

AI Java Labs: Building Intelligent Java Applications

From HTTP fundamentals to advanced RAG systems →



Contact Info

Ken Kousen

Kousen IT, Inc.

- ken.kousen@kousenit.com
- <http://www.kousenit.com>
- <http://kousenit.org> (blog)
- Social Media:
 - [@kenkousen](#) (Twitter)
 - [@kousenit.com](#) (Bluesky)
 - <https://www.linkedin.com/in/kenkousen/> (LinkedIn)
- *Tales from the jar side* (free newsletter)
 - <https://kenkousen.substack.com>
 - <https://youtube.com/@talesfromthejarside>

What You'll Learn



What You'll Learn

- **HTTP Fundamentals:** Raw API integration patterns



What You'll Learn

- **HTTP Fundamentals:** Raw API integration patterns
- **LangChain4j Framework:** High-level AI abstractions



What You'll Learn

- **HTTP Fundamentals:** Raw API integration patterns
- **LangChain4j Framework:** High-level AI abstractions
- **Text Generation:** OpenAI and Ollama models



What You'll Learn

- **HTTP Fundamentals:** Raw API integration patterns
- **LangChain4j Framework:** High-level AI abstractions
- **Text Generation:** OpenAI and Ollama models
- **Streaming Responses:** Real-time token-by-token output



What You'll Learn

- **HTTP Fundamentals:** Raw API integration patterns
- **LangChain4j Framework:** High-level AI abstractions
- **Text Generation:** OpenAI and Ollama models
- **Streaming Responses:** Real-time token-by-token output
- **Multimodal AI:** Vision, audio, and image generation



What You'll Learn

- **HTTP Fundamentals:** Raw API integration patterns
- **LangChain4j Framework:** High-level AI abstractions
- **Text Generation:** OpenAI and Ollama models
- **Streaming Responses:** Real-time token-by-token output
- **Multimodal AI:** Vision, audio, and image generation
- **RAG Systems:** Document-based question answering



What You'll Learn

- **HTTP Fundamentals:** Raw API integration patterns
- **LangChain4j Framework:** High-level AI abstractions
- **Text Generation:** OpenAI and Ollama models
- **Streaming Responses:** Real-time token-by-token output
- **Multimodal AI:** Vision, audio, and image generation
- **RAG Systems:** Document-based question answering
- **Cost Management:** Efficient testing and model selection



What You'll Learn

- **HTTP Fundamentals:** Raw API integration patterns
- **LangChain4j Framework:** High-level AI abstractions
- **Text Generation:** OpenAI and Ollama models
- **Streaming Responses:** Real-time token-by-token output
- **Multimodal AI:** Vision, audio, and image generation
- **RAG Systems:** Document-based question answering
- **Cost Management:** Efficient testing and model selection
- **Modern Java:** Records, sealed interfaces, pattern matching



Repository Structure

```
1  AiJavaLabs/
2  |— labs.md           # 15 progressive lab exercises
3  |— src/
4  |   |— main/java/com/kousenit/
5  |   |   |— demos/
6  |   |   |   |— QuickChatDemo.java    # Fast OpenAI demo
7  |   |   |   |— TextToSpeechDemo.java # TTS generation
8  |   |   |   |— LocalOllamaDemo.java  # Local AI models
9  |   |   |   |— MultiModelDemo.java   # Provider comparison
10 |   |   |   |— ResponsesApiDemo.java  # Raw HTTP examples
11 |   |   |   |— ApiComparisonDemo.java # Framework vs raw
12 |   |   |   |— StreamingDemo.java     # Real-time responses
13 |   |   |— EasyRAGDemo.java           # Document Q&A
14 |   |   |— *Service.java              # Core implementations
15 |   |   |— *Records.java              # Data models
16 |   |— test/java/                     # Comprehensive test suite
17 |— build.gradle.kts                   # Test categories for cost control
18 |— slides.md                         # This presentation
```

Repository Structure

```
1  AiJavaLabs/
2  |— labs.md           # 15 progressive lab exercises
3  |— src/
4  |   |— main/java/com/kousenit/
5  |   |   |— demos/
6  |   |   |   |— QuickChatDemo.java    # Fast OpenAI demo
7  |   |   |   |— TextToSpeechDemo.java  # TTS generation
8  |   |   |   |— LocalOllamaDemo.java   # Local AI models
9  |   |   |   |— MultiModelDemo.java    # Provider comparison
10 |   |   |   |— ResponsesApiDemo.java   # Raw HTTP examples
11 |   |   |   |— ApiComparisonDemo.java  # Framework vs raw
12 |   |   |   |— StreamingDemo.java      # Real-time responses
13 |   |   |— EasyRAGDemo.java            # Document Q&A
14 |   |   |— *Service.java               # Core implementations
15 |   |   |— *Records.java               # Data models
16 |   |— test/java/                      # Comprehensive test suite
17 |— build.gradle.kts                   # Test categories for cost control
18 |— slides.md                         # This presentation
```

- **8 Live Demos:** Ready-to-run examples for each topic

Repository Structure

```
1  AiJavaLabs/
2  |— labs.md           # 15 progressive lab exercises
3  |— src/
4  |   |— main/java/com/kousenit/
5  |   |   |— demos/
6  |   |   |   |— QuickChatDemo.java    # Fast OpenAI demo
7  |   |   |   |— TextToSpeechDemo.java # TTS generation
8  |   |   |   |— LocalOllamaDemo.java  # Local AI models
9  |   |   |   |— MultiModelDemo.java   # Provider comparison
10 |   |   |   |— ResponsesApiDemo.java  # Raw HTTP examples
11 |   |   |   |— ApiComparisonDemo.java # Framework vs raw
12 |   |   |   |— StreamingDemo.java     # Real-time responses
13 |   |   |— EasyRAGDemo.java           # Document Q&A
14 |   |   |— *Service.java              # Core implementations
15 |   |   |— *Records.java              # Data models
16 |   |— test/java/                     # Comprehensive test suite
17 |— build.gradle.kts                   # Test categories for cost control
18 |— slides.md                         # This presentation
```

- **8 Live Demos:** Ready-to-run examples for each topic
- **Cost-Controlled Testing:** Free local tests, cheap API tests

