

Large Language Models Operations (LLMOps)

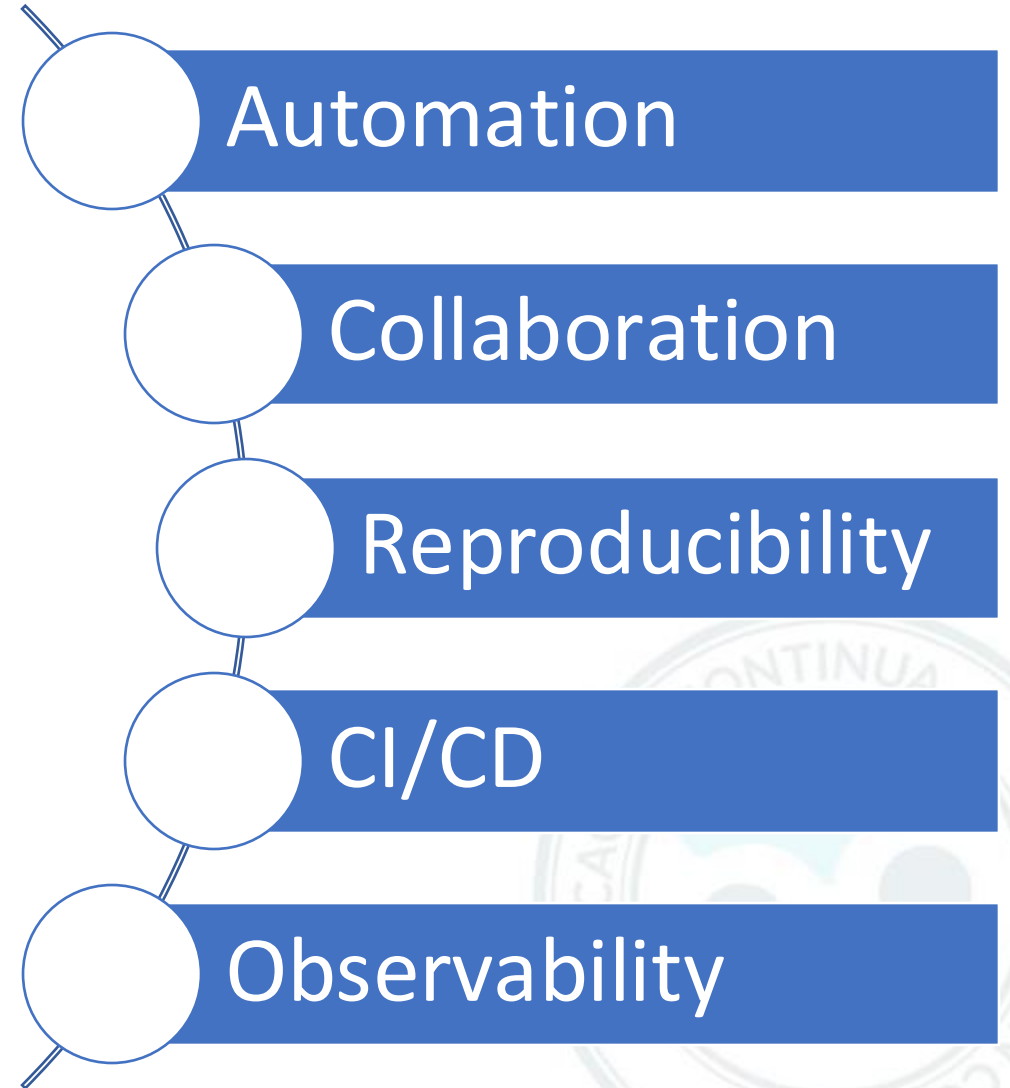
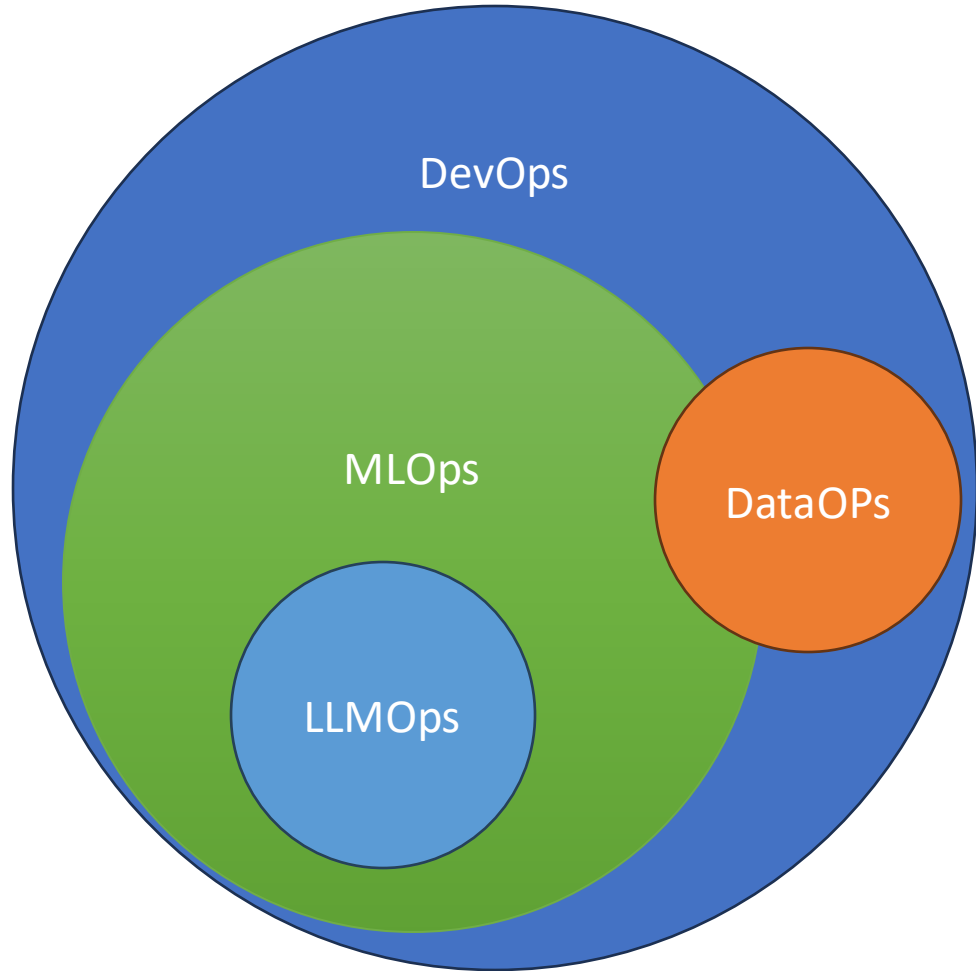
- LLMOPs introduction
- Deployment and scalability of LLMs
- Monitoring and maintenance of models in production
- Performance evaluation and continuous improvement
- Ethical considerations and privacy



