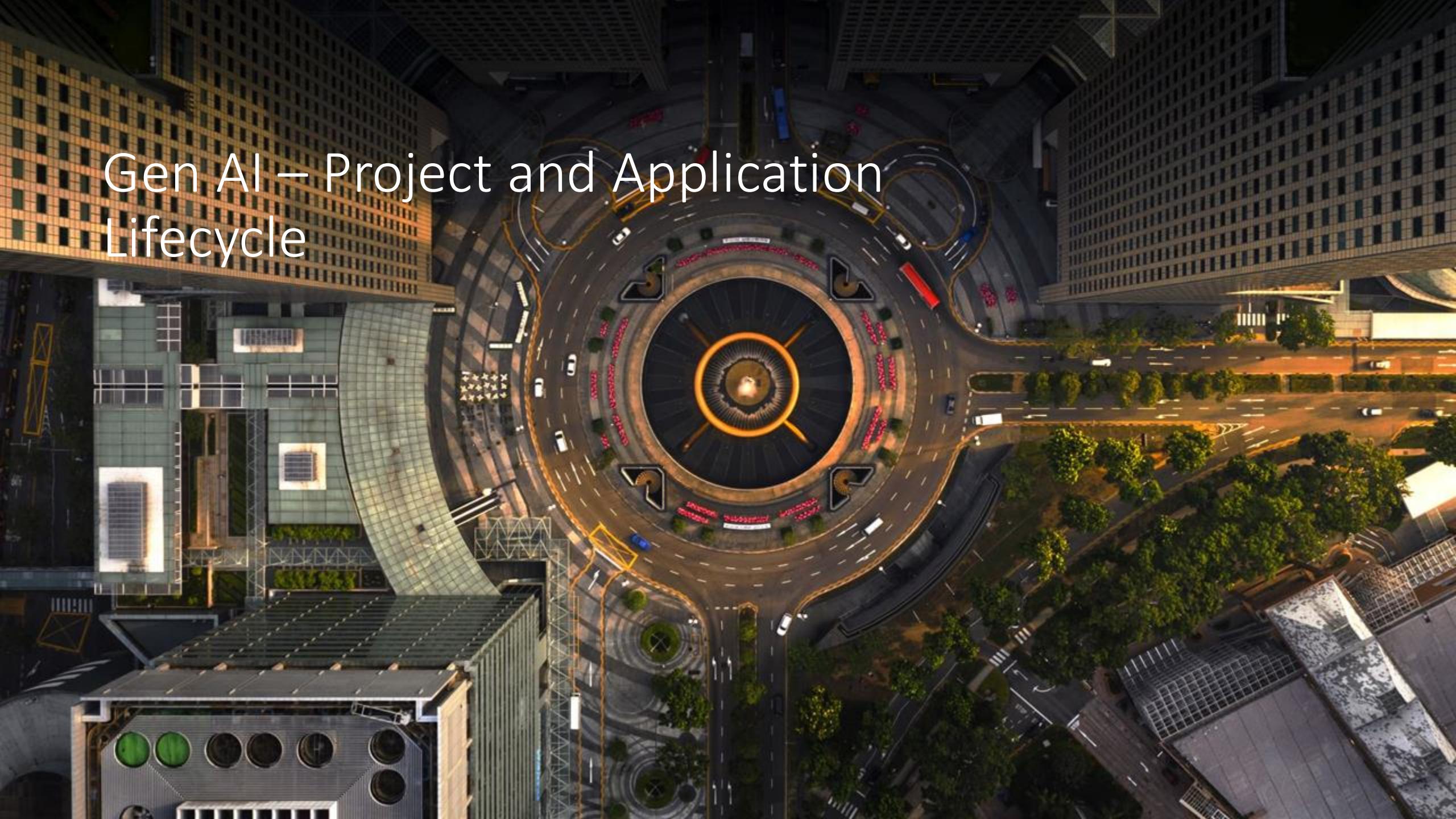


# LLMs – Characteristics

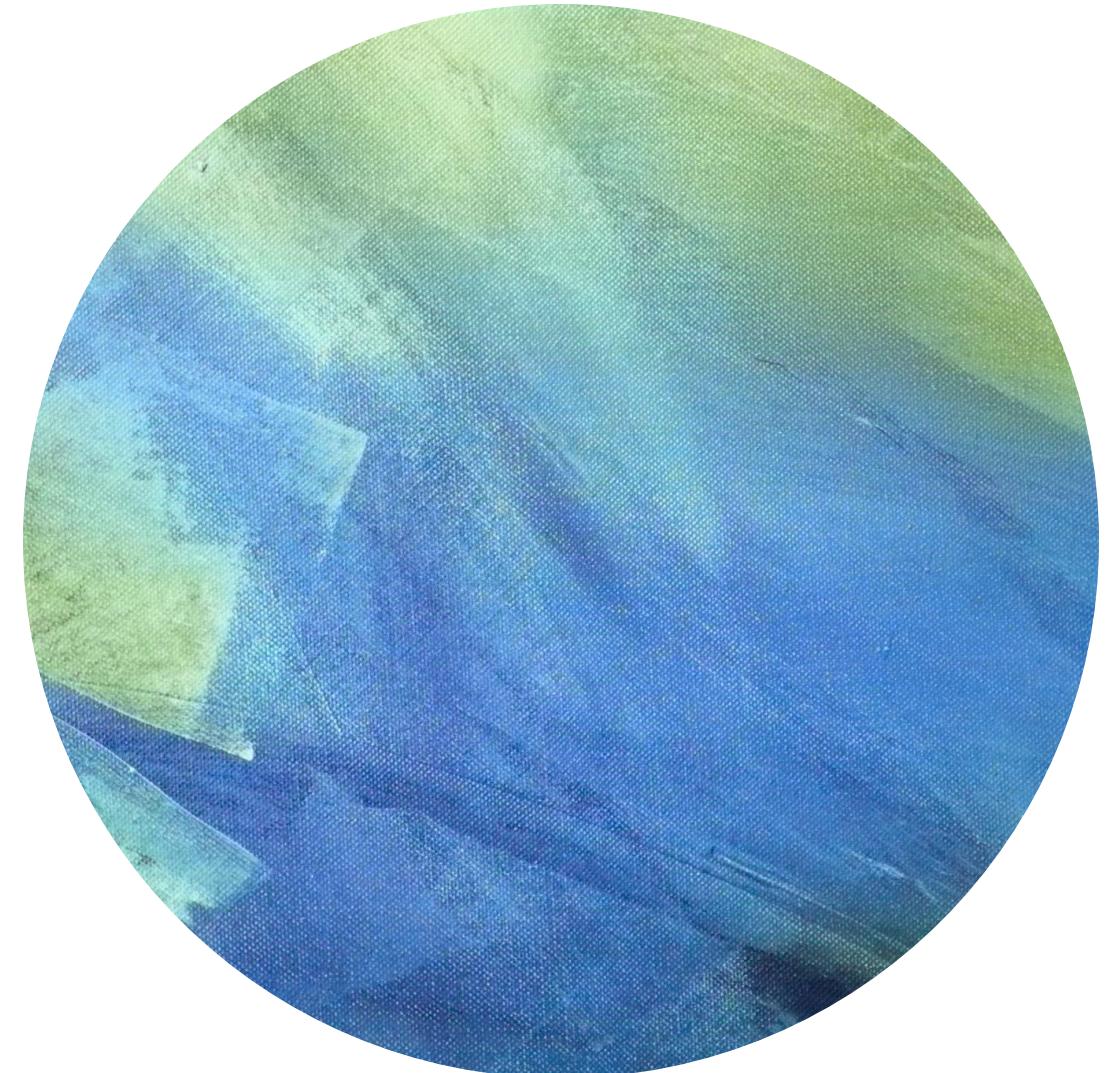
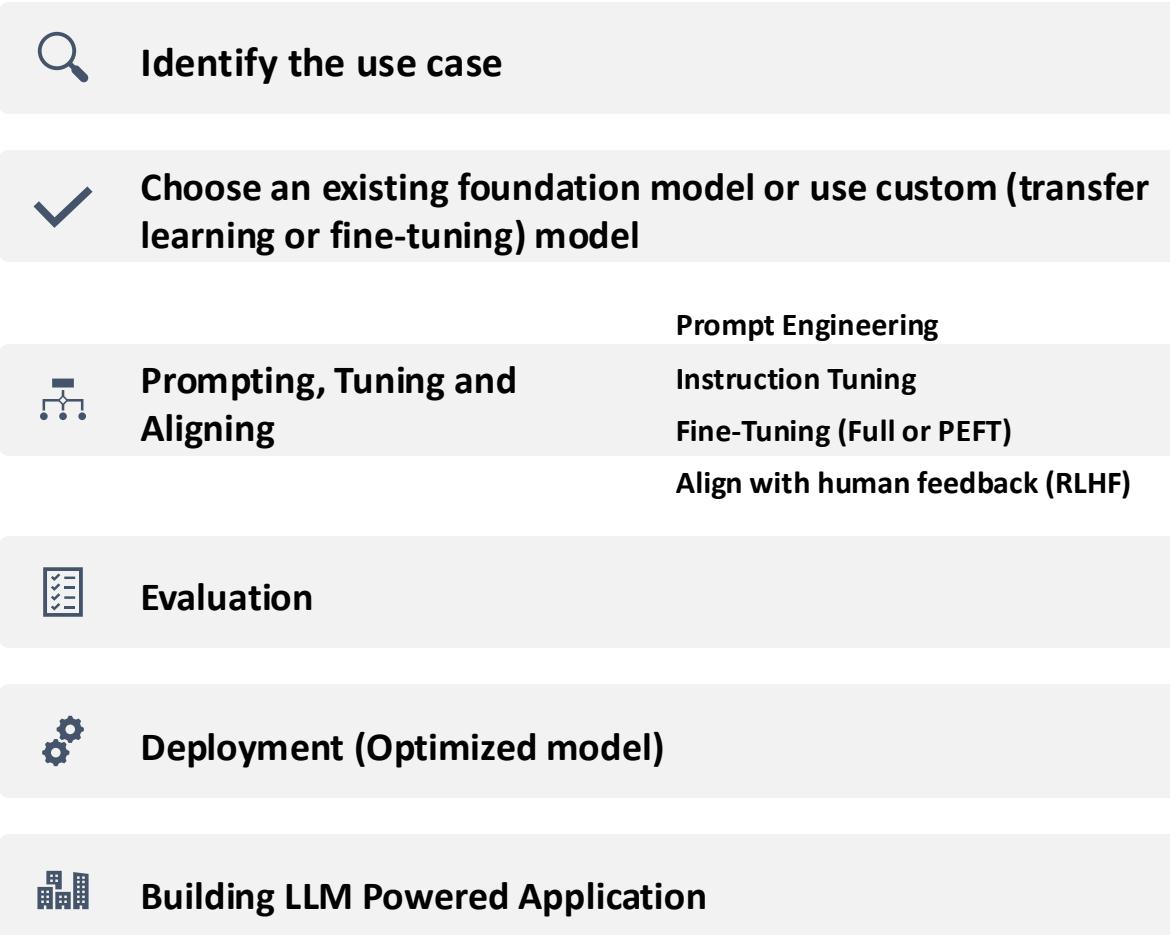
- General purpose models
- Stateless
- Stochastic