Repository Structure

```
1  AiJavaLabs/
2  |— labs.md           # 15 progressive lab exercises
3  |— src/
4  |   |— main/java/com/kousenit/
5  |   |   |— demos/
6  |   |   |   |— QuickChatDemo.java    # Fast OpenAI demo
7  |   |   |   |— TextToSpeechDemo.java # TTS generation
8  |   |   |   |— LocalOllamaDemo.java  # Local AI models
9  |   |   |   |— MultiModelDemo.java   # Provider comparison
10 |   |   |   |— ResponsesApiDemo.java  # Raw HTTP examples
11 |   |   |   |— ApiComparisonDemo.java # Framework vs raw
12 |   |   |   |— StreamingDemo.java     # Real-time responses
13 |   |   |— EasyRAGDemo.java           # Document Q&A
14 |   |   |— *Service.java              # Core implementations
15 |   |   |— *Records.java              # Data models
16 |   |— test/java/                     # Comprehensive test suite
17 |— build.gradle.kts                   # Test categories for cost control
18 |— slides.md                         # This presentation
```

- **8 Live Demos:** Ready-to-run examples for each topic
- **Cost-Controlled Testing:** Free local tests, cheap API tests
- **Modern Java 21:** Records, sealed interfaces, pattern matching

Repository Structure

```
1  AiJavaLabs/
2  |— labs.md           # 15 progressive lab exercises
3  |— src/
4  |   |— main/java/com/kousenit/
5  |   |   |— demos/
6  |   |   |   |— QuickChatDemo.java    # Fast OpenAI demo
7  |   |   |   |— TextToSpeechDemo.java # TTS generation
8  |   |   |   |— LocalOllamaDemo.java  # Local AI models
9  |   |   |   |— MultiModelDemo.java   # Provider comparison
10 |   |   |   |— ResponsesApiDemo.java  # Raw HTTP examples
11 |   |   |   |— ApiComparisonDemo.java # Framework vs raw
12 |   |   |   |— StreamingDemo.java     # Real-time responses
13 |   |   |— EasyRAGDemo.java           # Document Q&A
14 |   |   |— *Service.java              # Core implementations
15 |   |   |— *Records.java              # Data models
16 |   |— test/java/                     # Comprehensive test suite
17 |— build.gradle.kts                   # Test categories for cost control
18 |— slides.md                         # This presentation
```

- **8 Live Demos:** Ready-to-run examples for each topic
- **Cost-Controlled Testing:** Free local tests, cheap API tests
- **Modern Java 21:** Records, sealed interfaces, pattern matching

Course Philosophy: Understanding Both Levels

Why Learn Raw HTTP First? Then Master Frameworks

Course Philosophy: Understanding Both Levels

Why Learn Raw HTTP First? Then Master Frameworks

- **True Understanding:** Know what frameworks do under the hood

Course Philosophy: Understanding Both Levels

Why Learn Raw HTTP First? Then Master Frameworks

- **True Understanding:** Know what frameworks do under the hood
- **Debugging Skills:** When frameworks fail, you can troubleshoot

Course Philosophy: Understanding Both Levels

Why Learn Raw HTTP First? Then Master Frameworks

- **True Understanding:** Know what frameworks do under the hood
- **Debugging Skills:** When frameworks fail, you can troubleshoot
- **Flexibility:** Not locked into any single framework

Course Philosophy: Understanding Both Levels

Why Learn Raw HTTP First? Then Master Frameworks

- **True Understanding:** Know what frameworks do under the hood
- **Debugging Skills:** When frameworks fail, you can troubleshoot
- **Flexibility:** Not locked into any single framework
- **Performance:** Optimize for your specific needs

Course Philosophy: Understanding Both Levels

Why Learn Raw HTTP First?

- **True Understanding**: Know what frameworks do under the hood
- **Debugging Skills**: When frameworks fail, you can troubleshoot
- **Flexibility**: Not locked into any single framework
- **Performance**: Optimize for your specific needs

Then Master Frameworks

- **Productivity**: LangChain4j handles boilerplate

Course Philosophy: Understanding Both Levels

Why Learn Raw HTTP First?

- **True Understanding**: Know what frameworks do under the hood
- **Debugging Skills**: When frameworks fail, you can troubleshoot
- **Flexibility**: Not locked into any single framework
- **Performance**: Optimize for your specific needs

Then Master Frameworks

- **Productivity**: LangChain4j handles boilerplate
- **Best Practices**: Proven patterns and abstractions

Course Philosophy: Understanding Both Levels

Why Learn Raw HTTP First?

- **True Understanding**: Know what frameworks do under the hood
- **Debugging Skills**: When frameworks fail, you can troubleshoot
- **Flexibility**: Not locked into any single framework
- **Performance**: Optimize for your specific needs

Then Master Frameworks

- **Productivity**: LangChain4j handles boilerplate
- **Best Practices**: Proven patterns and abstractions
- **Advanced Features**: RAG, streaming, memory

Course Philosophy: Understanding Both Levels

Why Learn Raw HTTP First?

- **True Understanding**: Know what frameworks do under the hood
- **Debugging Skills**: When frameworks fail, you can troubleshoot
- **Flexibility**: Not locked into any single framework
- **Performance**: Optimize for your specific needs

Then Master Frameworks

- **Productivity**: LangChain4j handles boilerplate
- **Best Practices**: Proven patterns and abstractions
- **Advanced Features**: RAG, streaming, memory
- **Rapid Development**: Focus on business logic

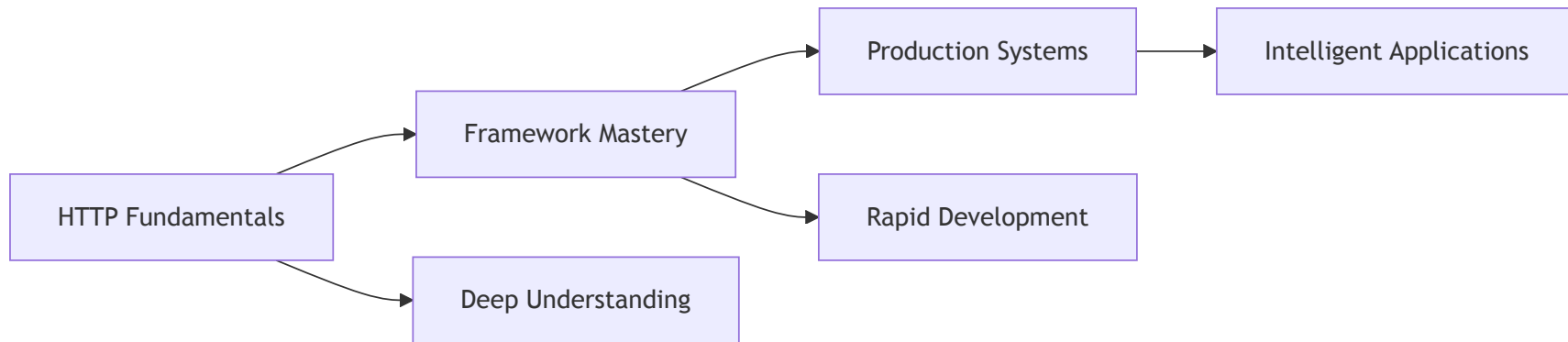
Course Philosophy: Understanding Both Levels

Why Learn Raw HTTP First?

- **True Understanding:** Know what frameworks do under the hood
- **Debugging Skills:** When frameworks fail, you can troubleshoot
- **Flexibility:** Not locked into any single framework
- **Performance:** Optimize for your specific needs

Then Master Frameworks

- **Productivity:** LangChain4j handles boilerplate
- **Best Practices:** Proven patterns and abstractions
- **Advanced Features:** RAG, streaming, memory
- **Rapid Development:** Focus on business logic



Prerequisites

Technical Requirements

Environment Setup

```
1  # Required API keys
2  export OPENAI_API_KEY=your_key
3
4  # Optional: Local AI models
5  curl -fsSL https://ollama.com/install.sh | sh
6  ollama pull gemma3
7  ollama pull moondream # For vision
8
9  # Clone and start
10 git clone <repo-url>
11 ./gradlew build
12 ./gradlew testDemo # Quick validation
```

Prerequisites

Technical Requirements

- **Java 21+** (Records, sealed interfaces, pattern matching)

Environment Setup

```
1  # Required API keys
2  export OPENAI_API_KEY=your_key
3
4  # Optional: Local AI models
5  curl -fsSL https://ollama.com/install.sh | sh
6  ollama pull gemma3
7  ollama pull moondream # For vision
8
9  # Clone and start
10 git clone <repo-url>
11 ./gradlew build
12 ./gradlew testDemo # Quick validation
```

Prerequisites

Technical Requirements

- **Java 21+** (Records, sealed interfaces, pattern matching)
- **Gradle 8.4+** (Kotlin DSL, test categories)

Environment Setup

```
1  # Required API keys
2  export OPENAI_API_KEY=your_key
3
4  # Optional: Local AI models
5  curl -fsSL https://ollama.com/install.sh | sh
6  ollama pull gemma3
7  ollama pull moondream # For vision
8
9  # Clone and start
10 git clone <repo-url>
11 ./gradlew build
12 ./gradlew testDemo # Quick validation
```

Prerequisites

Technical Requirements

- **Java 21+** (Records, sealed interfaces, pattern matching)
- **Gradle 8.4+** (Kotlin DSL, test categories)
- **LangChain4j 1.4.0** (Latest AI framework)

Environment Setup

```
1  # Required API keys
2  export OPENAI_API_KEY=your_key
3
4  # Optional: Local AI models
5  curl -fsSL https://ollama.com/install.sh | sh
6  ollama pull gemma3
7  ollama pull moondream # For vision
8
9  # Clone and start
10 git clone <repo-url>
11 ./gradlew build
12 ./gradlew testDemo # Quick validation
```

Prerequisites

Technical Requirements

- **Java 21+** (Records, sealed interfaces, pattern matching)
- **Gradle 8.4+** (Kotlin DSL, test categories)
- **LangChain4j 1.4.0** (Latest AI framework)
- **Git** for repository management

Environment Setup

```
1  # Required API keys
2  export OPENAI_API_KEY=your_key
3
4  # Optional: Local AI models
5  curl -fsSL https://ollama.com/install.sh | sh
6  ollama pull gemma3
7  ollama pull moondream # For vision
8
9  # Clone and start
10 git clone <repo-url>
11 ./gradlew build
12 ./gradlew testDemo # Quick validation
```

Prerequisites

Technical Requirements

- **Java 21+** (Records, sealed interfaces, pattern matching)
- **Gradle 8.4+** (Kotlin DSL, test categories)
- **LangChain4j 1.4.0** (Latest AI framework)
- **Git** for repository management
- **Ollama** (optional, for local AI models)

Environment Setup

```
1  # Required API keys
2  export OPENAI_API_KEY=your_key
3
4  # Optional: Local AI models
5  curl -fsSL https://ollama.com/install.sh | sh
6  ollama pull gemma3
7  ollama pull moondream # For vision
8
9  # Clone and start
10 git clone <repo-url>
11 ./gradlew build
12 ./gradlew testDemo # Quick validation
```


Cost Management Strategy

Test Categories

Gradle Tasks

Model Costs

```
1  ./gradlew testLocal      # $0
2  ./gradlew testCheap     # ~$0.0001
3  ./gradlew testDemo      # ~$0.0001
4  ./gradlew testNotExpensive
5  ./gradlew testOpenAI
```

Cost Management Strategy

Test Categories

- `@Tag("local")` - Free Ollama tests

Gradle Tasks

```
1  ./gradlew testLocal      # $0
2  ./gradlew testCheap      # ~$0.0
3  ./gradlew testDemo       # ~$0.0
4  ./gradlew testNotExpensive
5  ./gradlew testOpenAI
```

Model Costs

Cost Management Strategy

Test Categories

- `@Tag("local")` - **Free** Ollama tests
- `@Tag("cheap")` - **Low-cost** nano models

Gradle Tasks

```
1  ./gradlew testLocal      # $0
2  ./gradlew testCheap      # ~$0.0
3  ./gradlew testDemo       # ~$0.0
4  ./gradlew testNotExpensive
5  ./gradlew testOpenAI
```

Model Costs

Cost Management Strategy

Test Categories

- `@Tag("local")` - **Free** Ollama tests
- `@Tag("cheap")` - **Low-cost** nano models
- `@Tag("demo")` - **Fast** live demos

Gradle Tasks

```
1  ./gradlew testLocal      # $0
2  ./gradlew testCheap     # ~$0.0
3  ./gradlew testDemo      # ~$0.0
4  ./gradlew testNotExpensive
5  ./gradlew testOpenAI
```

Model Costs

Cost Management Strategy

Test Categories

- `@Tag("local")` - **Free** Ollama tests
- `@Tag("cheap")` - **Low-cost** nano models
- `@Tag("demo")` - **Fast** live demos
- `@Tag("expensive")` - **Avoided** by default

Gradle Tasks

```
1  ./gradlew testLocal      # $0
2  ./gradlew testCheap      # ~$0.0
3  ./gradlew testDemo       # ~$0.0
4  ./gradlew testNotExpensive
5  ./gradlew testOpenAI
```