Intro to LLMOPS

XOP

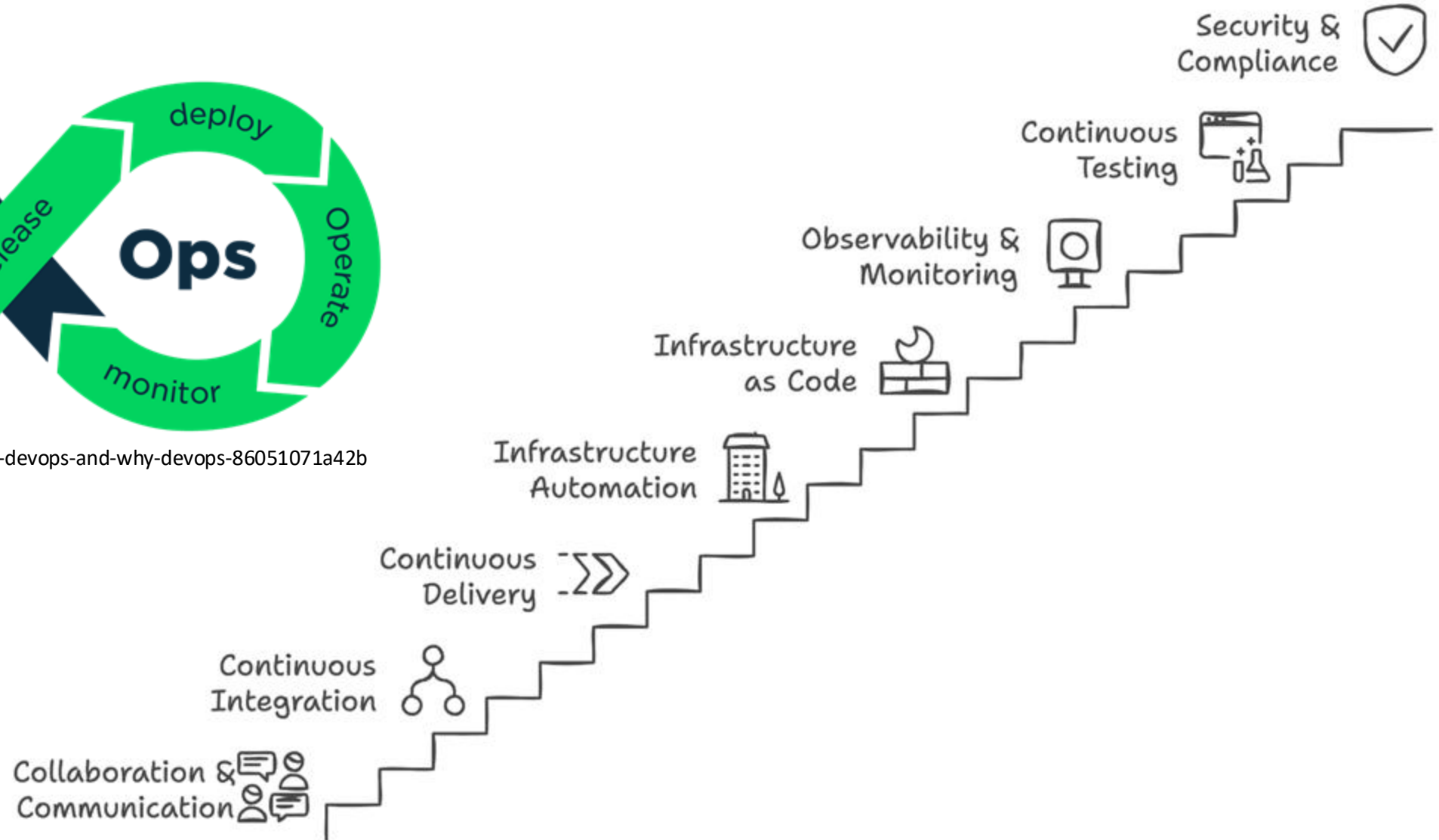
S



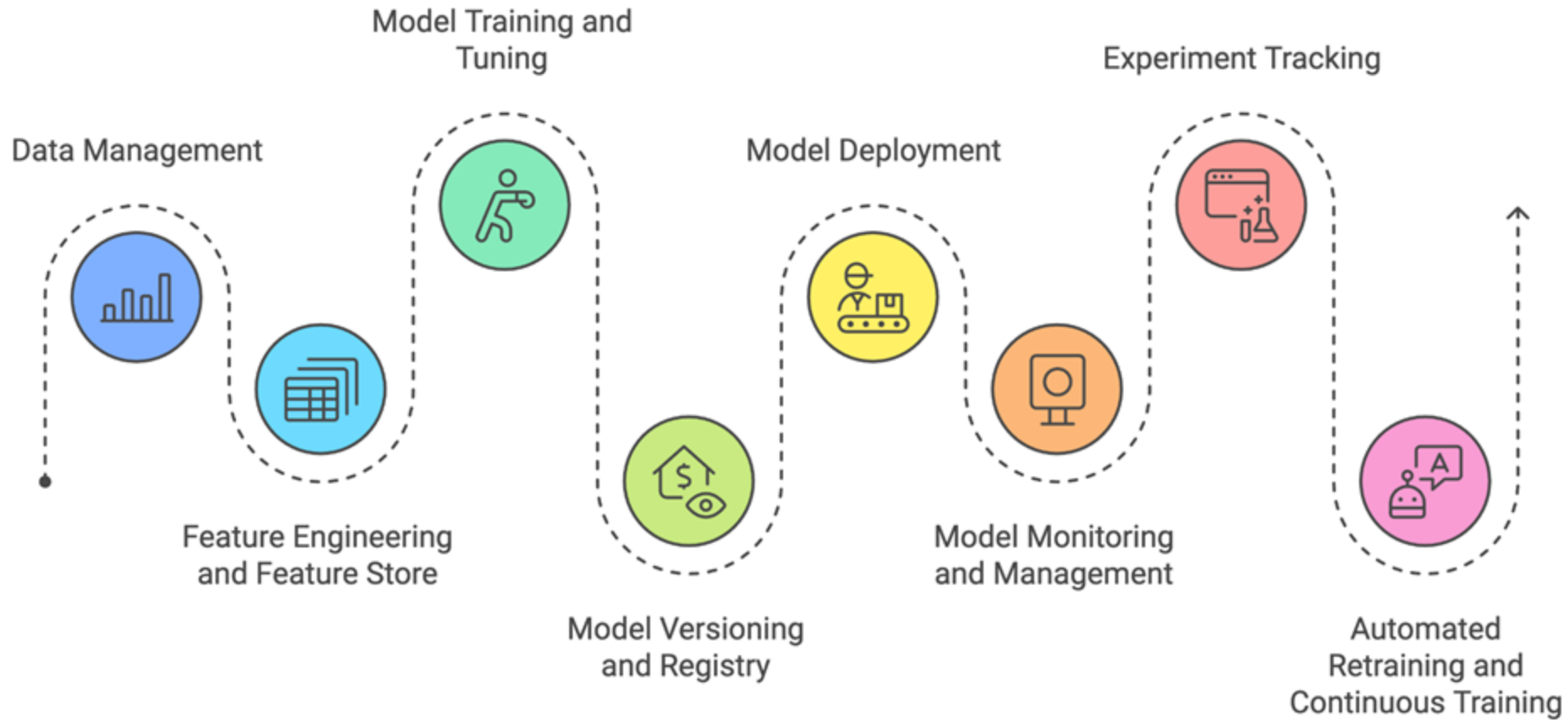
DevOps



<https://medium.com/@rituserke86/what-is-devops-and-why-devops-86051071a42b>



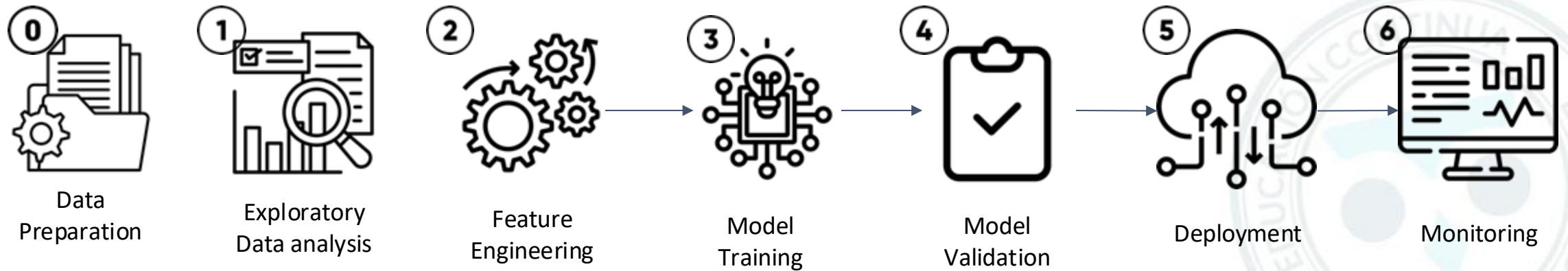
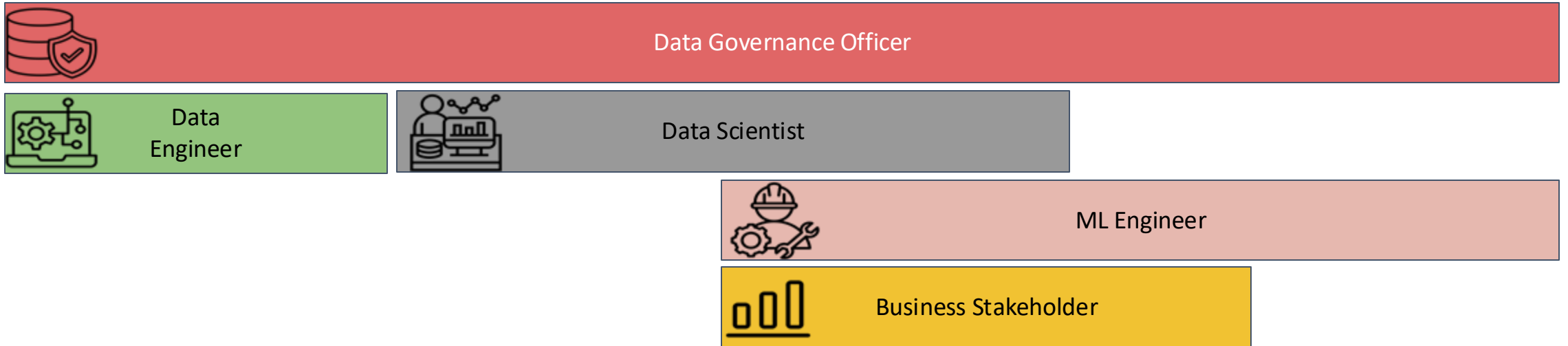
MLOps