# Gen AI – Project and Application Lifecycle



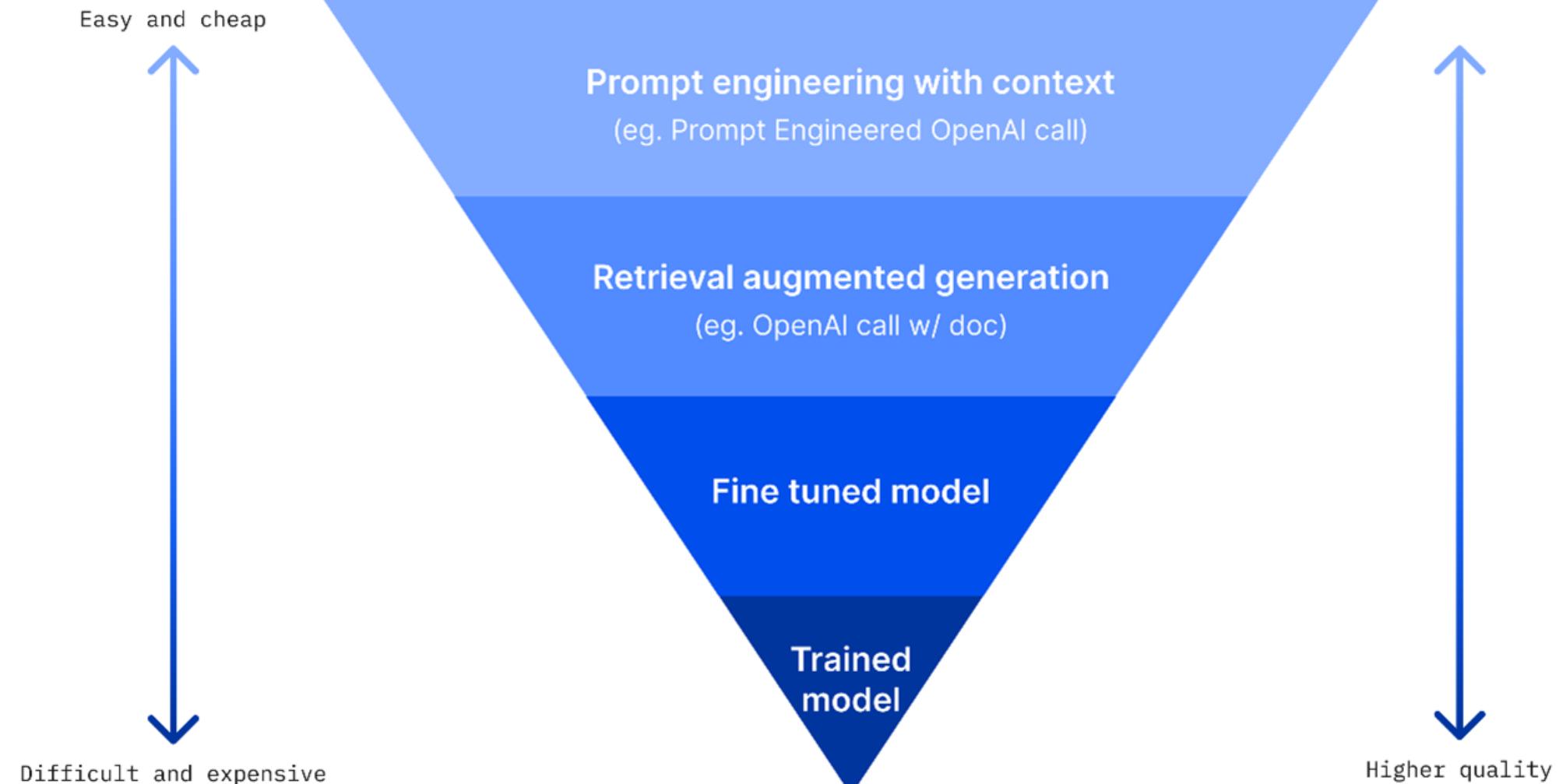
# Gen AI Project/Product Lifecycle



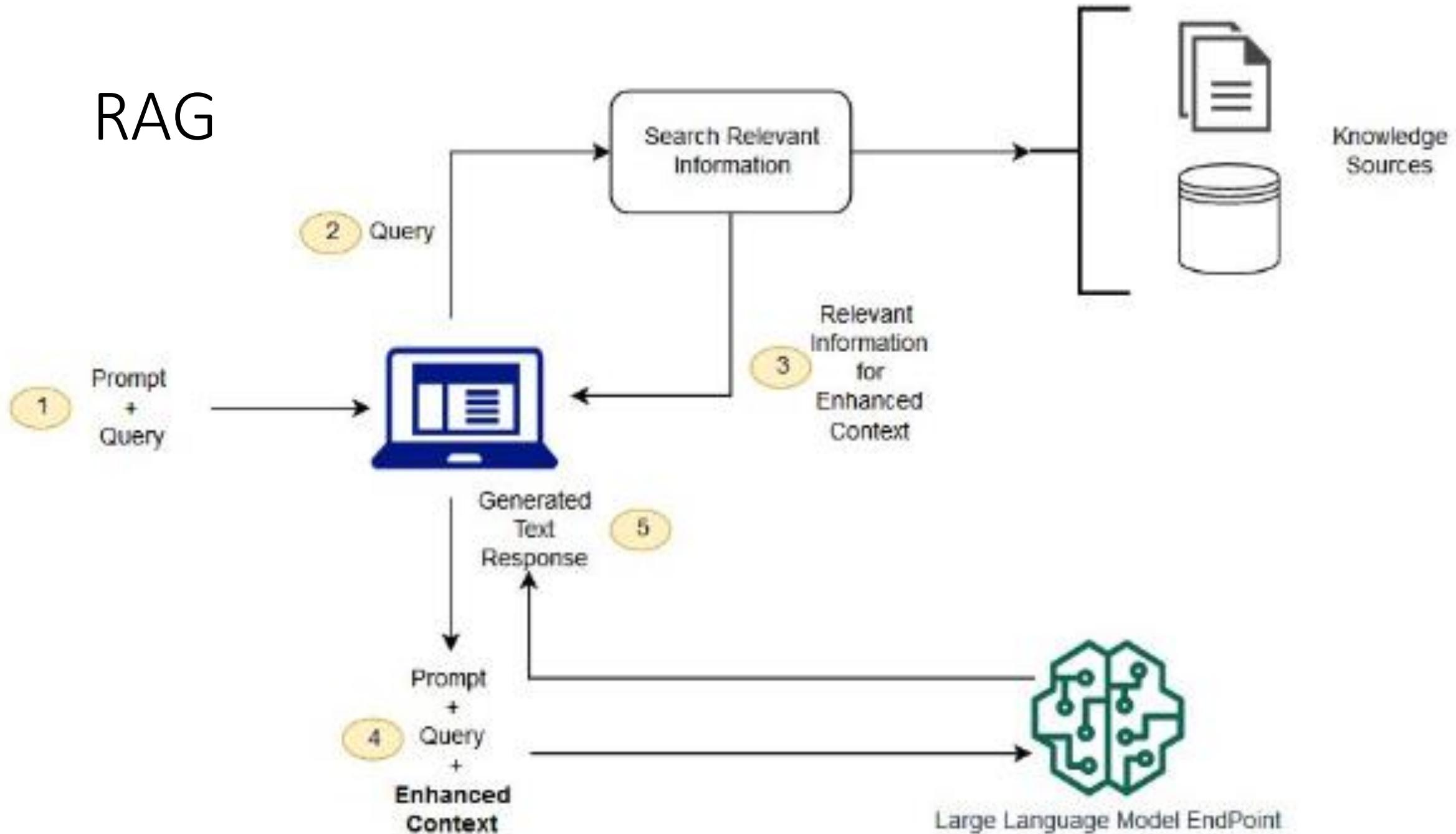
**Iterative in nature**

# Cheat Sheet - Time and effort in the lifecycle

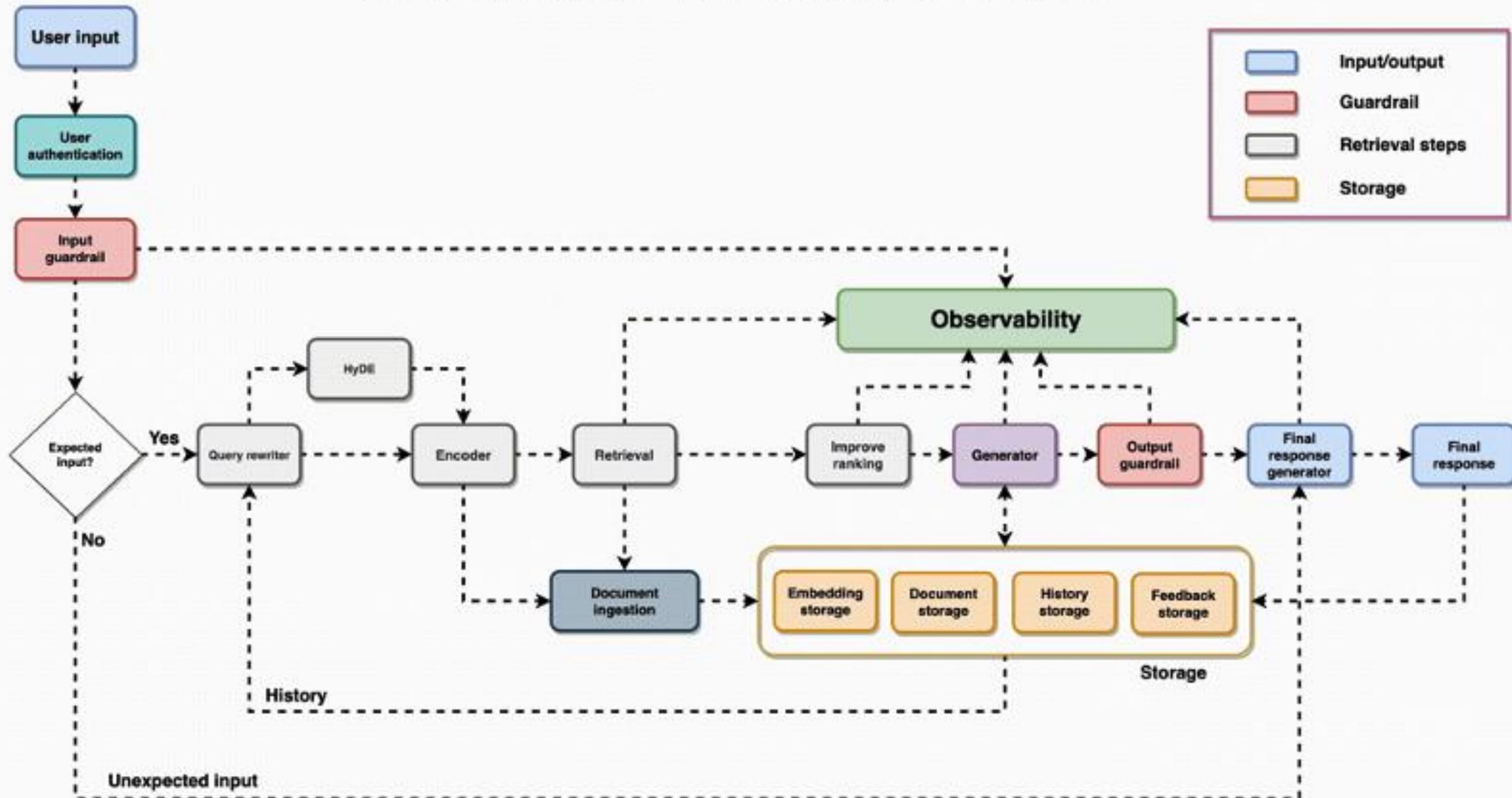
	Pre-training	Prompt engineering	Prompt tuning and fine-tuning	Reinforcement learning/human feedback	Compression/optimization/deployment
Training duration	Days to weeks to months	Not required	Minutes to hours	Minutes to hours similar to fine-tuning	Minutes to hours
Customization	Determine model architecture, size and tokenizer.  Choose vocabulary size and # of tokens for input/context  Large amount of domain training data	No model weights  Only prompt customization	Tune for specific tasks  Add domain-specific data  Update LLM model or adapter weights	Need separate reward model to align with human goals (helpful, honest, harmless)  Update LLM model or adapter weights	Reduce model size through model pruning, weight quantization, distillation  Smaller size, faster inference
Objective	Next-token prediction	Increase task performance	Increase task performance	Increase alignment with human preferences	Increase inference performance
Expertise	High	Low	Medium	Medium-High	Medium