Model Costs

Cost Management Strategy

Test Categories

- `@Tag("local")` - **Free** Ollama tests
- `@Tag("cheap")` - **Low-cost** nano models
- `@Tag("demo")` - **Fast** live demos
- `@Tag("expensive")` - **Avoided** by default

Gradle Tasks

```
1 ./gradlew testLocal      # $0
2 ./gradlew testCheap      # ~$0.0
3 ./gradlew testDemo       # ~$0.0
4 ./gradlew testNotExpensive
5 ./gradlew testOpenAI
```

Model Costs

- **gpt-4.1-nano**: \$0.150/1M input tokens

Cost Management Strategy

Test Categories

- `@Tag("local")` - **Free** Ollama tests
- `@Tag("cheap")` - **Low-cost** nano models
- `@Tag("demo")` - **Fast** live demos
- `@Tag("expensive")` - **Avoided** by default

Gradle Tasks

```
1  ./gradlew testLocal      # $0
2  ./gradlew testCheap      # ~$0.0
3  ./gradlew testDemo       # ~$0.0
4  ./gradlew testNotExpensive
5  ./gradlew testOpenAI
```

Model Costs

- **gpt-4.1-nano**: \$0.150/1M input tokens
- **gemma3 (local)**: Free with Ollama

Cost Management Strategy

Test Categories

- `@Tag("local")` - **Free** Ollama tests
- `@Tag("cheap")` - **Low-cost** nano models
- `@Tag("demo")` - **Fast** live demos
- `@Tag("expensive")` - **Avoided** by default

Gradle Tasks

```
1 ./gradlew testLocal      # $0
2 ./gradlew testCheap      # ~$0.0
3 ./gradlew testDemo       # ~$0.0
4 ./gradlew testNotExpensive
5 ./gradlew testOpenAI
```

Model Costs

- **gpt-4.1-nano**: \$0.150/1M input tokens
- **gemma3 (local)**: Free with Ollama
- **dall-e-3**: \$0.040 per image

Cost Management Strategy

Test Categories

- `@Tag("local")` - **Free** Ollama tests
- `@Tag("cheap")` - **Low-cost** nano models
- `@Tag("demo")` - **Fast** live demos
- `@Tag("expensive")` - **Avoided** by default

Gradle Tasks

```
1 ./gradlew testLocal      # $0
2 ./gradlew testCheap      # ~$0.0
3 ./gradlew testDemo       # ~$0.0
4 ./gradlew testNotExpensive
5 ./gradlew testOpenAI
```

Model Costs

- **gpt-4.1-nano**: \$0.150/1M input tokens
- **gemma3 (local)**: Free with Ollama
- **dall-e-3**: \$0.040 per image
- **tts-1**: \$0.015 per 1K chars

Cost Management Strategy

Test Categories

- `@Tag("local")` - **Free** Ollama tests
- `@Tag("cheap")` - **Low-cost** nano models
- `@Tag("demo")` - **Fast** live demos
- `@Tag("expensive")` - **Avoided** by default

Gradle Tasks

```
1 ./gradlew testLocal      # $0
2 ./gradlew testCheap      # ~$0.0
3 ./gradlew testDemo       # ~$0.0
4 ./gradlew testNotExpensive
5 ./gradlew testOpenAI
```

Model Costs

- **gpt-4.1-nano**: \$0.150/1M input tokens
- **gemma3 (local)**: Free with Ollama
- **dall-e-3**: \$0.040 per image
- **tts-1**: \$0.015 per 1K chars

Smart Strategy: Develop with free local models, deploy with optimal cloud models

Demo 1-2: HTTP Fundamentals

Understanding Raw API Integration

Demo 1: Text-to-Speech with Raw HTTP

```
1 // Raw HTTP approach - Understanding the fundamentals
2 public class TextToSpeechService {
3     private static final String OPENAI_API_KEY = System.getenv("OPENAI_API_KEY");
4     private static final HttpClient client = HttpClient.newHttpClient();
5
6     public Path generateMp3(String model, String input, String voice) {
7         // TODO: How do we construct the request?
8     }
9 }
```

Demo 1: Text-to-Speech with Raw HTTP

```
1 // Step 1: Create JSON payload manually
2 String payload = ""
3     {
4         "model": "%s",
5         "input": "%s",
6         "voice": "%s"
7     }
8     "".formatted(model, input.replaceAll("\\s+", " ").trim(), voice);
```

Demo 1: Text-to-Speech with Raw HTTP

```
1 // Step 2: Build HTTP request with proper headers
2 HttpRequest request = HttpRequest.newBuilder()
3   .uri(URI.create("https://api.openai.com/v1/audio/speech"))
4   .header("Authorization", "Bearer %s".formatted(OPENAI_API_KEY))
5   .header("Content-Type", "application/json")
6   .header("Accept", "audio/mpeg")
7   .POST(HttpRequest.BodyPublishers.ofString(payload))
8   .build();
```

Demo 1: Text-to-Speech with Raw HTTP

```
1 // Step 3: Complete implementation with file handling
2 public Path generateMp3(String model, String input, String voice) {
3     String payload = ""
4     {
5         "model": "%s",
6         "input": "%s",
7         "voice": "%s"
8     }
9     """".formatted(model, input.replaceAll("\\s+", " ").trim(), voice);
10
11     HttpRequest request = HttpRequest.newBuilder()
12         .uri(URI.create("https://api.openai.com/v1/audio/speech"))
13         .header("Authorization", "Bearer %s".formatted(OPENAI_API_KEY))
14         .header("Content-Type", "application/json")
15         .header("Accept", "audio/mpeg")
16         .POST(HttpRequest.BodyPublishers.ofString(payload))
17         .build();
18
19     try {
20         HttpResponse<Path> response =
21             client.send(request, HttpResponse.BodyHandlers.ofFile(getFilePath()));
22         return response.body();
23     } catch (IOException | InterruptedException e) {
24         throw new RuntimeException(e);
25     }
```

