



Java Me Smarter

Unleashing AI Power with Quarkus



China 2024

Daniel Oh



X @danieloh30

► @danieloh30

Github danieloh30

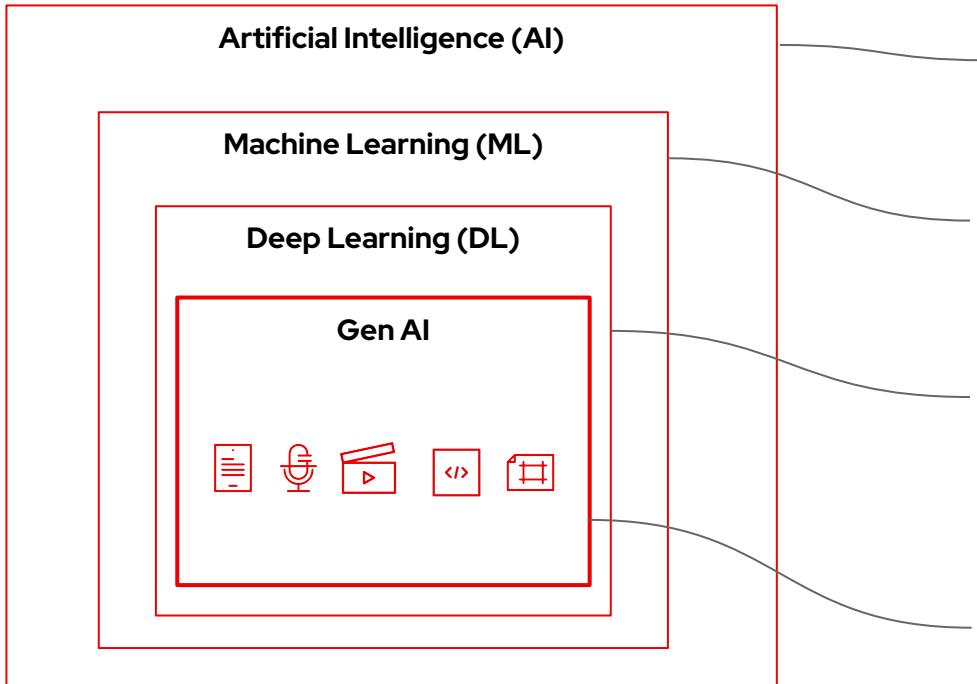


- Developer Advocate at Red Hat
 - Cloud Native Runtimes
 - AI, GitOps, and Serverless
- JAVA Champion
- Cloud Native Computing Foundation Ambassador
- App Development WG Co-Chair in CNCF TAG App Delivery
- Public Speaker & Published Author
- Living in Boston, MA USA



The Journey to Generative AI and LLMs

What really is Generative AI?



Artificial Intelligence (AI)

A multidisciplinary field of Computer Science that aims to create systems capable of emulating and surpassing human-level intelligence

Machine Learning (ML)

A subfield for training computers to learn and make decisions from data, without explicit programming

Deep Learning (DL)

A subset from ML involved neural networks enabling complex pattern recognition and analysis from data

Gen AI Learning

Subfield focusing on generating new content such as text, images, code and more

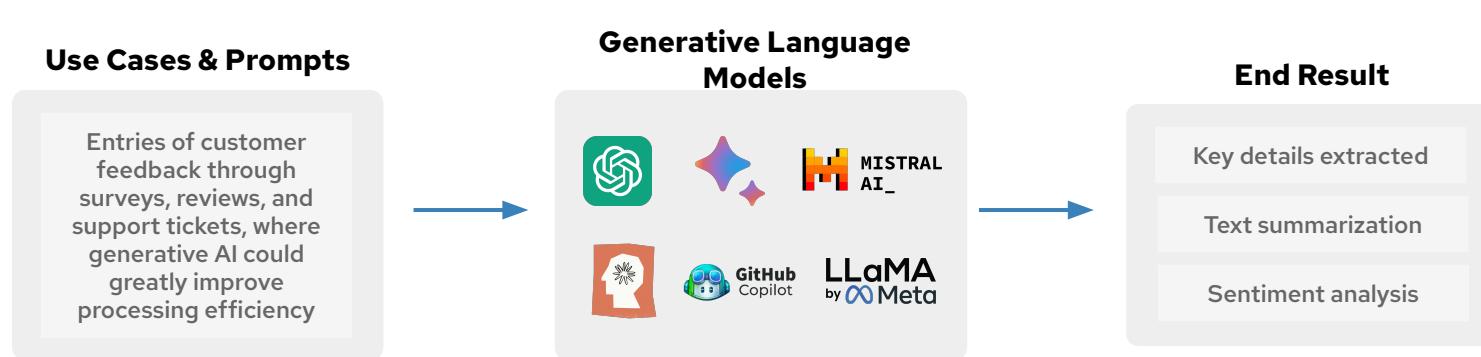
How does Generative AI really work?