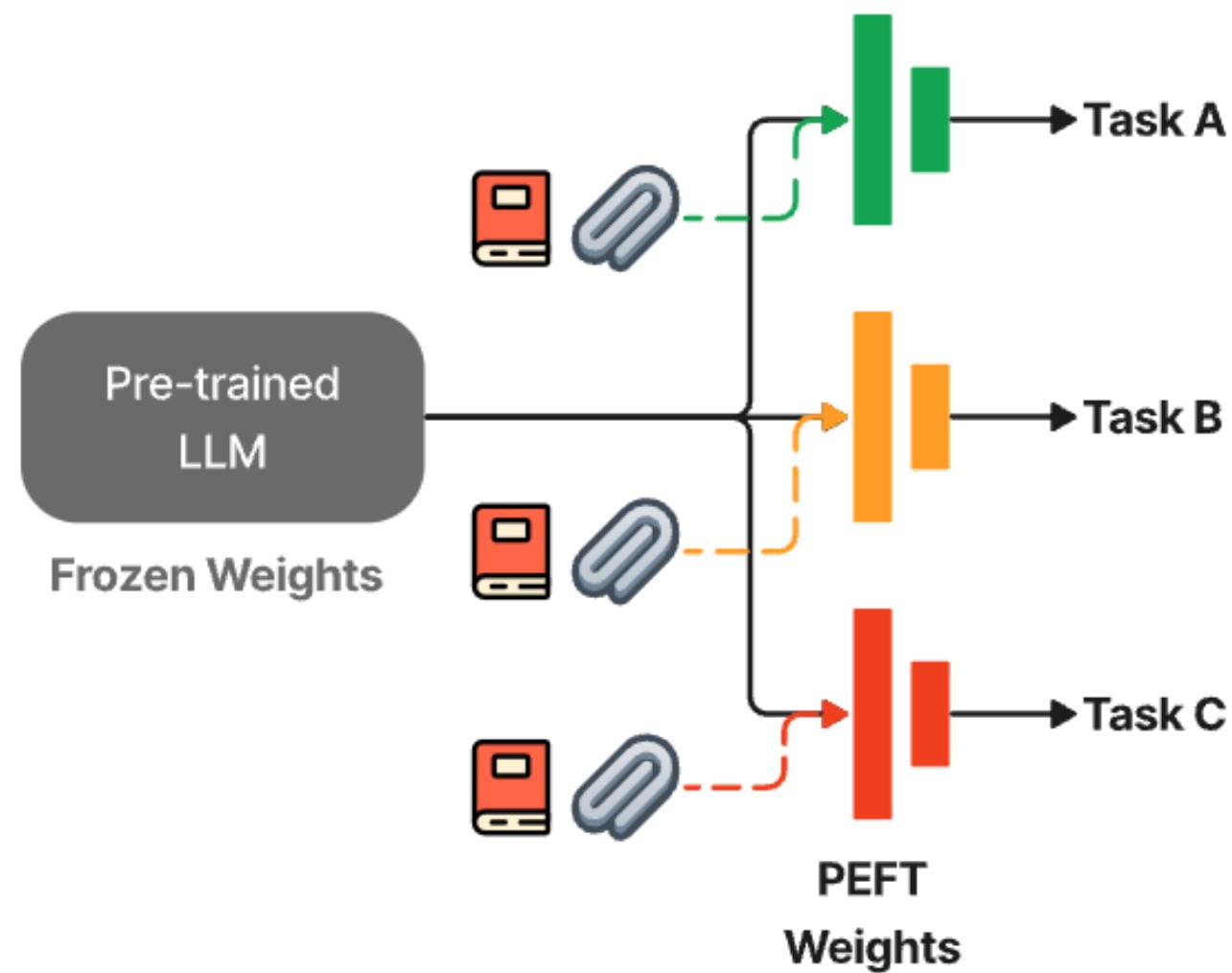
# RAG



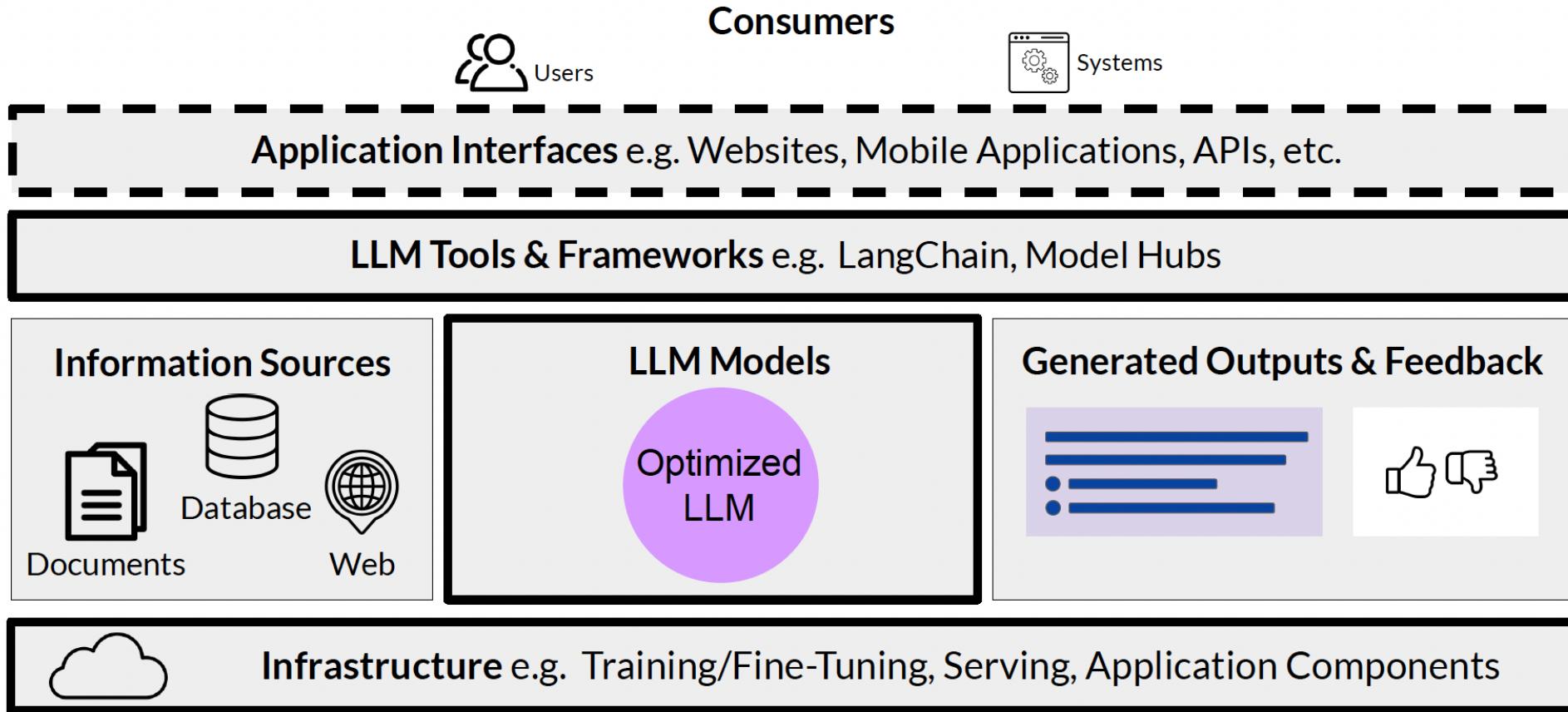
# Architecture For Enterprise RAG