Demo 1: Text-to-Speech with Raw HTTP

```
1 // Step 3: Complete implementation with file handling
2 public Path generateMp3(String model, String input, String voice) {
3     String payload = ""
4     {
5         "model": "%s",
6         "input": "%s",
7         "voice": "%s"
8     }
9     """".formatted(model, input.replaceAll("\\s+", " ").trim(), voice);
10
11     HttpRequest request = HttpRequest.newBuilder()
12         .uri(URI.create("https://api.openai.com/v1/audio/speech"))
13         .header("Authorization", "Bearer %s".formatted(OPENAI_API_KEY))
14         .header("Content-Type", "application/json")
15         .header("Accept", "audio/mpeg")
16         .POST(HttpRequest.BodyPublishers.ofString(payload))
17         .build();
18
19     try {
20         HttpResponse<Path> response =
21             client.send(request, HttpResponse.BodyHandlers.ofFile(getFilePath()));
22         return response.body();
23     } catch (IOException | InterruptedException e) {
24         throw new RuntimeException(e);
25     }
```


Demo 2: Model Discovery & JSON Parsing

```
1 // Understanding API responses with Java records
2 public class OpenAiRecords {
3     public record ModelList(List<Model> data) {
4         public record Model(
5             String id,
6             long created,
7             @SerializedName("owned_by") String ownedBy) {
8         }
9     }
10 }
```

Demo 2: Model Discovery & JSON Parsing

```
1 // Service implementation with Gson parsing
2 public class OpenAiService {
3     private static final HttpClient client = HttpClient.newHttpClient();
4     private final Gson gson = new GsonBuilder()
5         .setFieldNamingPolicy(FieldNamingPolicy.LOWER_CASE_WITH_UNDERSCORES)
6         .create();
7
8     public ModelList listModels() throws IOException, InterruptedException {
9         HttpRequest request = HttpRequest.newBuilder()
10             .uri(URI.create("https://api.openai.com/v1/models"))
11             .header("Authorization", "Bearer " + System.getenv("OPENAI_API_KEY"))
12             .GET()
13             .build();
14
15         HttpResponse<String> response = client.send(request,
16             HttpResponse.BodyHandlers.ofString());
17         return gson.fromJson(response.body(), ModelList.class);
18     }
19 }
20
21     .build();
22
23     try {
24         HttpResponse<String> response =
25             client.send(request, HttpResponse.BodyHandlers.ofString());
```

Demo 2: Model Discovery & JSON Parsing

```
1 // Service implementation with Gson parsing
2 public class OpenAiService {
3     private static final HttpClient client = HttpClient.newHttpClient();
4     private final Gson gson = new GsonBuilder()
5         .setFieldNamingPolicy(FieldNamingPolicy.LOWER_CASE_WITH_UNDERSCORES)
6         .create();
7
8     public ModelList listModels() throws IOException, InterruptedException {
9         HttpRequest request = HttpRequest.newBuilder()
10             .uri(URI.create("https://api.openai.com/v1/models"))
11             .header("Authorization", "Bearer " + System.getenv("OPENAI_API_KEY"))
12             .GET()
13             .build();
14
15         HttpResponse<String> response = client.send(request,
16             HttpResponse.BodyHandlers.ofString());
17         return gson.fromJson(response.body(), ModelList.class);
18     }
19 }
20
21     .build();
22
23     try {
24         HttpResponse<String> response =
25             client.send(request, HttpResponse.BodyHandlers.ofString());
```

JSON Parsing: Gson vs Jackson

Gson Approach

```
1 // Using JsonElement tree navigation
2 JsonElement root = JsonParser.parseString(json);
3 JsonObject obj = root.getAsJsonObject();
4 JsonArray outputs = obj.getAsJsonArray("output");
5
6 for (JsonElement elem : outputs) {
7     JsonObject item = elem.getAsJsonObject();
8     if ("message".equals(item.get("type").getString())) {
9         // Extract nested content...
10    }
11 }
```

Jackson Approach

```
1 // Using JsonNode with at() method
2 ObjectMapper mapper = new ObjectMapper();
3 JsonNode root = mapper.readTree(json);
4
5 // Direct path with JSON Pointer!
6 String text = root.at("/output/0/content/0/text")
7     .asText();
8
9 // Or safe navigation with path()
10 root.path("output")
11     .path(0)
12     .path("content")
13     .path(0)
14     .path("text");
```

JSON Parsing: Gson vs Jackson

Gson Approach

```
1  // Using JsonElement tree navigation
2  JsonElement root = JsonParser.parseString(json);
3  JsonObject obj = root.getAsJsonObject();
4  JsonArray outputs = obj.getAsJsonArray("output");
5
6  for (JsonElement elem : outputs) {
7      JsonObject item = elem.getAsJsonObject();
8      if ("message".equals(item.get("type").getString))
9          // Extract nested content...
10 }
11 }
```

- Type-safe navigation
- Explicit type conversions
- Manual null checking needed

Jackson Approach

```
1  // Using JsonNode with at() method
2  ObjectMapper mapper = new ObjectMapper();
3  JsonNode root = mapper.readTree(json);
4
5  // Direct path with JSON Pointer!
6  String text = root.at("/output/0/content/0/text")
7                  .asText();
8
9  // Or safe navigation with path()
10 root.path("output")
11     .path(0)
12     .path("content")
13     .path(0)
14     .path("text");
```

JSON Parsing: Gson vs Jackson

Gson Approach

```
1 // Using JsonElement tree navigation
2 JsonElement root = JsonParser.parseString(json);
3 JsonObject obj = root.getAsJsonObject();
4 JsonArray outputs = obj.getAsJsonArray("output");
5
6 for (JsonElement elem : outputs) {
7     JsonObject item = elem.getAsJsonObject();
8     if ("message".equals(item.get("type").getString())) {
9         // Extract nested content...
10    }
11 }
```

- Type-safe navigation
- Explicit type conversions
- Manual null checking needed

Jackson Approach

```
1 // Using JsonNode with at() method
2 ObjectMapper mapper = new ObjectMapper();
3 JsonNode root = mapper.readTree(json);
4
5 // Direct path with JSON Pointer!
6 String text = root.at("/output/0/content/0/text")
7                 .asText();
8
9 // Or safe navigation with path()
10 root.path("output")
11     .path(0)
12     .path("content")
13     .path(0)
14     .path("text");
```

- JSON Pointer navigation (`at()`)
- Safe chaining with `path()`
- No null pointer exceptions

JSON Pointer: Elegant Deep Navigation

RFC 6901 JSON Pointer Syntax

JSON Pointer: Elegant Deep Navigation

RFC 6901 JSON Pointer Syntax

- **Direct path access:** `/output/0/content/0/text`

JSON Pointer: Elegant Deep Navigation

RFC 6901 JSON Pointer Syntax

- **Direct path access:** `/output/0/content/0/text`
- **Array indexing:** Use numbers for array elements