Generative AI

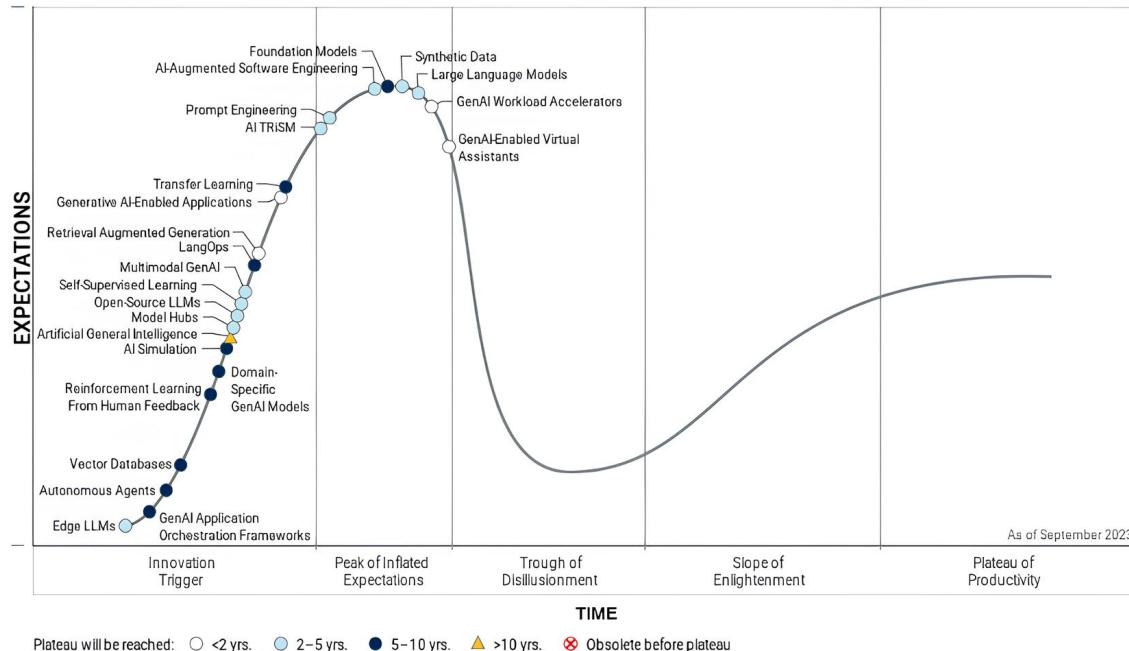
Generative AI utilizes AI models to create text, images, videos, code, or data in response to prompts

LLMs

Large Language Models are a specialized type of AI, trained on massive datasets to generate outputs in natural, conversational language