LLMOps



Training and deployment of Large Language Models



What specific scaling challenges exist for LLMs?



Initial training: trillions of tokens, hundreds to thousands of GBs & very long run times.



Fine-tuning: updating model weights based on your own data, still requires relatively large data and long training times. Plus lots of evaluation!



Storage: the models contain billions of parameters, often hundreds of GB. This can be prohibitive for some devices.



Latency: Running inference on these models can be very costly in terms of time.



Cost: All of this costs \$\$\$!



What specific scaling challenges exist for LLMs?



Initial training: don't do it!



Fine-tuning: Optimize and use a scalable framework, such as Ray.



Storage: Quantization, memorization, caching ...



Latency: Quantization, memorization, caching, hardware and memory bandwidth optimization...



Cost: Above plus use 'open source' models



Types of use-cases for Enterprises

How willing are enterprises to use LLMs for different use cases?



(% of enterprises experimenting with given use case who have deployed to production)



Source: a16z survey of 70 enterprise AI decision makers

LLM application archetypes

Prompt
Engineering

Retrieval
Augmented
Generation

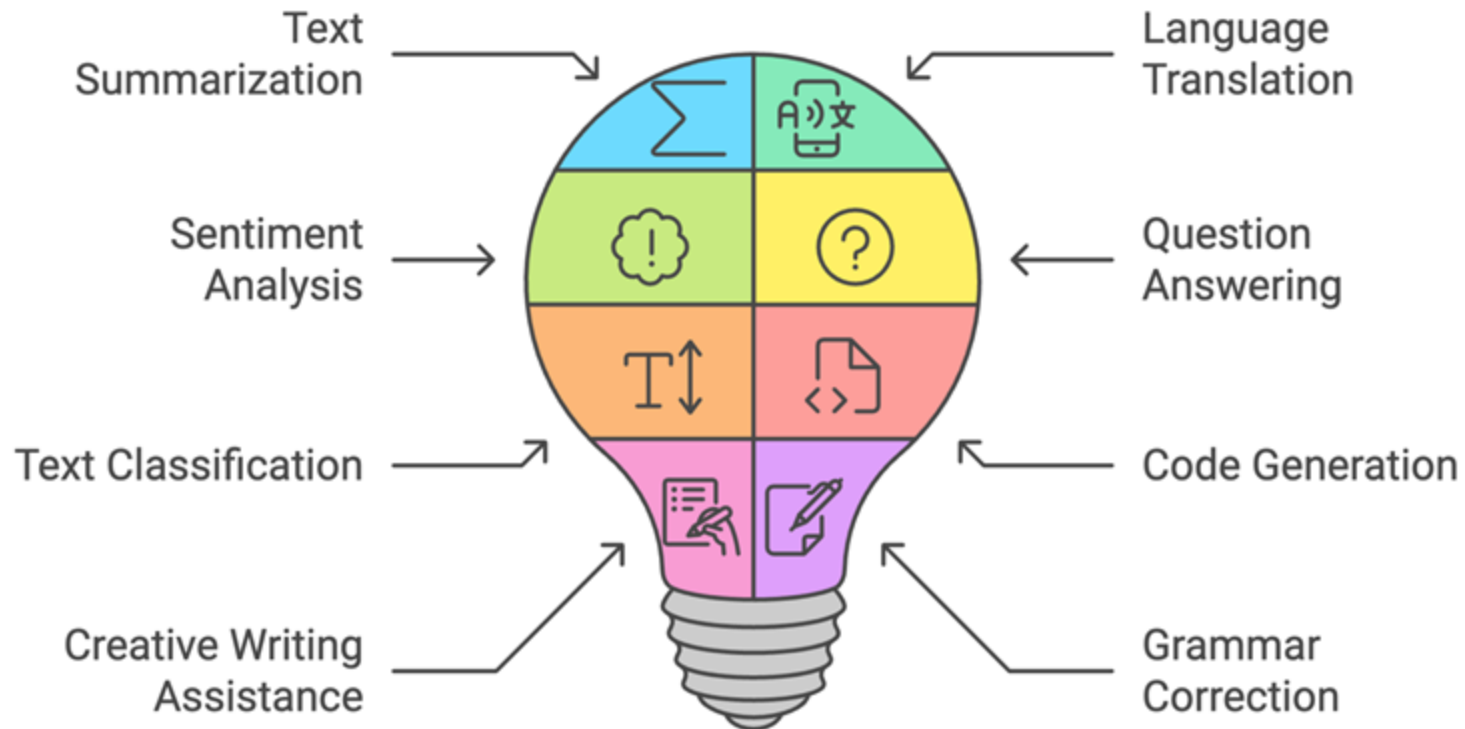
Agents

Multi-Agents

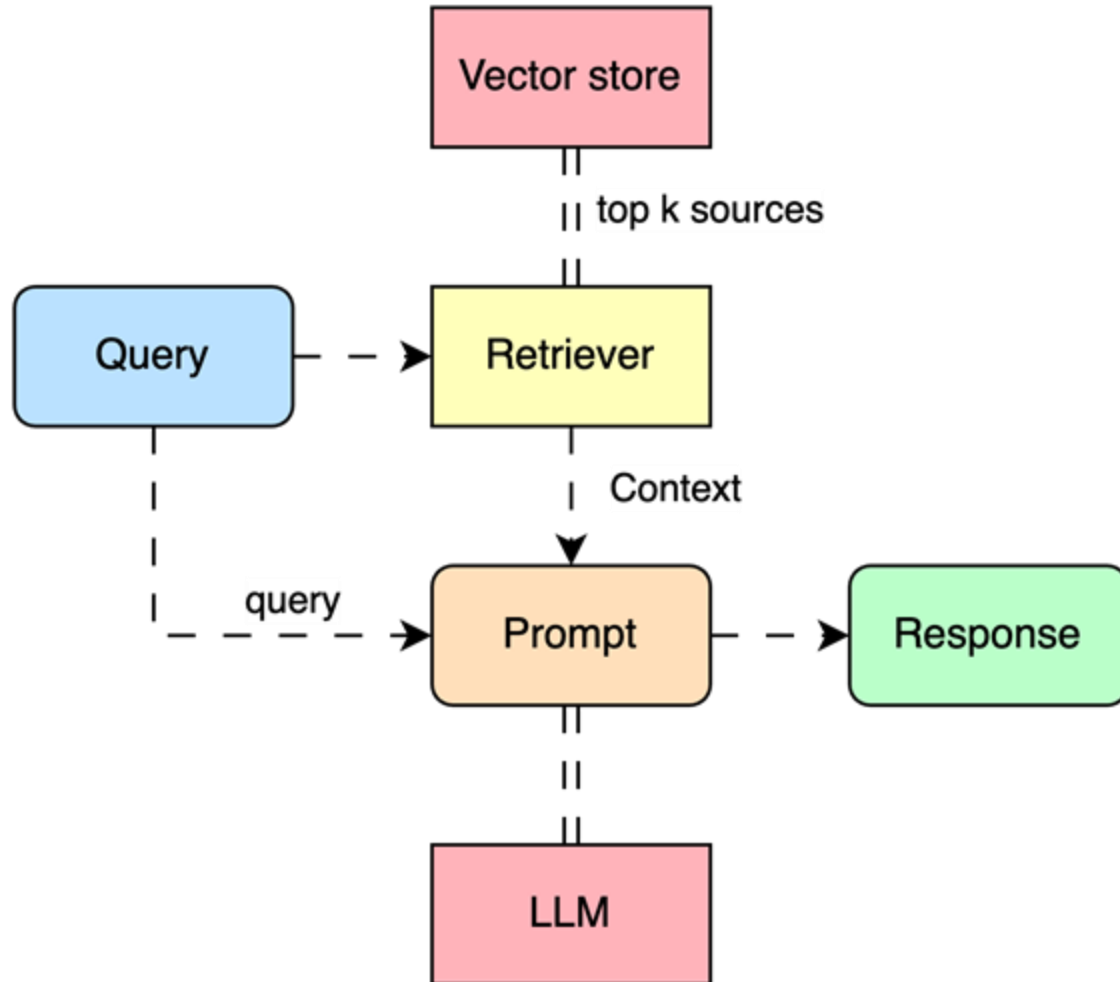
Fine-tuning



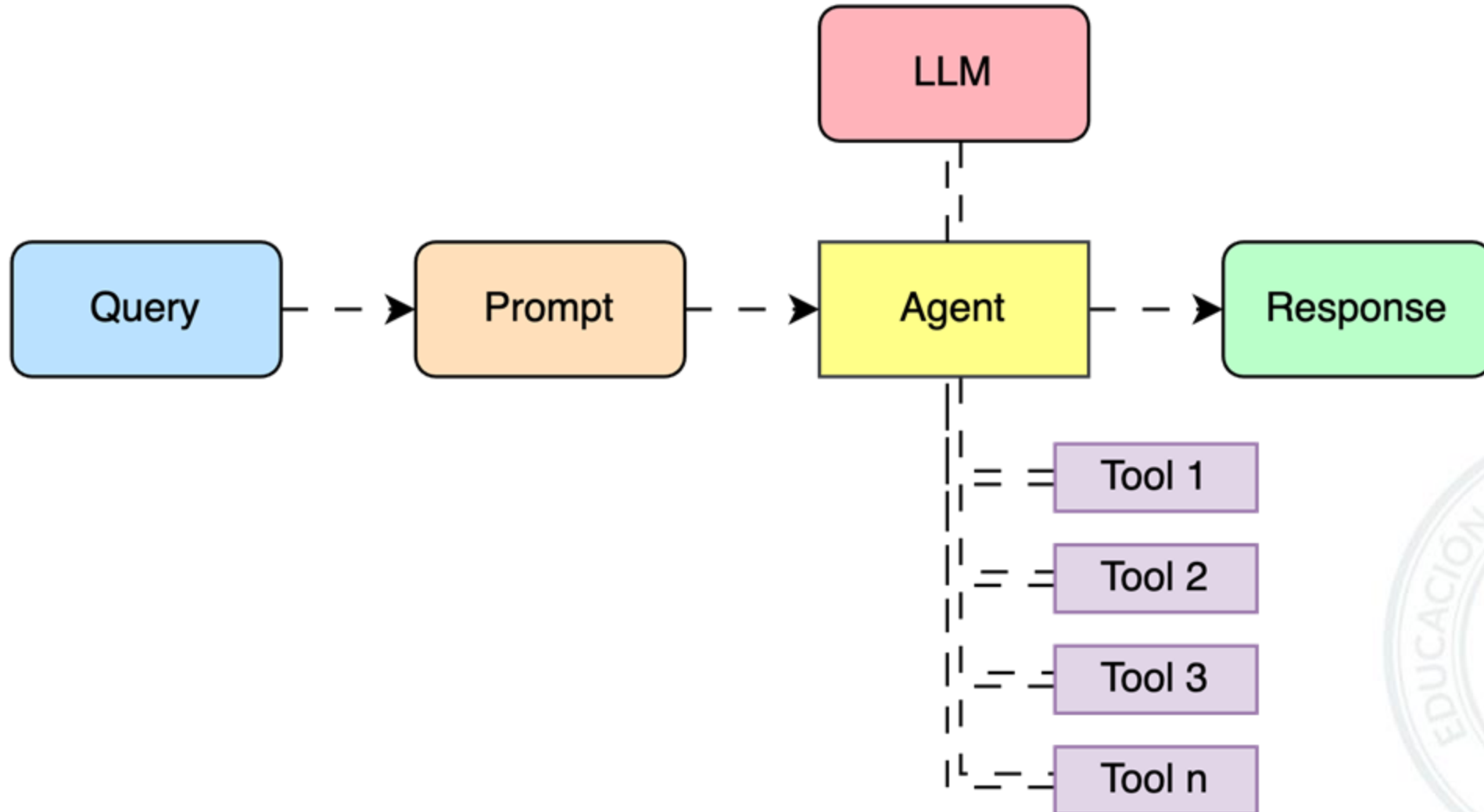
Prompt Engineering Applications



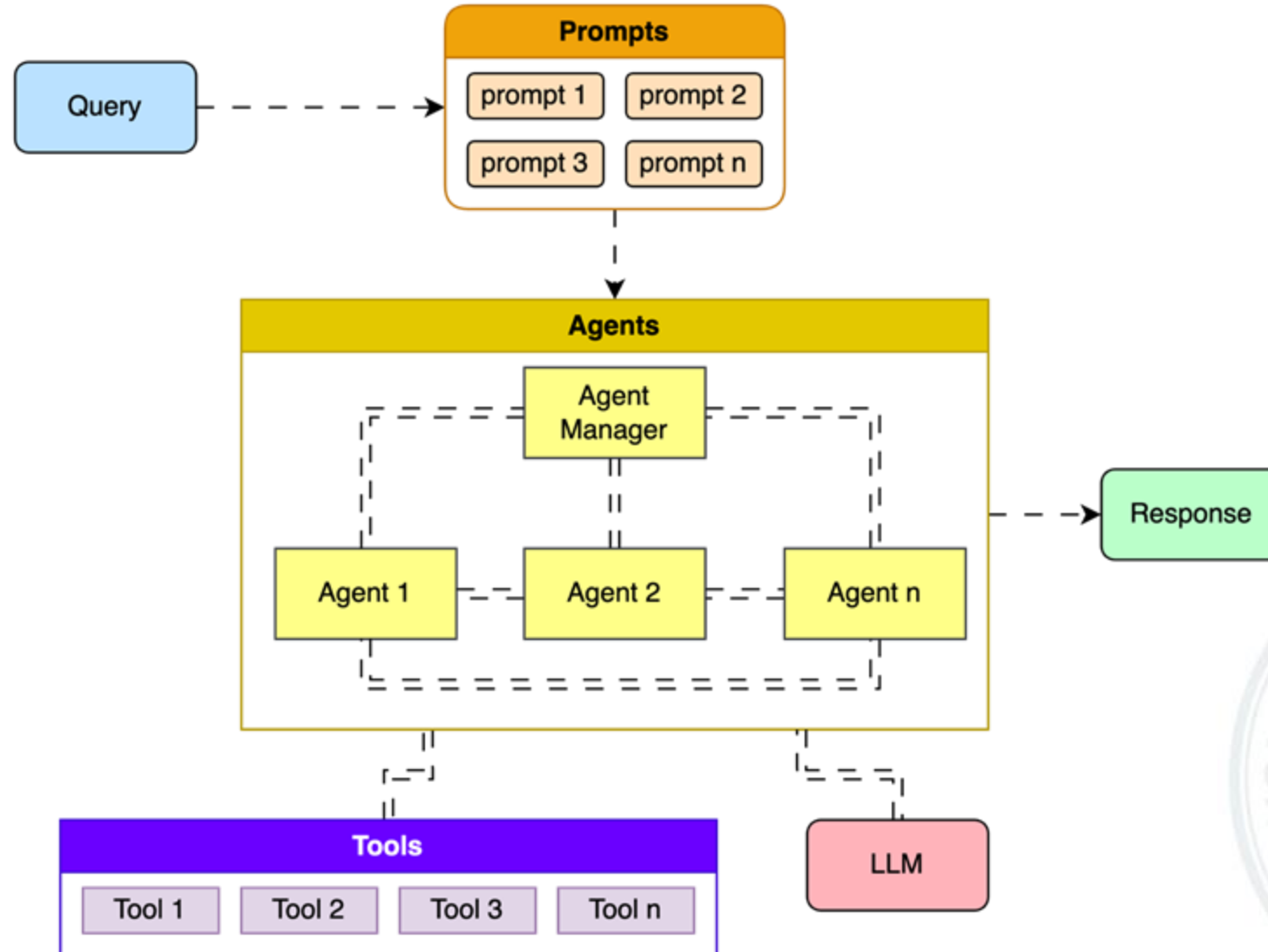
Retrieval Augmented Generation Applications