JSON Pointer: Elegant Deep Navigation

RFC 6901 JSON Pointer Syntax

- **Direct path access:** `/output/0/content/0/text`
- **Array indexing:** Use numbers for array elements
- **Nested objects:** Chain with `/` separator

JSON Pointer: Elegant Deep Navigation

RFC 6901 JSON Pointer Syntax

- **Direct path access:** `/output/0/content/0/text`
- **Array indexing:** Use numbers for array elements
- **Nested objects:** Chain with `/` separator

```
1 // Complex nested structure navigation - one line!
2 JsonNode root = mapper.readTree(response.body());
3
4 // Instead of multiple loops and null checks:
5 String text = root.at("/data/results/0/attributes/name").asText();
6
7 // With fallback:
8 String value = root.at("/missing/path").asText("default");
```

JSON Pointer: Elegant Deep Navigation

RFC 6901 JSON Pointer Syntax

- **Direct path access:** `/output/0/content/0/text`
- **Array indexing:** Use numbers for array elements
- **Nested objects:** Chain with `/` separator

```
1 // Complex nested structure navigation - one line!
2 JsonNode root = mapper.readTree(response.body());
3
4 // Instead of multiple loops and null checks:
5 String text = root.at("/data/results/0/attributes/name").asText();
6
7 // With fallback:
8 String value = root.at("/missing/path").asText("default");
```

When to Use Which?

- **Gson:** Already in project, simpler API, Google ecosystem
- **Jackson:** Spring Boot default, JSON Pointer support, more features
- **Both:** Can coexist - use what fits your needs!

HTTP vs Framework Comparison

Raw HTTP Approach

```
1 // Manual request construction
2 HttpRequest request = HttpRequest.newBuilder()
3   .uri(URI.create(url))
4   .header("Authorization", "Bearer " + key)
5   .header("Content-Type", "application/json")
6   .POST(HttpRequest.BodyPublishers.ofString(
7     gson.toJson(payload)))
8   .build();
9
10 // Manual response handling
11 HttpResponse<String> response =
12   client.send(request, HttpResponse.BodyHandlers.ofString())
13 return gson.fromJson(response.body(), ResponseClass.class);
```

LangChain4j Framework

```
1 // High-level abstraction
2 ChatModel model = OpenAiChatModel.builder()
3   .apiKey(System.getenv("OPENAI_API_KEY"))
4   .modelName("gpt-4.1-nano")
5   .build();
6
7 // Simple method call
8 String response = model.chat("Hello, world!");
```

HTTP vs Framework Comparison

Raw HTTP Approach

```
1 // Manual request construction
2 HttpRequest request = HttpRequest.newBuilder()
3   .uri(URI.create(url))
4   .header("Authorization", "Bearer " + key)
5   .header("Content-Type", "application/json")
6   .POST(HttpRequest.BodyPublishers.ofString(
7     gson.toJson(payload)))
8   .build();
9
10 // Manual response handling
11 HttpResponse<String> response =
12   client.send(request, HttpResponse.BodyHandlers.ofString())
13 return gson.fromJson(response.body(), ResponseClass.class);
```

LangChain4j Framework

```
1 // High-level abstraction
2 ChatModel model = OpenAiChatModel.builder()
3   .apiKey(System.getenv("OPENAI_API_KEY"))
4   .modelName("gpt-4.1-nano")
5   .build();
6
7 // Simple method call
8 String response = model.chat("Hello, world!");
```

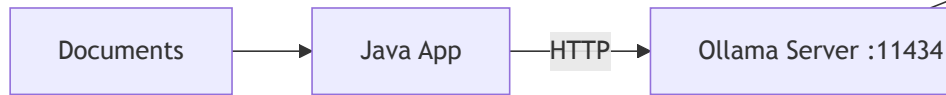
Both are valuable: Raw HTTP for understanding, frameworks for productivity

Demo 3-4: Local AI with Ollama

Privacy-First AI Integration

Demo 3: Ollama Setup and Architecture

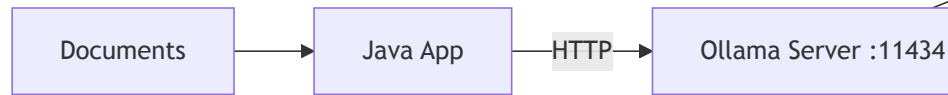
Why Ollama?



Demo 3: Ollama Setup and Architecture

Why Ollama?

- **Privacy:** No data leaves your machine



Demo 3: Ollama Setup and Architecture

Why Ollama?

- **Privacy:** No data leaves your machine
- **Cost:** Zero API costs after setup



Demo 3: Ollama Setup and Architecture

Why Ollama?

- **Privacy**: No data leaves your machine
- **Cost**: Zero API costs after setup
- **Control**: Full model customization



Demo 3: Ollama Setup and Architecture

Why Ollama?

- **Privacy**: No data leaves your machine
- **Cost**: Zero API costs after setup
- **Control**: Full model customization
- **Speed**: Local inference, no network delays



Demo 3: Ollama Setup and Architecture

Why Ollama?

- **Privacy**: No data leaves your machine
- **Cost**: Zero API costs after setup
- **Control**: Full model customization
- **Speed**: Local inference, no network delays



Setup Commands:

```
1  ollama pull gemma3      # Text generation
2  ollama pull moondream   # Vision analysis
3  ollama serve            # Start server (automatic on install)
```

Demo 4: Modern Java with Sealed Interfaces

```
1 // Step 1: Model different request types with records
2 public class OllamaRecords {
3     public record OllamaTextRequest(
4         String model,
5         String prompt,
6         boolean stream) {}
7
8     public record OllamaVisionRequest(
9         String model,
10        String prompt,
11        boolean stream,
12        List<String> images) {}
13 }
```

Demo 4: Modern Java with Sealed Interfaces

```
1 // Step 2: Use sealed interfaces for type safety (Java 17+)
2 public sealed interface OllamaRequest
3     permits OllamaTextRequest, OllamaVisionRequest {}
4
5 public record OllamaTextRequest(
6     String model, String prompt, boolean stream)
7     implements OllamaRequest {}
8
9 public record OllamaVisionRequest(
10     String model, String prompt, boolean stream, List<String> images)
11     implements OllamaRequest {
12
13     // Compact constructor with Base64 encoding
14     public OllamaVisionRequest {
15         images = images.stream()
16             .map(this::encodeImage)
17             .collect(Collectors.toList());
18     }
19 }
```