The AI adoption journey with the rise of Gen AI



Gartner

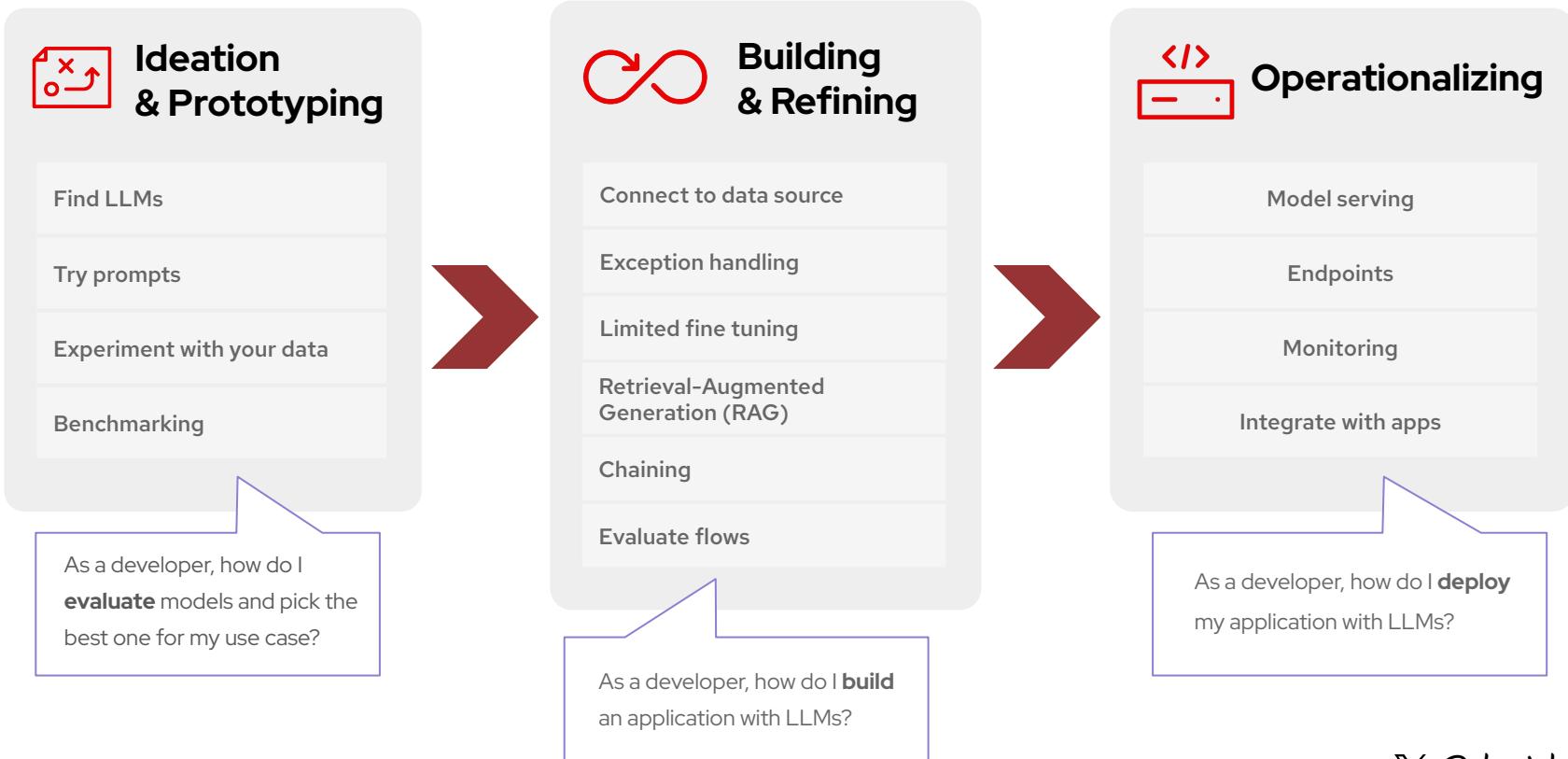
" More Than 80% of Enterprises Will Have Used Generative AI APIs or Deployed Generative AI-Enabled Applications by 2026. "