Agentic Applications

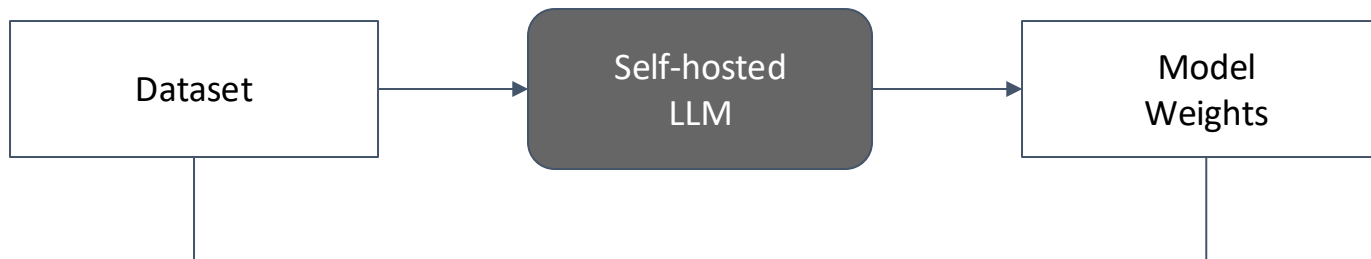


Multi-Agent Applications



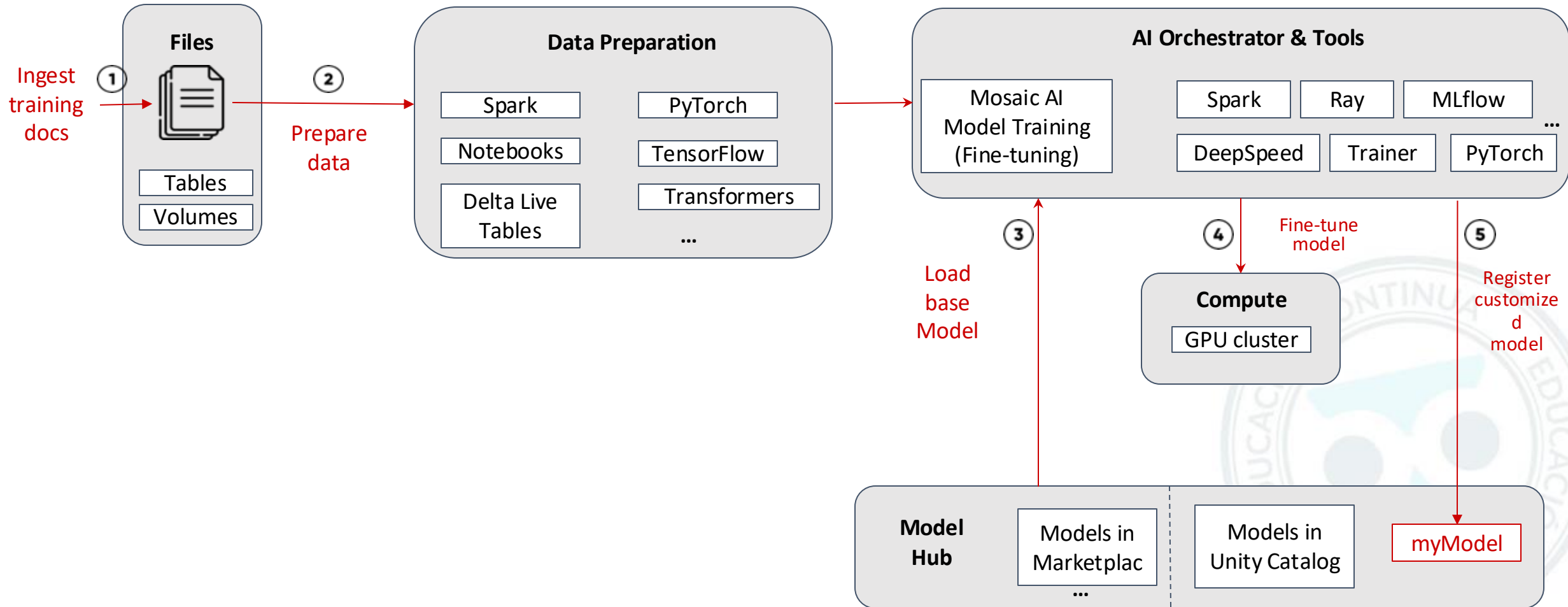
Operationalizing Fine-tuning Pipelines

- Automation, version control, reproducibility
- Distributed training infrastructure
 - DeepSpeed, PEFT, GPUs
- Similar to classical model training serving
- Optimization techniques
 - vLLM, MLC, CudaGraph, MQA, Quantization, TensorRT
- Deeper understanding of GPUs, TTFT, TPOT