Demo 4: Modern Java with Sealed Interfaces

```
1 // Step 3: Pattern matching with switch expressions (Java 21)
2 public OllamaResponse generate(OllamaRequest request) {
3     switch (request) {
4         case OllamaTextRequest textRequest -> {
5             logger.log(INFO, "Text request: {0}", textRequest.prompt());
6         }
7         case OllamaVisionRequest visionRequest -> {
8             logger.log(INFO, "Vision request with {0} images",
9                 visionRequest.images().size());
10        }
11        // Exhaustive - no default needed with sealed types!
12    }
13
14    // Same HTTP logic for both request types
15    return sendRequest(request);
16 }
```


Demo 4: Modern Java with Sealed Interfaces

```
1 // Step 3: Pattern matching with switch expressions (Java 21)
2 public OllamaResponse generate(OllamaRequest request) {
3     switch (request) {
4         case OllamaTextRequest textRequest -> {
5             logger.log(INFO, "Text request: {0}", textRequest.prompt());
6         }
7         case OllamaVisionRequest visionRequest -> {
8             logger.log(INFO, "Vision request with {0} images",
9                 visionRequest.images().size());
10        }
11        // Exhaustive - no default needed with sealed types!
12    }
13
14    // Same HTTP logic for both request types
15    return sendRequest(request);
16 }
```

Modern Java Benefits: Type safety, exhaustive pattern matching, compact constructors

Demo 5-6: Framework Integration

LangChain4j Power and Simplicity

Demo 5: Quick Chat with LangChain4j

```
1 // From 50+ lines of HTTP code to this:
2 import dev.langchain4j.model.chat.ChatModel;
3 import dev.langchain4j.model.openai.OpenAiChatModel;
4
5 public class QuickChatDemo {
6     public static void main(String[] args) {
7         // TODO: How do we create a chat model?
8     }
9 }
```

Demo 5: Quick Chat with LangChain4j

```
1 // Step 1: Build the model with fluent API
2 ChatModel model = OpenAiChatModel.builder()
3     .apiKey(System.getenv("OPENAI_API_KEY"))
4     .modelName("gpt-4.1-nano")
5     .logRequests(false) // Toggle for debugging
6     .logResponses(false) // Toggle for debugging
7     .build();
```

Demo 5: Quick Chat with LangChain4j

```
1 // Step 2: Complete implementation - just one method call!
2 public class QuickChatDemo {
3     public static void main(String[] args) {
4         ChatModel model = OpenAiChatModel.builder()
5             .apiKey(System.getenv("OPENAI_API_KEY"))
6             .modelName("gpt-4.1-nano")
7             .build();
8
9         String prompt = "Explain what LangChain4j is in 2-3 sentences.";
10        System.out.println("Answer: " + model.chat(prompt));
11    }
12 }
```

Demo 5: Quick Chat with LangChain4j

```
1 // Step 2: Complete implementation - just one method call!
2 public class QuickChatDemo {
3     public static void main(String[] args) {
4         ChatModel model = OpenAiChatModel.builder()
5             .apiKey(System.getenv("OPENAI_API_KEY"))
6             .modelName("gpt-4.1-nano")
7             .build();
8
9         String prompt = "Explain what LangChain4j is in 2-3 sentences.";
10        System.out.println("Answer: " + model.chat(prompt));
11    }
12 }
```

Live Demo: Run `QuickChatDemo.java` - from complex HTTP to simple method call!

Demo 6: Multi-Model Provider Support

```
1 // Step 1: OpenAI implementation
2 ChatModel openai = OpenAiChatModel.builder()
3     .apiKey(System.getenv("OPENAI_API_KEY"))
4     .modelName("gpt-4.1-nano")
5     .build();
6
7 String response = openai.chat("What is the capital of France?");
8 System.out.println("OpenAI: " + response);
```

Demo 6: Multi-Model Provider Support

```
1 // Step 2: Switch to Google Gemini - same interface!
2 ChatModel gemini = GoogleAiGeminiChatModel.builder()
3     .apiKey(System.getenv("GOOGLEAI_API_KEY"))
4     .modelName("gemini-2.0-flash-exp")
5     .build();
6
7 String response = gemini.chat("What is the capital of France?");
8 System.out.println("Gemini: " + response);
```


Demo 6: Multi-Model Provider Support

```
1 // Step 3: Or use local Ollama - still same interface!
2 ChatModel ollama = OllamaChatModel.builder()
3     .baseUrl("http://localhost:11434")
4     .modelName("gemma3")
5     .build();
6
7 String response = ollama.chat("What is the capital of France?");
8 System.out.println("Ollama: " + response);
```

Demo 6: Multi-Model Provider Support

```
1 // Step 3: Or use local Ollama - still same interface!
2 ChatModel ollama = OllamaChatModel.builder()
3     .baseUrl("http://localhost:11434")
4     .modelName("gemma3")
5     .build();
6
7 String response = ollama.chat("What is the capital of France?");
8 System.out.println("Ollama: " + response);
```

Key Insight: Same `ChatModel` interface = vendor-agnostic code!

Provider Ecosystem

Cloud Providers

Local Models

**Specialized
Services**

Provider Ecosystem

Cloud Providers

- OpenAI • GPT-4.1-nano, DALL-

E 3

Local Models

Specialized Services

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM

Local Models

Specialized Services

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM
- **Anthropic** • Claude 3.5 Sonnet

Local Models

Specialized Services

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM
- **Anthropic** • Claude 3.5 Sonnet
- **Azure OpenAI** • Enterprise GPT

Local Models

Specialized Services

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM
- **Anthropic** • Claude 3.5 Sonnet
- **Azure OpenAI** • Enterprise GPT
- **AWS Bedrock** • Multiple models

Local Models

Specialized Services

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM
- **Anthropic** • Claude 3.5 Sonnet
- **Azure OpenAI** • Enterprise GPT
- **AWS Bedrock** • Multiple models

Local Models

- **Ollama** • gemma3, llama3.1, mistral

Specialized Services

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM
- **Anthropic** • Claude 3.5 Sonnet
- **Azure OpenAI** • Enterprise GPT
- **AWS Bedrock** • Multiple models

Local Models

- **Ollama** • gemma3, llama3.1, mistral
- **Hugging Face** • Open source models

Specialized Services

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM
- **Anthropic** • Claude 3.5 Sonnet
- **Azure OpenAI** • Enterprise GPT
- **AWS Bedrock** • Multiple models

Local Models

- **Ollama** • gemma3, llama3.1, mistral
- **Hugging Face** • Open source models
- **vLLM** • High-performance inference