Gartner

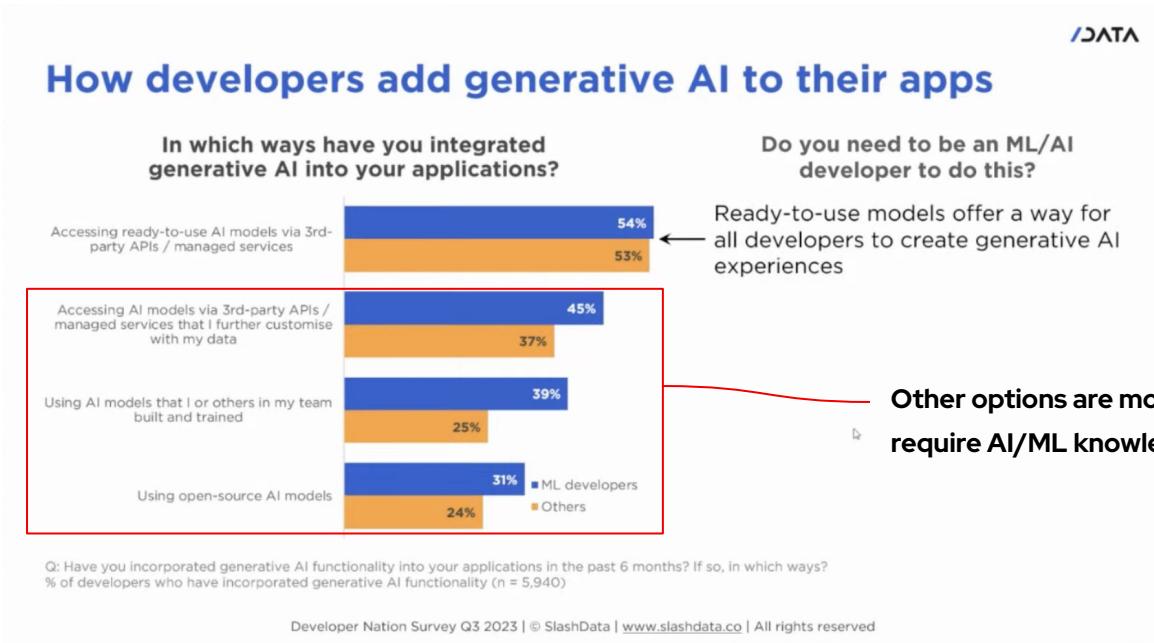
<https://www.gartner.com/en/newsroom/press-releases/2023-10-11-gartner-says-more-than-80-percent-of-enterprises-will-have-used-generative-ai-apis-or-deployed-generative-ai-enabled-applications-by-2026>

The Gen AI Journey for an Application Developer

The Journey of Adopting **Generative AI**

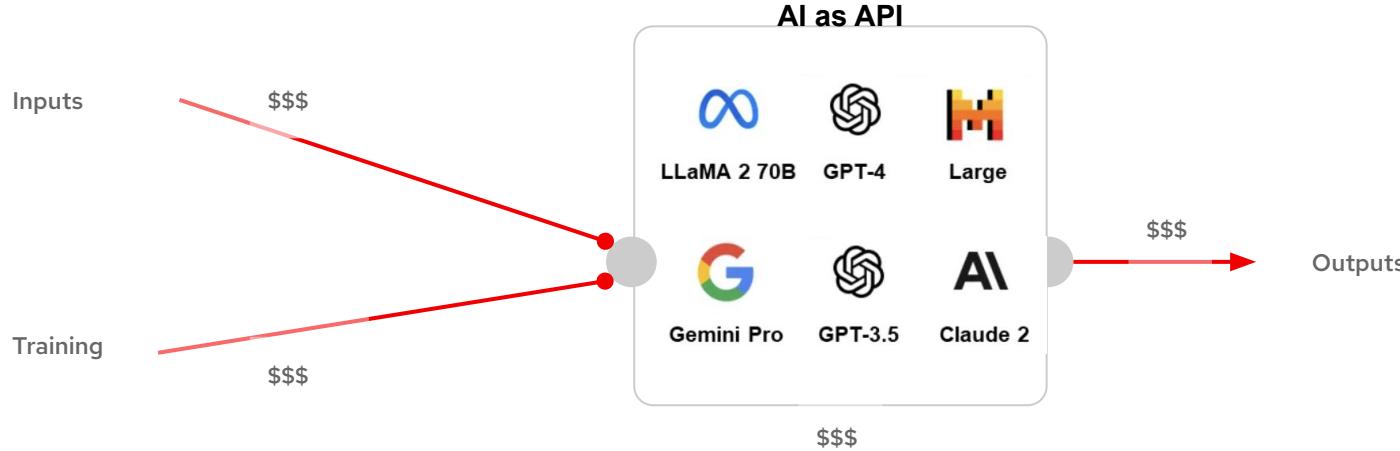


Developer's **don't** have to be AI/ML experts



AI/ML expertise is **not** a significant factor
for integrating ready-to-use models