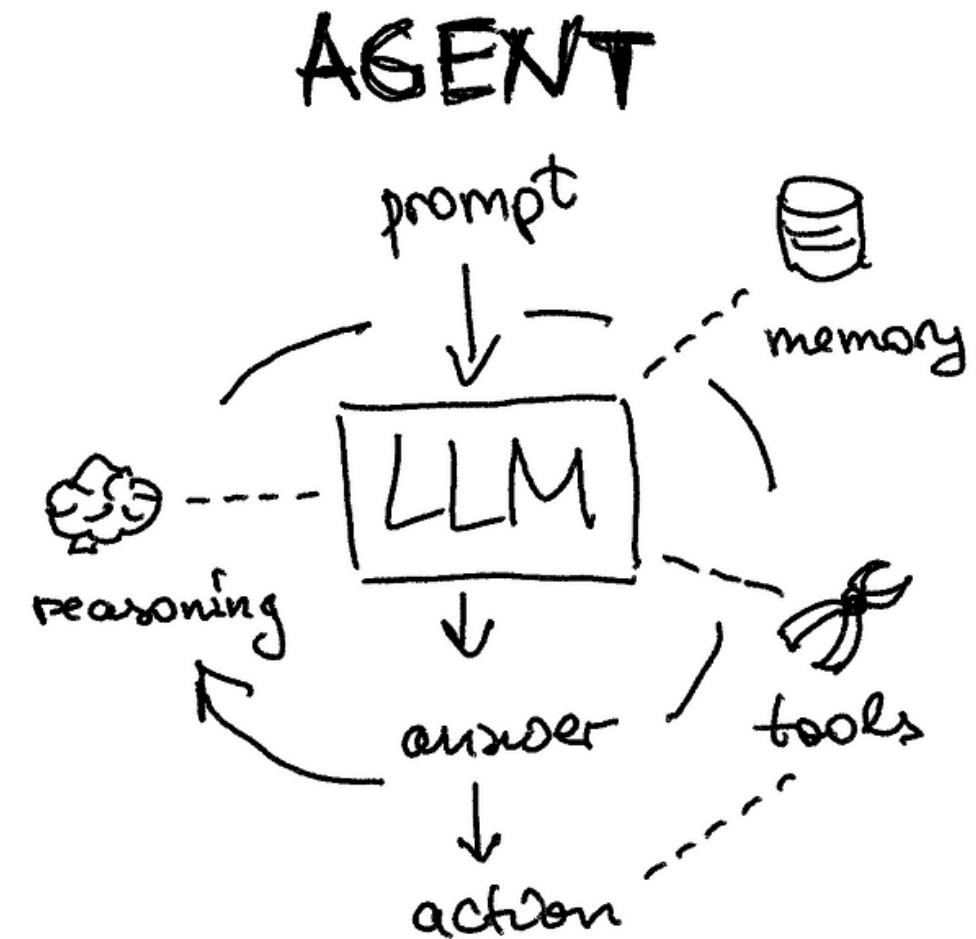
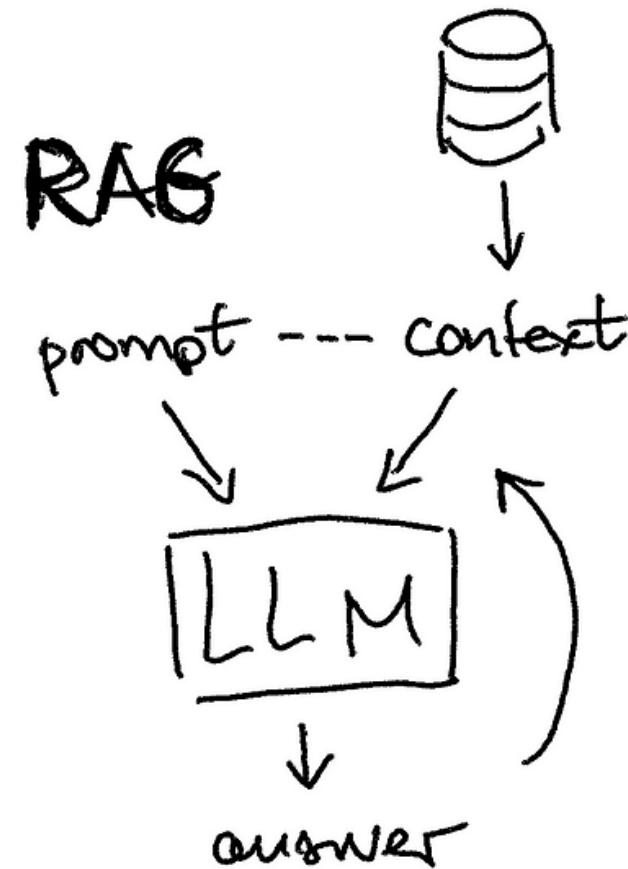
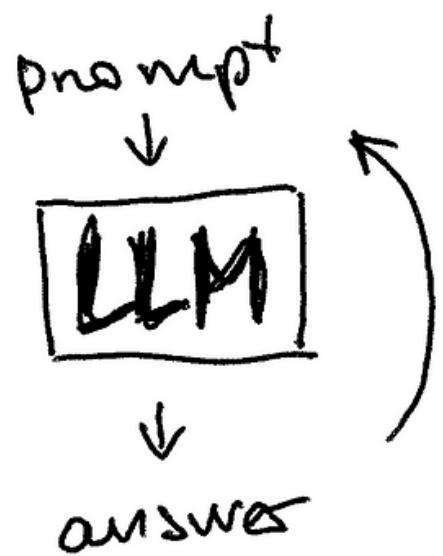