Types of LLM Applications

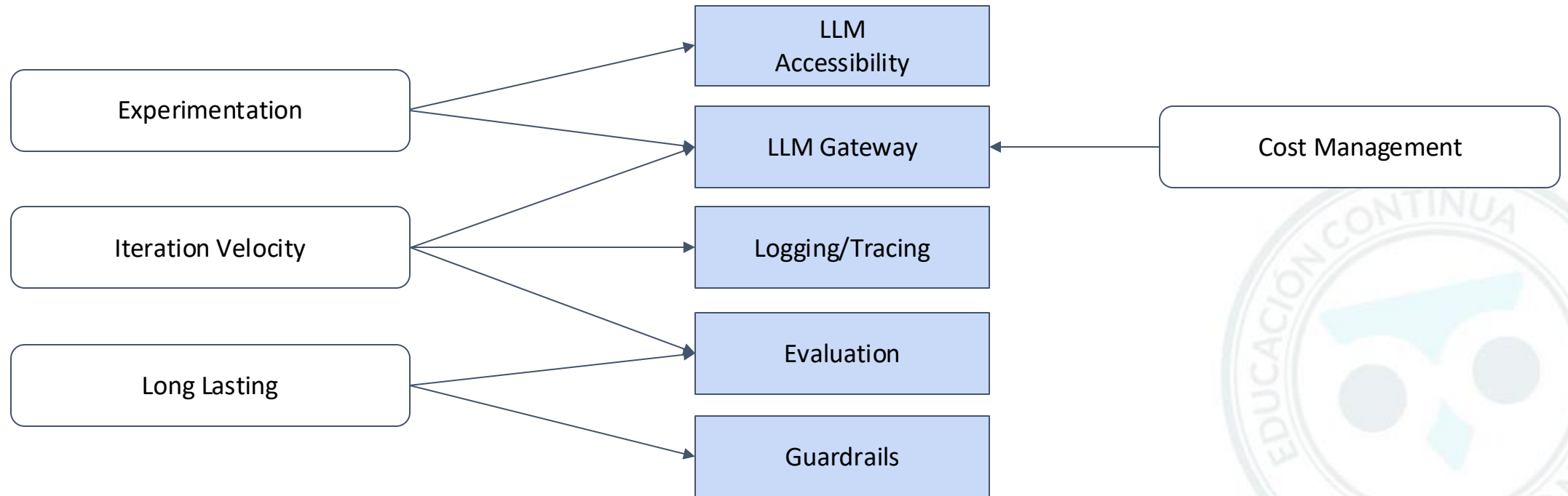
Architecture: fine-tuning

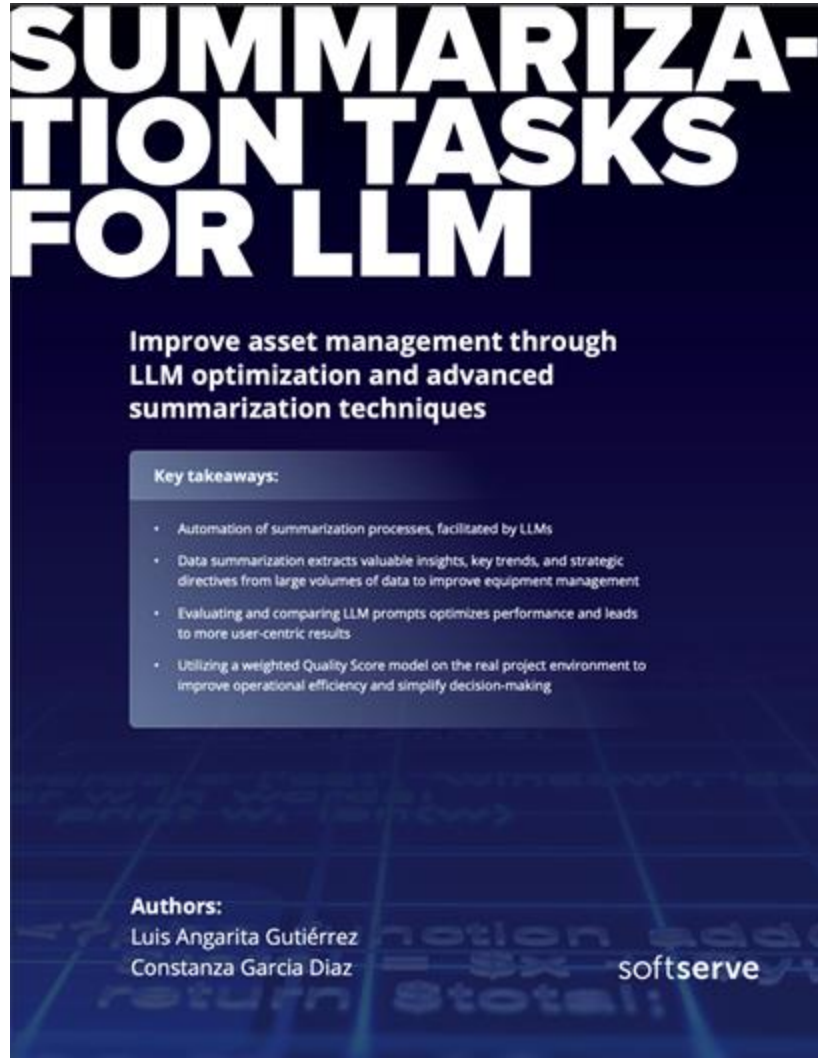


Types of LLM Applications

Practitioner's Perspective

LLMOps Starting Points





Improve Asset Management with LLM Evaluation Pipelines

Creating efficient AI-powered summaries, LLMs save up to 93% of the time taken by traditional manual processes. Our white paper explains how this allows you to optimize asset management operations.

<https://info.softserveinc.com/summarization-tasks-for-llm-white-paper>



Prompt engineering Example: Youtube Summarizer

GitHub Repo: <https://github.com/LGuillermoAngaritaG/llmops-youtube-summarizer>

Groq-API-Key: <https://console.groq.com/keys>