Those APIs are costly... and challenging to test against



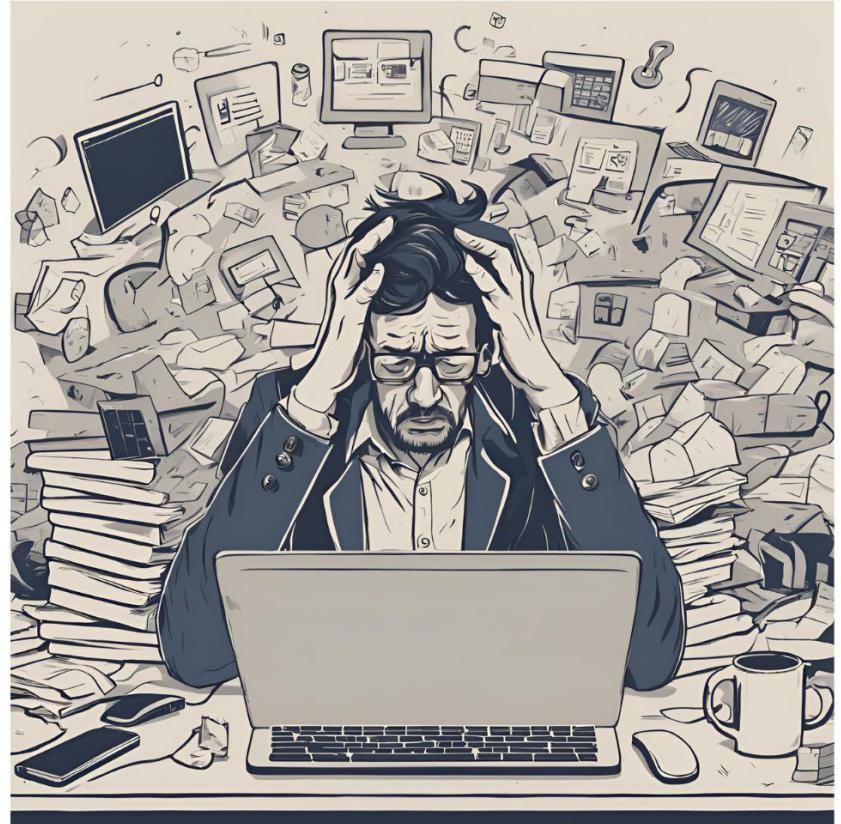
of tokens used and costs randomly exploded over night

OpenAI pricing				
GPT-4		Chat		
Model	Prompt	Completion	Input	Output
8K context	\$0.03 / 1K tokens	\$0.06 / 1K tokens		
32K context	\$0.06 / 1K tokens	\$0.12 / 1K tokens		
InstructGPT				
Model	Usage		Price	
Ada	\$0.0004 / 1K tokens		\$0.020 / image	
Babbage	\$0.0005 / 1K tokens		\$0.018 / image	
Curie	\$0.0020 / 1K tokens		\$0.016 / image	
Davinci	\$0.0200 / 1K tokens			
Image models				
Model	Resolution		Price	
	1024x1024		\$0.020 / image	
	512x512		\$0.018 / image	
	256x256		\$0.016 / image	
Embedding models				
Model	Usage		Price	
Ada v2	\$0.0001 / 1K tokens		\$0.0004 / 1K tokens	
Ada v1	\$0.0004 / 1K tokens		\$0.0010 / 1K tokens	
Babbage v1	\$0.0005 / 1K tokens		\$0.0008 / 1K tokens	
Curie v1	\$0.0020 / 1K tokens		\$0.0120 / 1K tokens	
Davinci v1	\$0.0200 / 1K tokens		\$0.1200 / 1K tokens	
Fine-tuning models				
Model	Training	Usage		
Ada	\$0.0004 / 1K tokens	\$0.0010 / 1K tokens		
Babbage	\$0.0006 / 1K tokens	\$0.0024 / 1K tokens		
Curie	\$0.0030 / 1K tokens	\$0.0120 / 1K tokens		
Davinci	\$0.0030 / 1K tokens	\$0.1200 / 1K tokens		
Audio models				
Model	Usage			
Whisper	\$0.006 / minute (rounded to the nearest second)			