## Parameter Efficient Fine-Tuning (PEFT)

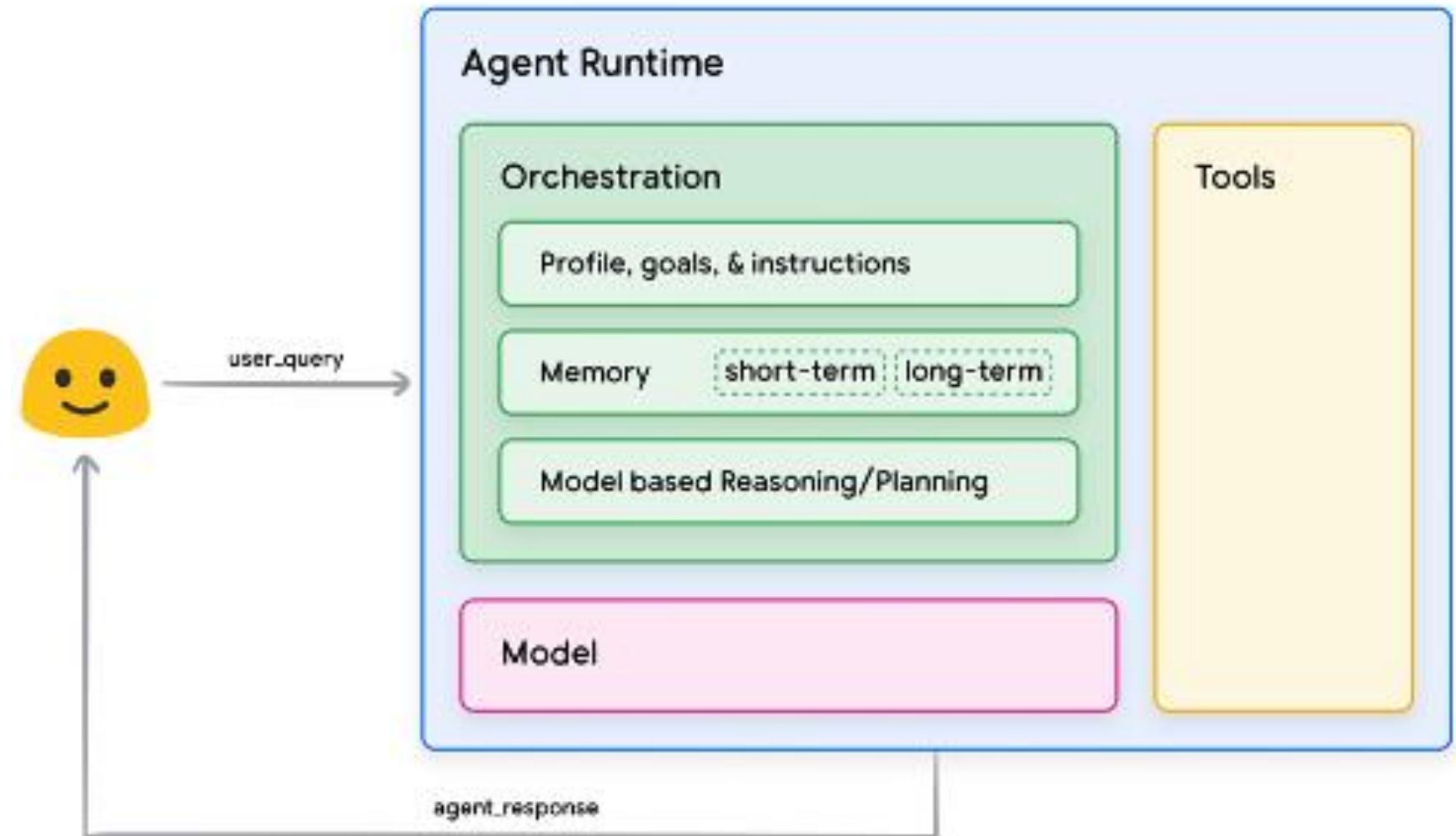


# Building generative applications





General  
agent  
architecture  
and  
components



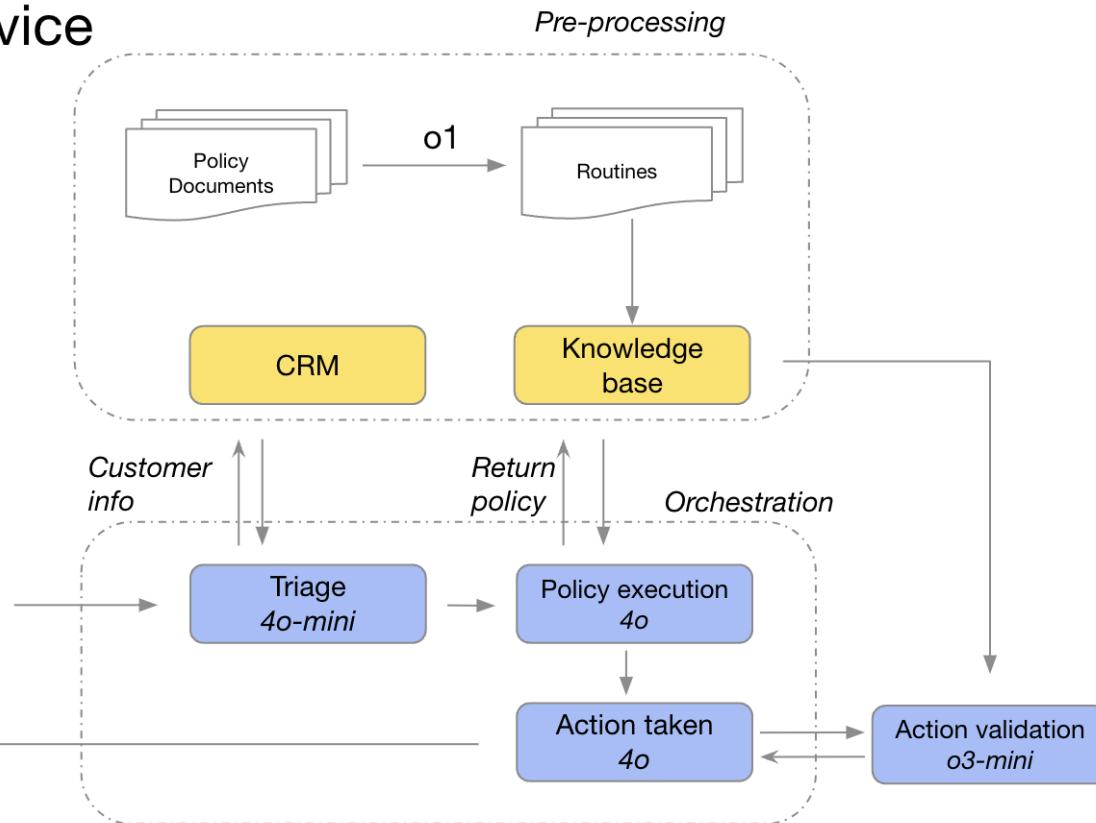
# How to choose?

GPT-4o and GPT-4o mini models triage order details with customer information, identify the order issues and the return policy, and then feed all of these data points into o3-mini to make the final decision about the viability of the return based on policy.

## Example: Customer service

- Customer
- Tools
- Agents

“My order was damaged.”



# Agentic AI is picking up where GenAI and RPA left off

## RPA

"I follow instructions exactly"

RPA is the **task automation** that replaces manual effort in routine, rule-based processes.

### Key Characteristics:

- Executes structured, rule-based processes
- Performs repetitive digital tasks with precision
- Operates within defined system boundaries
- Follows exact step-by-step procedures

## Low code?

## GenAI

"I can create based on prompts"

GenAI is a **productivity amplifier** that supports and enhances human work, transforming workflows without fully replacing human decision-making.

### Key Characteristics:

- Assists with specific tasks (writing, analysis, coding)
- Requires human direction and oversight
- Improves individual productivity
- Works within existing job roles

## Agentic AI

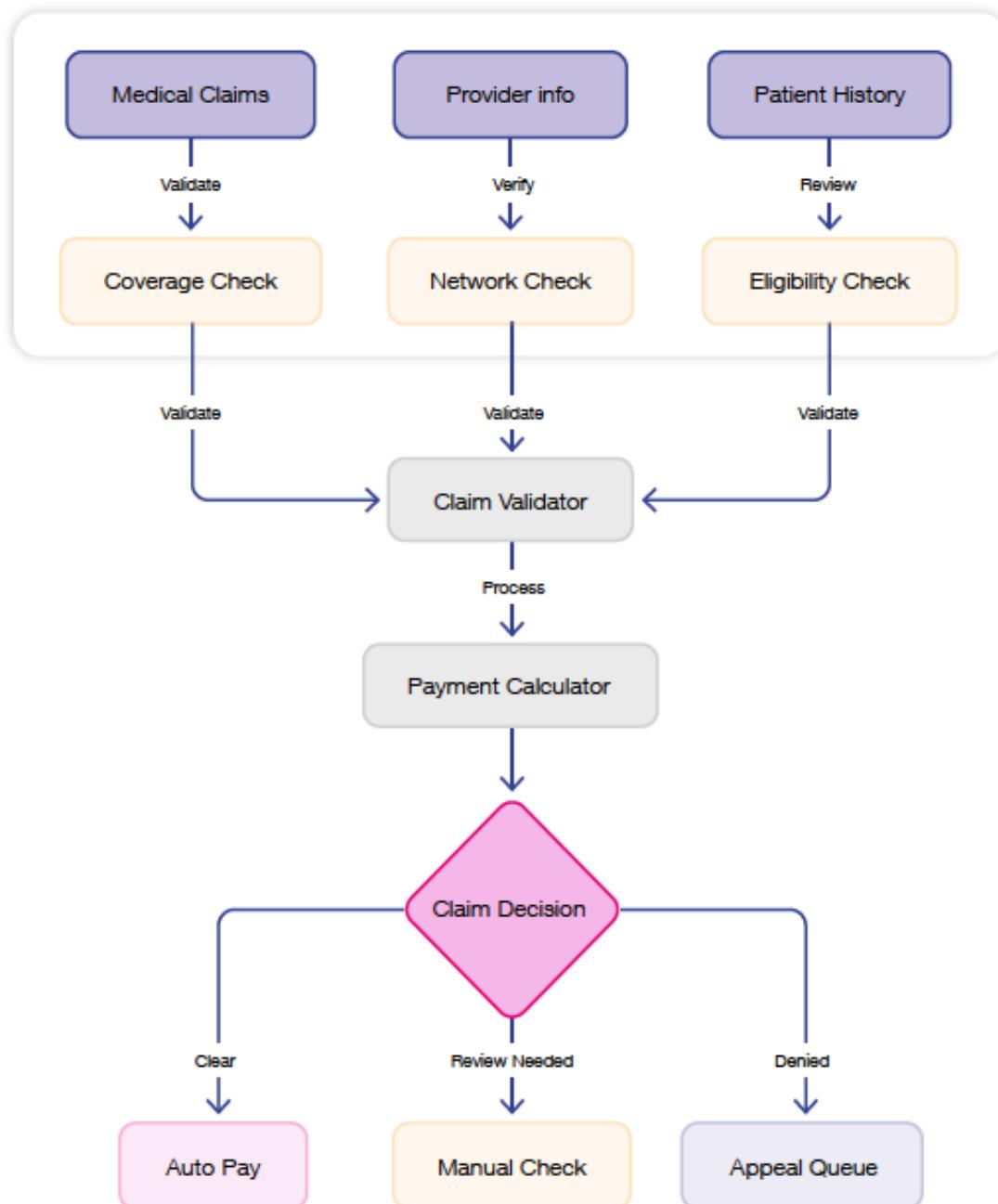
"I can understand goals and figure out how to achieve them"

Agentic AI is a **collaborative actor** that autonomously executes and coordinates complex tasks.