Cost for GPT failed requests:

- Issue from OpenAI side
- Timeout in Application

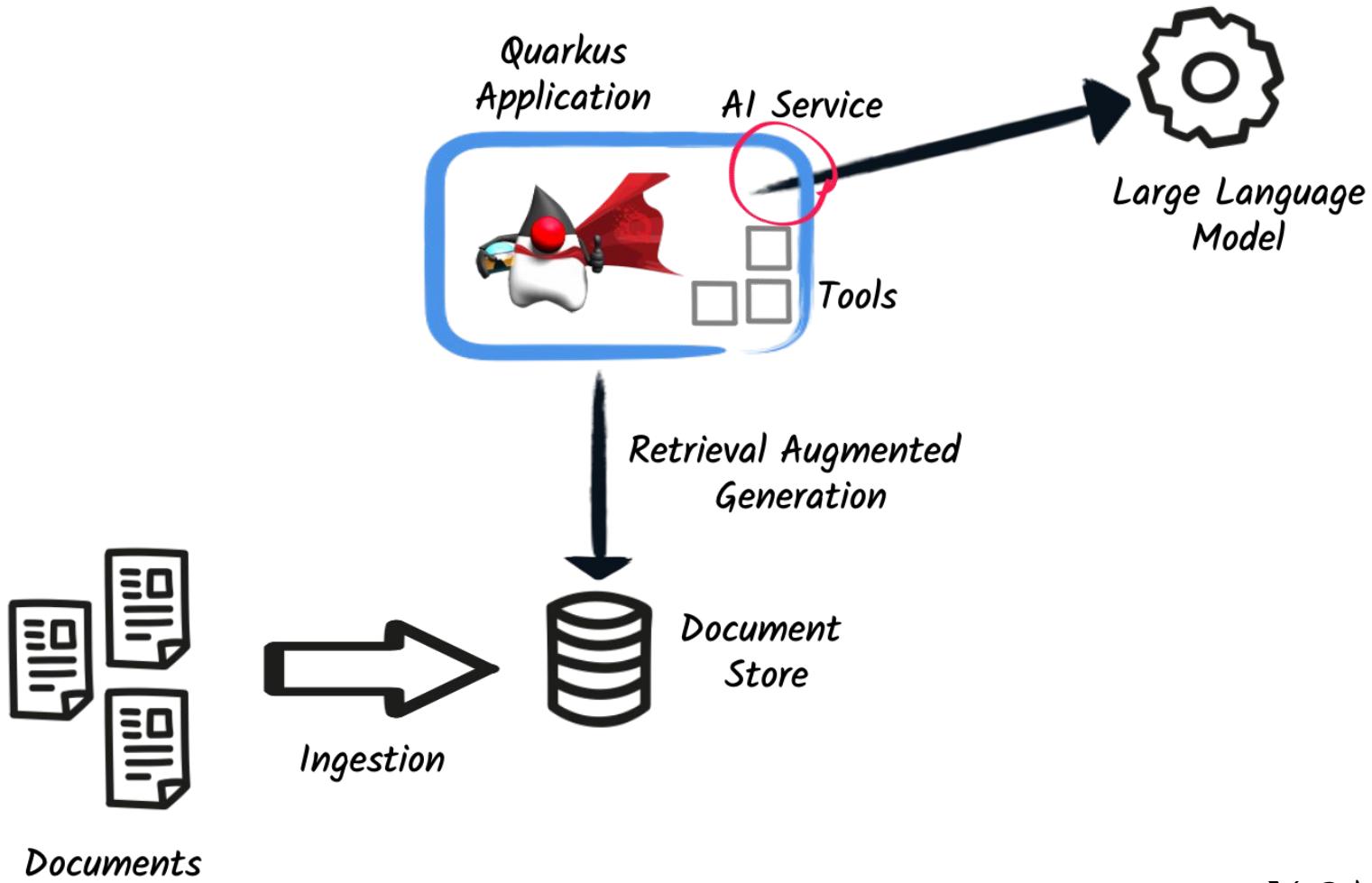
Still Feeling Stuck in a **Java** Rut?