### Key Characteristics:

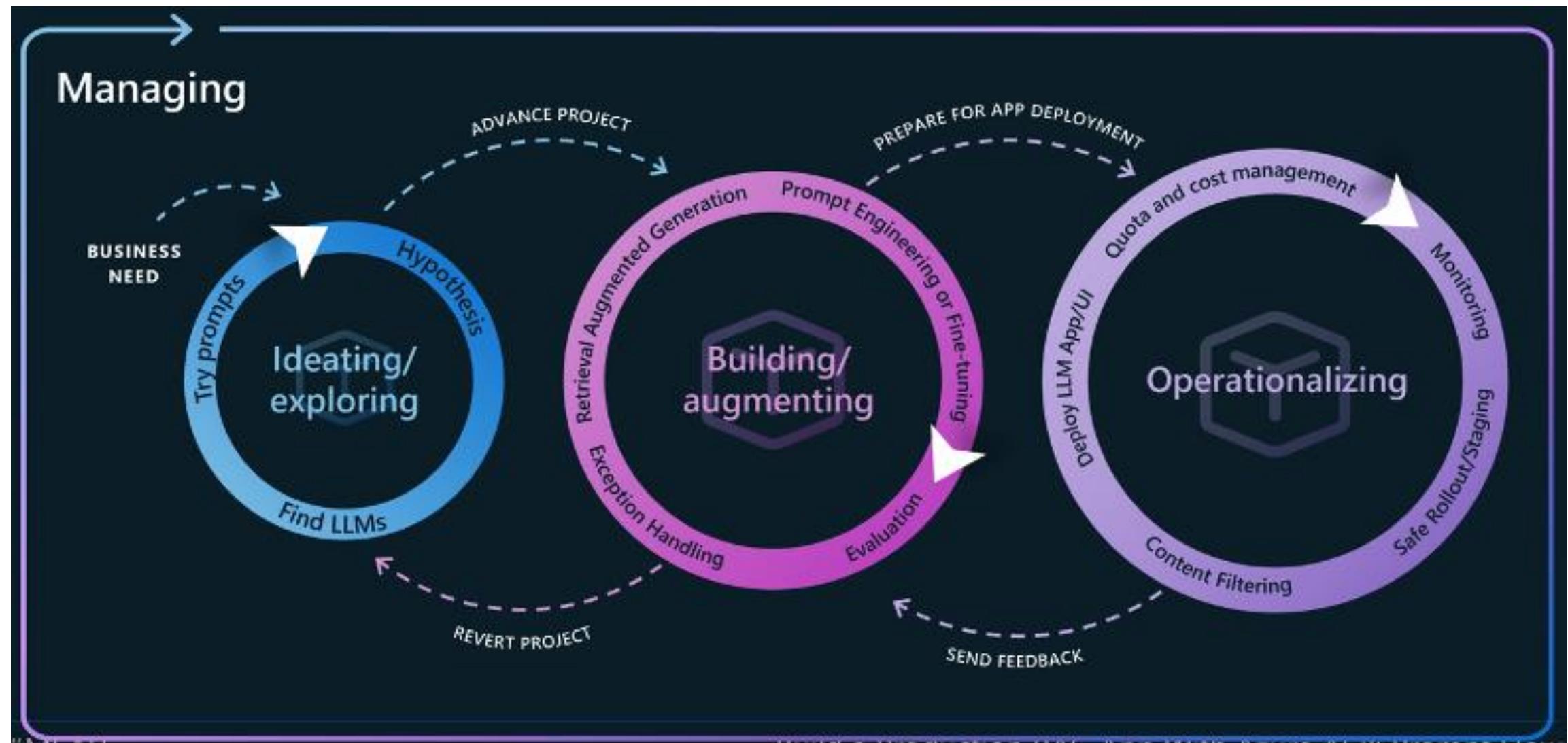
- Acts as virtual coworker completing end-to-end processes
- Self-directs and coordinates multiple tasks
- Transforms entire workflows
- Creates new organizational paradigms

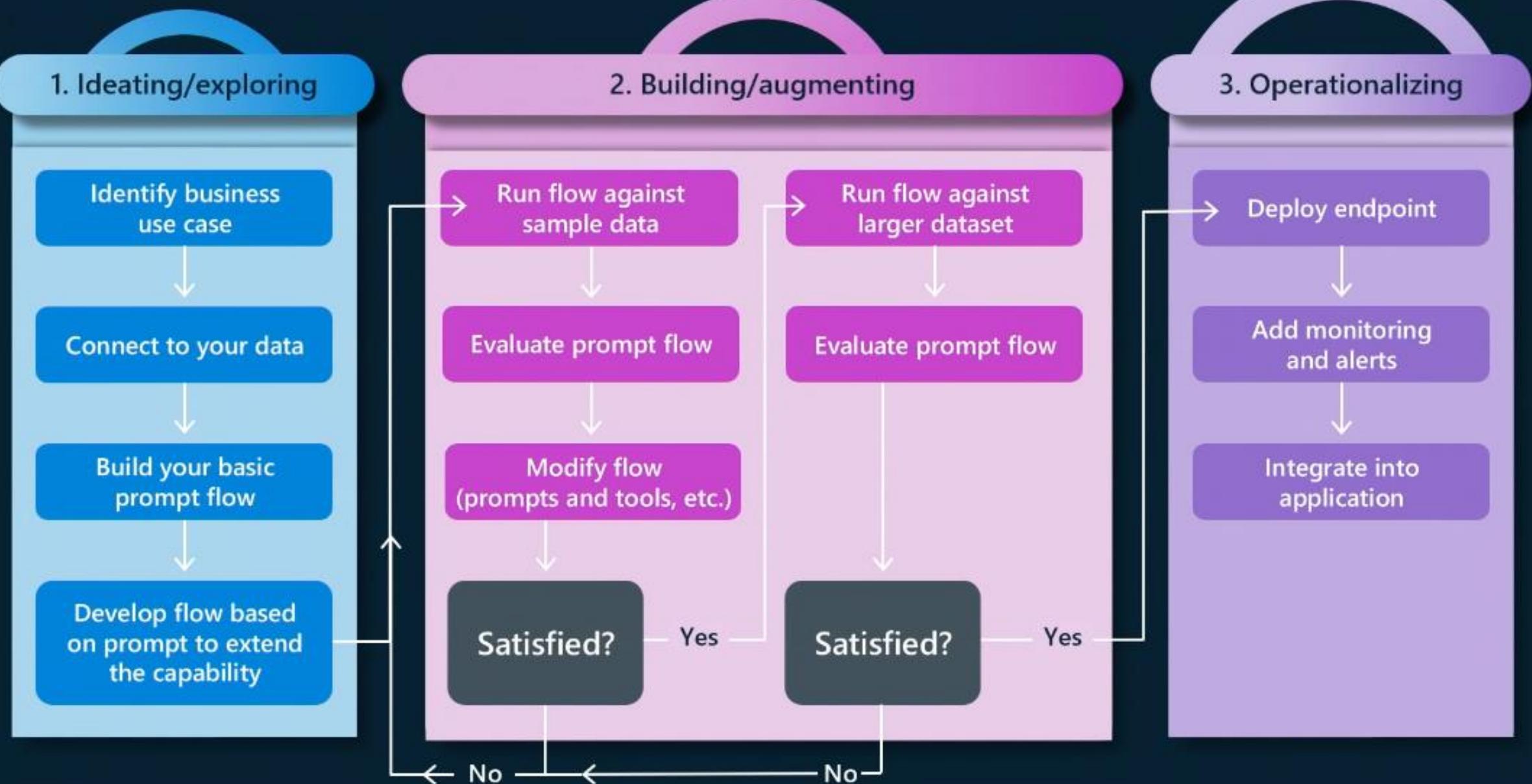
# Advancing the Claims Processing Agent



# AI Agent Challenges and Solutions

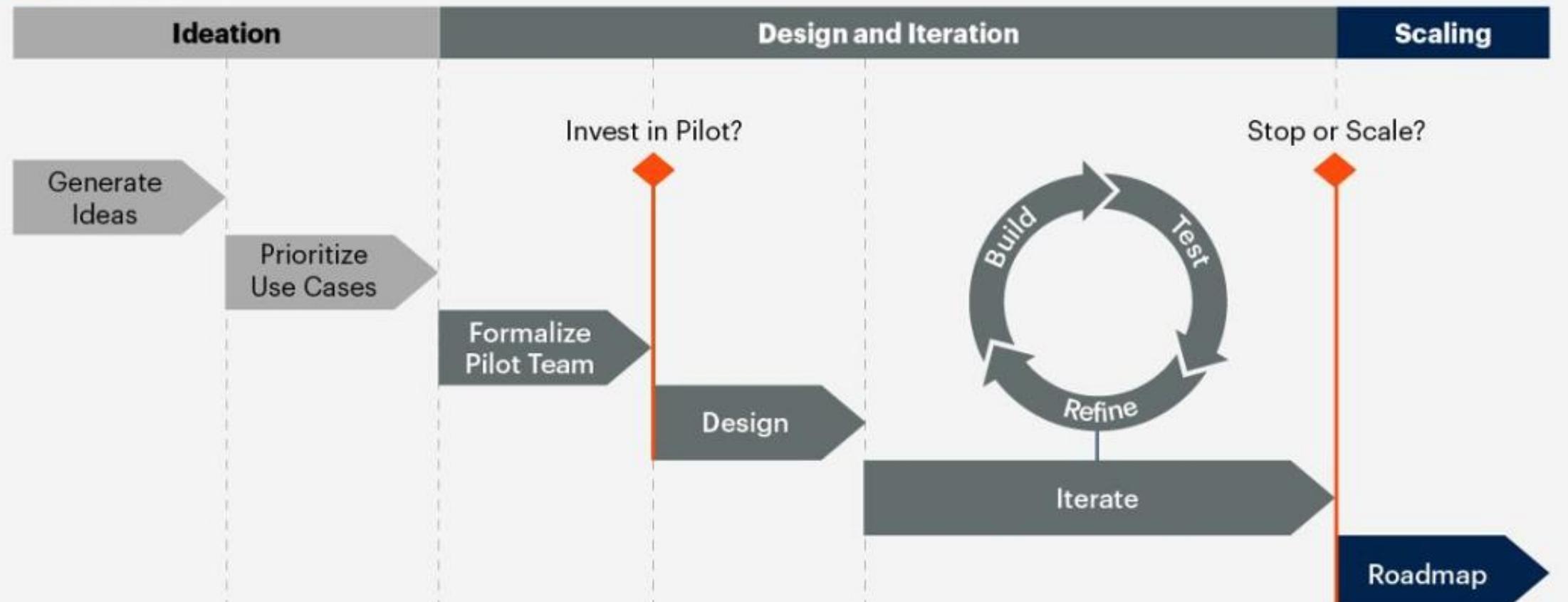
DEVELOPMENT ISSUES	LLM ISSUES	PRODUCTION ISSUES
Poorly Defined Prompts <ul style="list-style-type: none"><li>• Define Clear Objectives</li><li>• Craft Detailed Personas</li><li>• Use Effective Prompting</li></ul>	Difficult to Steer <ul style="list-style-type: none"><li>• Specialized Prompts</li><li>• Hierarchical Design</li><li>• Fine-Tuning Models</li></ul>	Guardrails <ul style="list-style-type: none"><li>• Rule-Based Filters &amp; Validation</li><li>• Human-in-the-Loop Oversight</li><li>• Ethical &amp; Compliance Frameworks</li></ul>
Evaluation Challenges <ul style="list-style-type: none"><li>• Continuous Evaluation</li><li>• Use Real-World Scenarios</li><li>• Incorporate Feedback Loops</li></ul>	High Cost of Running <ul style="list-style-type: none"><li>• Reduce Context Size</li><li>• Use Smaller Models</li><li>• Cloud-Based Solutions</li></ul>	Agent Scaling <ul style="list-style-type: none"><li>• Scalable Architectures</li><li>• Resource Management</li><li>• Monitor Performance</li></ul>
	Planning Failures <ul style="list-style-type: none"><li>• Task Decomposition</li><li>• Multi-Plan Selection</li><li>• Reflection and Refinement</li></ul>	Fault Tolerance <ul style="list-style-type: none"><li>• Redundancy</li><li>• Automated Recovery</li><li>• Stateful Recovery</li></ul>
	Reasoning Failures <ul style="list-style-type: none"><li>• Enhance Reasoning Capabilities</li><li>• Fine-Tune LLMs with Feedback</li><li>• Use Specialized Agents</li></ul>	Infinite Looping <ul style="list-style-type: none"><li>• Clear Termination Conditions</li><li>• Enhance Reasoning &amp; Planning</li><li>• Monitor Agent Behavior</li></ul>
	Tool Calling Failures <ul style="list-style-type: none"><li>• Define Clear Parameters</li><li>• Validate Tool Outputs</li><li>• Tool Selection Verification Loops</li></ul>	