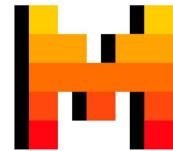
Langchain4j



QUARKUS



```
<dependency>
  <groupId>io.quarkiverse.langchain4j</groupId>
  <artifactId>quarkus-langchain4j-openai</artifactId>
  <version>0.17.0</version>
</dependency>
```



```
package io.quarkiverse.langchain4j.samples;

import dev.langchain4j.service.SystemMessage;
import dev.langchain4j.service.UserMessage;
import io.quarkiverse.langchain4j.RegisterAiService;

@registerAIService(
    tools = EmailService.class
)
public interface MyAiService {

    @SystemMessage("You are a professional poet")
    @UserMessage("""
        Write a poem about {topic}. The poem should be {lines}
lines long. Then send this poem by email.
    """)
    String writeAPoem(String topic, int lines);
}
```



DEMO



Agentic AI applications

Prompts

- Interacting with the model for asking questions
- Interpreting messages to get important information
- Populating Java classes from natural language
- Structuring output

Agent and Tools

- Mixing business code with model
- Delegating to external services
- Describing tools in a manner most beneficial to the agent

Embedding Documents (RAG)

- Adding specific knowledge to the model
- Asking questions about supplied documents
- Natural queries

Document Store

- Redis Store
- Chroma Store
- Infinispan Store
- Pinecone Store
- PgVector (PostgreSQL) Store
- In-Process Embeddings
- Loading CSV files
- Neo4j Store

Observability

- Metrics
- Tracing
- Auditing



Advanced Topics

- Fault Tolerance
- WebSockets
- Enabling and Disabling Integrations

youtube.com/@danieloh30

bit.ly/danielohtv

Subscribe



Gen AI and MLOps ► Play all

- [Streamline Inner Loop for AI App Development with... #10:48](#)
Daniel Oh
47 views • 12 days ago
- [Supercharging Static Code Analysis: Konveyor AI & LLMs #13:48](#)
Daniel Oh
86 views • 2 weeks ago
- [Build truly open source LLMs with InstructLab #15:57](#)
Daniel Oh
127 views • 3 weeks ago
- [Java Me Smarter Quarkus + AI? #24:02](#)
Daniel Oh
206 views • 4 months ago

KUBERNETES LEARN BY EXAMPLE ► Play all

- [Kubernetes Learn by Example #11 #5:37](#)
Daniel Oh
282 views • 11 months ago
- [KUBERNETES LEARN BY EXAMPLE #9 Persistent Volumes #10:31](#)
Daniel Oh
394 views • 2 years ago
- [KUBERNETES LEARN BY EXAMPLE #8 ConfigMaps #6:45](#)
Daniel Oh
285 views • 2 years ago
- [KUBERNETES LEARN BY EXAMPLE #7 Managing Secrets #7:12](#)
Daniel Oh
209 views • 2 years ago
- [KUBERNETES LEARN BY EXAMPLE #6 DaemonSet #7:00](#)
Daniel Oh
202 views • 2 years ago
- [KUBERNETES LEARN BY EXAMPLE #5 StatefulSets #8:54](#)
Daniel Oh
595 views • 2 years ago

QUARKUS ► Play all

- [Quarkus + Testcontainers: Local or Remote? #10:10](#)
- [Java Me Smarter Quarkus + AI? #24:02](#)
- [Quarkus Dev Services for Kubernetes #13:30](#)
- [Quarkus Dev Services Inside a Container #10:25](#)
- [Trigger a Remote Serverless Java Deployment - Part I #11:10](#)
- [Jumpstart New Cloud Dev Environment with Dev Spaces #5:04](#)



THANK YOU! QUESTIONS?



China 2024