# Generative AI Pilot Phases and Decision Points

◆ Decision Point

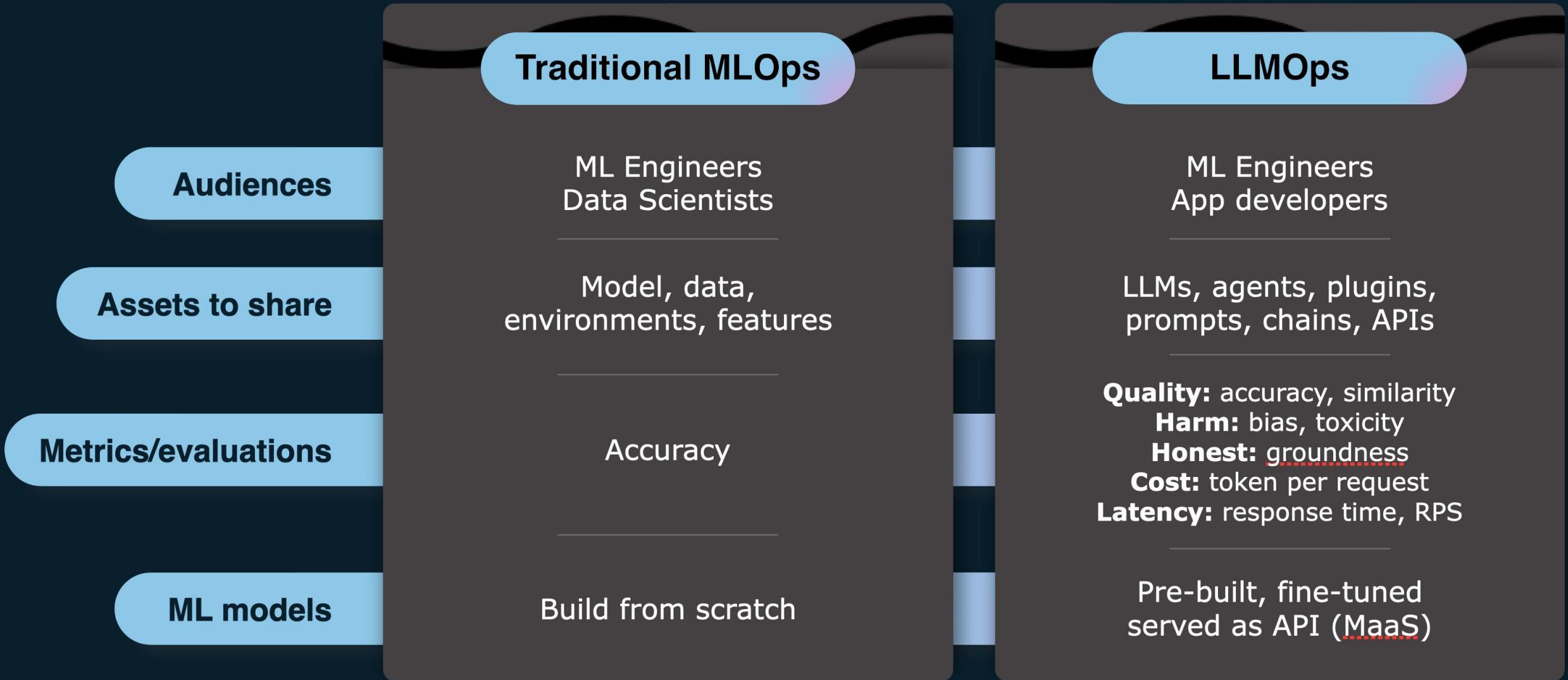


# GenAI - Paradigm Shift

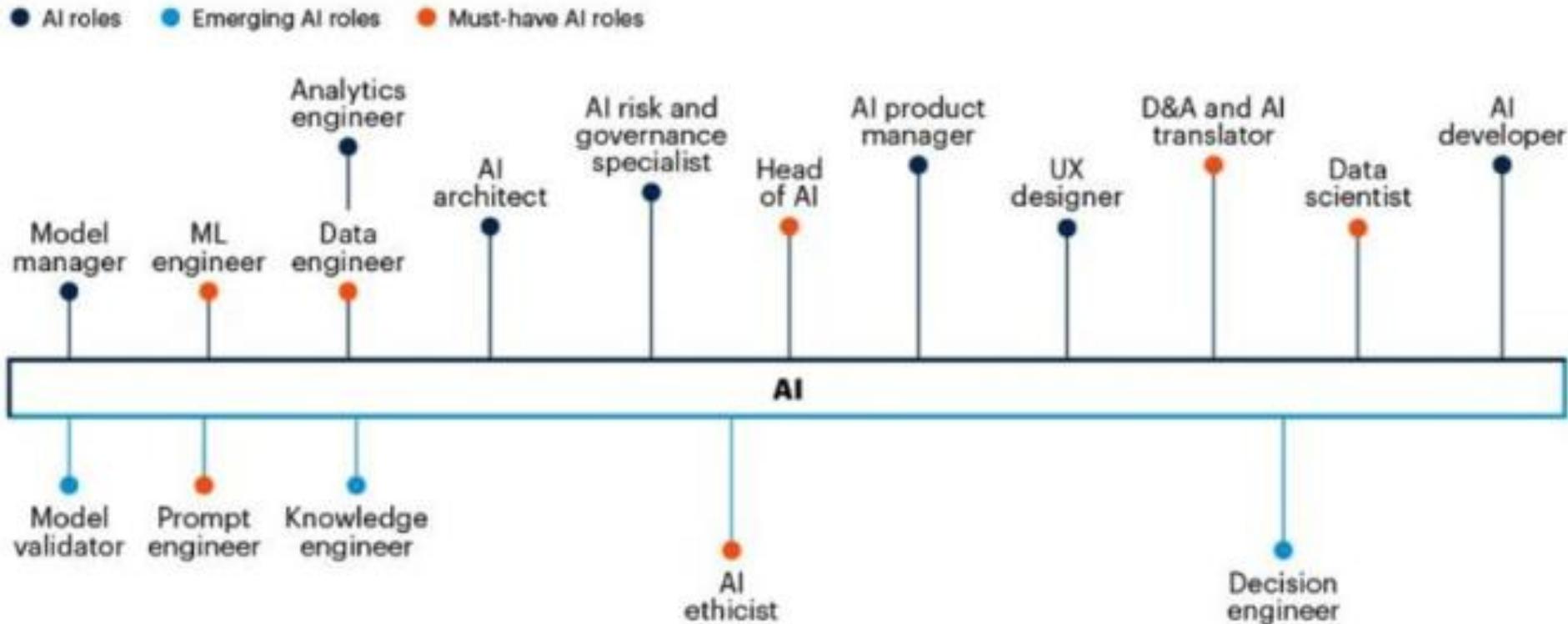
## Traditional software vs LLMs

	Traditional Software	LLM-based Applications
Behavior	Predefined rules	Probability + prediction
Output	Deterministic same input → same output	Non-deterministic Same input → many possible outputs
Testing	testing: 1 input, 1 correct output	1 input, many correct (and incorrect) outputs
Criteria	Evaluate as "right" or "wrong"	Evaluate on: <ul style="list-style-type: none"><li>• Accuracy</li><li>• Quality</li><li>• Consistency</li><li>• Bias</li><li>• Toxicity</li><li>• and more...</li></ul>

# The paradigm shift—from MLOps to LLMOps

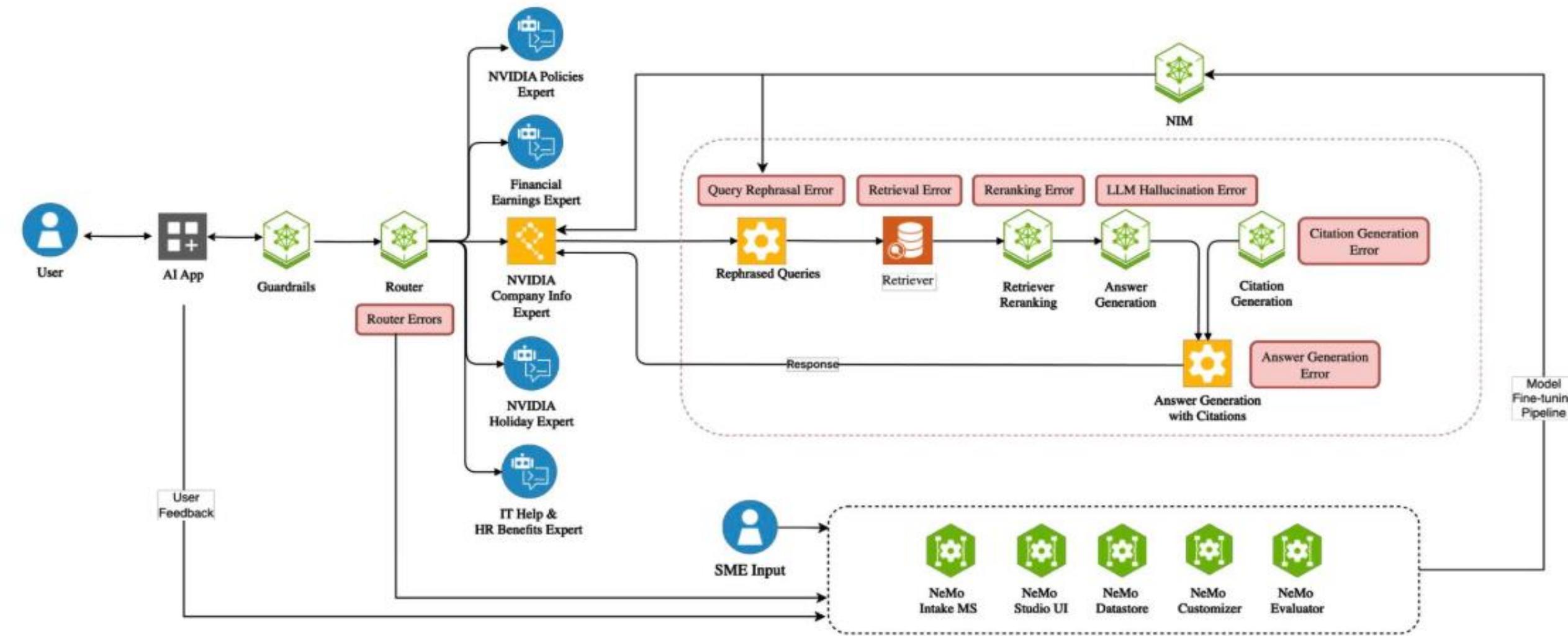


# NEW ROLES FOR AI



Rise of Problem Solvers & Business

# Architecture



## RAG going wrong

---

- Air Canada Loses Court Case After Its Chatbot Hallucinated Fake Policies To a Customer
- The airline argued that the chatbot itself was liable. The court disagreed.



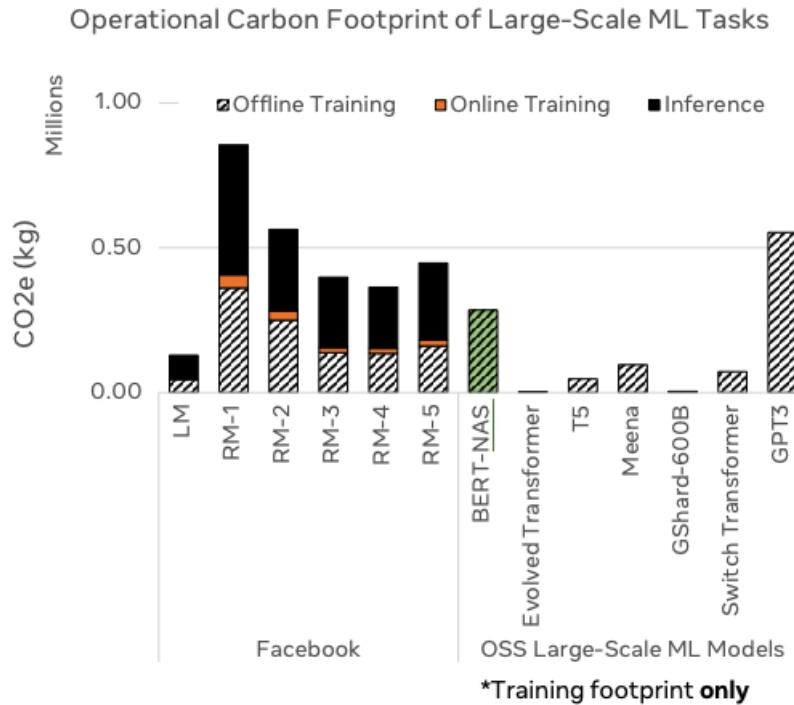
AIR CANA D

# Bias

AI was asked to  
create a picture  
of Mother  
Theresa fighting  
against poverty



# Sustainable AI: Carbon Footprint



## Common carbon footprint benchmarks

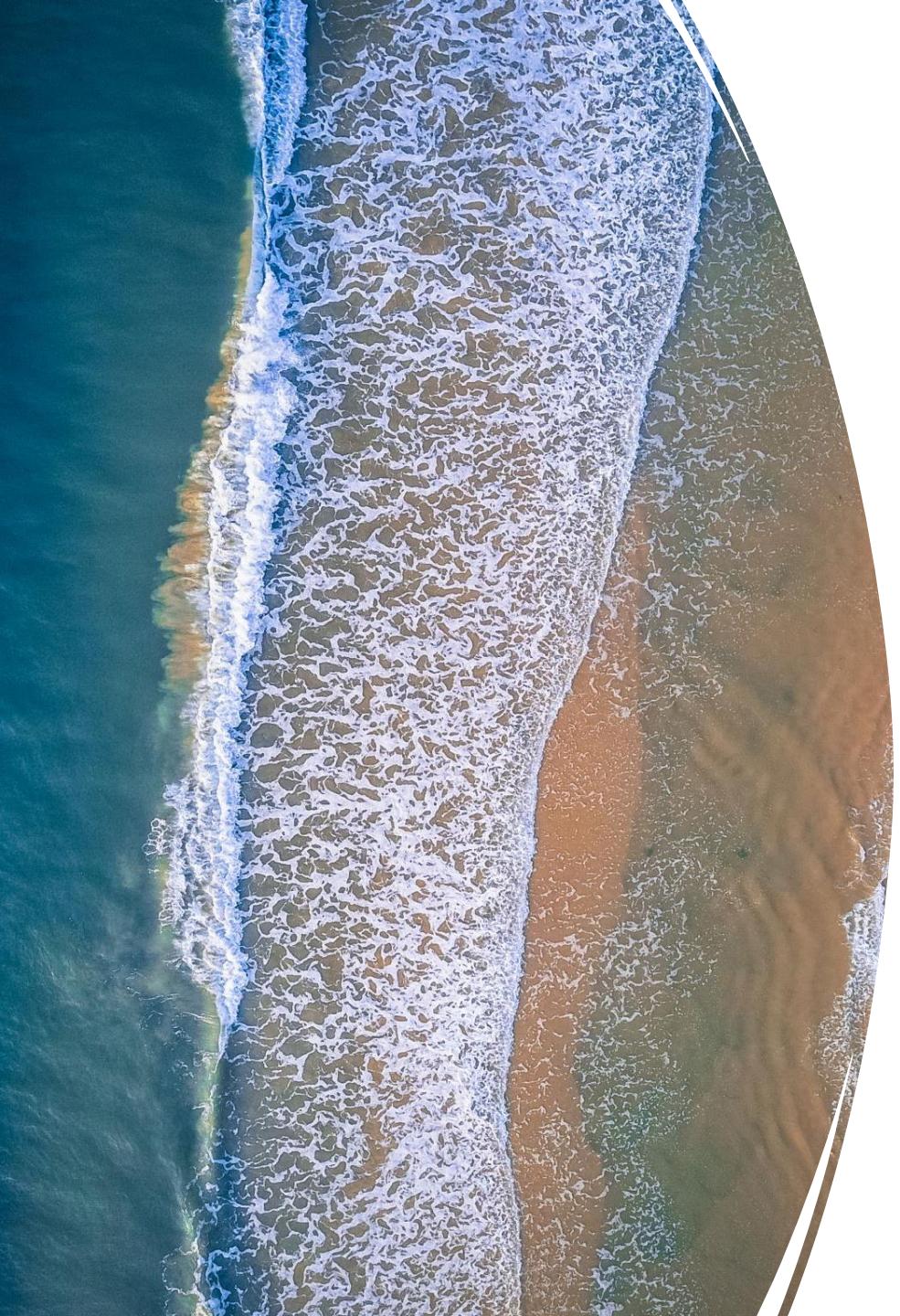
in lbs of CO<sub>2</sub> equivalent

Roundtrip flight b/w NY and SF (1 passenger)	1,984
Human life (avg. 1 year)	11,023
American life (avg. 1 year)	36,156
US car including fuel (avg. 1 lifetime)	126,000
Transformer (213M parameters) w/ neural architecture search	626,155

# MICROSOFT IS DRAINING AN ARIZONA TOWN'S WATER SUPPLY FOR ITS AI

## COULD WE JUST... NOT DO THAT?

Estimates commissioned by Microsoft itself, the 279-acre campus would approx. consume an annual **56 million gallons** of drinking water ( approximately the amount that a total of **670 Goodyear (Arizona) families** would consume in a year combined)

A circular inset image in the top left corner shows an aerial view of a coastline. The ocean is a deep blue, and white waves are crashing onto a light-colored sandy beach. The water has a textured, foamy appearance where it meets the shore.

# Sustainable AI: Beyond Carbon Footprint

- Building and operating a leading-edge semiconductor fabrication plant, or “fab”, to produce chips like in 5nm can require up to **four million gallons of pure water each day**.
- This water usage approaches what a city of **half a million people** would require for all needs.
- Sourcing this consistently places immense strain on local water tables and reservoirs, especially in already water-stressed regions which host many high-tech manufacturing hubs.



Who is she?

[https://www.reddit.com/r/ChatGPT/comments/1hu7i57/we\\_are\\_doomed/?rdt=40200](https://www.reddit.com/r/ChatGPT/comments/1hu7i57/we_are_doomed/?rdt=40200)