Specialized Services

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM
- **Anthropic** • Claude 3.5 Sonnet
- **Azure OpenAI** • Enterprise GPT
- **AWS Bedrock** • Multiple models

Local Models

- **Ollama** • gemma3, llama3.1, mistral
- **Hugging Face** • Open source models
- **vLLM** • High-performance inference
- **LocalAI** • OpenAI-compatible API

Specialized Services

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM
- **Anthropic** • Claude 3.5 Sonnet
- **Azure OpenAI** • Enterprise GPT
- **AWS Bedrock** • Multiple models

Local Models

- **Ollama** • gemma3, llama3.1, mistral
- **Hugging Face** • Open source models
- **vLLM** • High-performance inference
- **LocalAI** • OpenAI-compatible API

Specialized Services

- **Vision** • moondream, LLaVA

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM
- **Anthropic** • Claude 3.5 Sonnet
- **Azure OpenAI** • Enterprise GPT
- **AWS Bedrock** • Multiple models

Local Models

- **Ollama** • gemma3, llama3.1, mistral
- **Hugging Face** • Open source models
- **vLLM** • High-performance inference
- **LocalAI** • OpenAI-compatible API

Specialized Services

- **Vision** • moondream, LLaVA
- **Audio** • Whisper, TTS models

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM
- **Anthropic** • Claude 3.5 Sonnet
- **Azure OpenAI** • Enterprise GPT
- **AWS Bedrock** • Multiple models

Local Models

- **Ollama** • gemma3, llama3.1, mistral
- **Hugging Face** • Open source models
- **vLLM** • High-performance inference
- **LocalAI** • OpenAI-compatible API

Specialized Services

- **Vision** • moondream, LLaVA
- **Audio** • Whisper, TTS models
- **Embeddings** • sentence-transformers

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM
- **Anthropic** • Claude 3.5 Sonnet
- **Azure OpenAI** • Enterprise GPT
- **AWS Bedrock** • Multiple models

Local Models

- **Ollama** • gemma3, llama3.1, mistral
- **Hugging Face** • Open source models
- **vLLM** • High-performance inference
- **LocalAI** • OpenAI-compatible API

Specialized Services

- **Vision** • moondream, LLaVA
- **Audio** • Whisper, TTS models
- **Embeddings** • sentence-transformers
- **Code** • CodeLlama, StarCoder

Provider Ecosystem

Cloud Providers

- **OpenAI** • GPT-4.1-nano, DALL-E 3
- **Google AI** • Gemini 2.0, PaLM
- **Anthropic** • Claude 3.5 Sonnet
- **Azure OpenAI** • Enterprise GPT
- **AWS Bedrock** • Multiple models

Local Models

- **Ollama** • gemma3, llama3.1, mistral
- **Hugging Face** • Open source models
- **vLLM** • High-performance inference
- **LocalAI** • OpenAI-compatible API

Specialized Services

- **Vision** • moondream, LLaVA
- **Audio** • Whisper, TTS models
- **Embeddings** • sentence-transformers
- **Code** • CodeLlama, StarCoder

LangChain4j Advantage: Switch providers with just configuration changes!

Demo 7: Streaming Responses

Real-Time AI Interactions

Demo 7: Clean Streaming with Lambda Handlers

```
1 // Old way: Complex anonymous inner classes
2 public void oldStreamingApproach() {
3     model.chat("Tell me a story", new StreamingChatResponseHandler() {
4         @Override
5         public void onPartialResponse(String token) {
6             System.out.print(token);
7         }
8
9         @Override
10        public void onError(Throwable error) {
11            System.err.println("Error: " + error);
12        }
13
14        @Override
15        public void onCompleted() {
16            System.out.println("\nDone!");
17        }
18    });
19 }
```

Demo 7: Clean Streaming with Lambda Handlers

```
1 // LangChain4j utilities: Clean lambda approach
2 import static dev.langchain4j.model.LambdaStreamingResponseHandler.*;
3
4 public class StreamingDemo {
5     public static void main(String[] args) {
6         var ollama = OllamaStreamingChatModel.builder()
7             .baseUrl("http://localhost:11434")
8             .modelName("gemma3")
9             .build();
10
11         // One-liner for simple streaming
12         ollama.chat("Tell me a haiku about Java",
13             onPartialResponse(System.out::print));
14     }
15 }
```

Demo 7: Clean Streaming with Lambda Handlers

```
1 // With error handling - still clean and functional
2 var openai = OpenAiStreamingChatModel.builder()
3     .apiKey(System.getenv("OPENAI_API_KEY"))
4     .modelName("gpt-4o-mini")
5     .build();
6
7 // Handle both tokens and errors elegantly
8 openai.chat("Why is the sky blue?",
9     onPartialResponseAndError(
10         System.out::print, // Handle each token
11         error -> System.err.println("Error: " + error.getMessage())
12     ));
```

Demo 7: Clean Streaming with Lambda Handlers

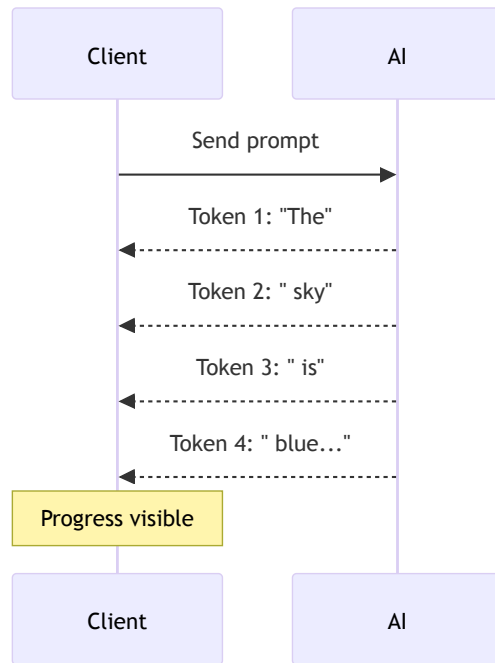
```
1 // With error handling - still clean and functional
2 var openai = OpenAiStreamingChatModel.builder()
3     .apiKey(System.getenv("OPENAI_API_KEY"))
4     .modelName("gpt-4o-mini")
5     .build();
6
7 // Handle both tokens and errors elegantly
8 openai.chat("Why is the sky blue?",
9     onPartialResponseAndError(
10         System.out::print, // Handle each token
11         error -> System.err.println("Error: " + error.getMessage())
12     ));
```

Live Demo: Run `StreamingDemo.java` to see tokens appear in real-time!

Streaming Architecture Benefits

User Experience

Technical Benefits

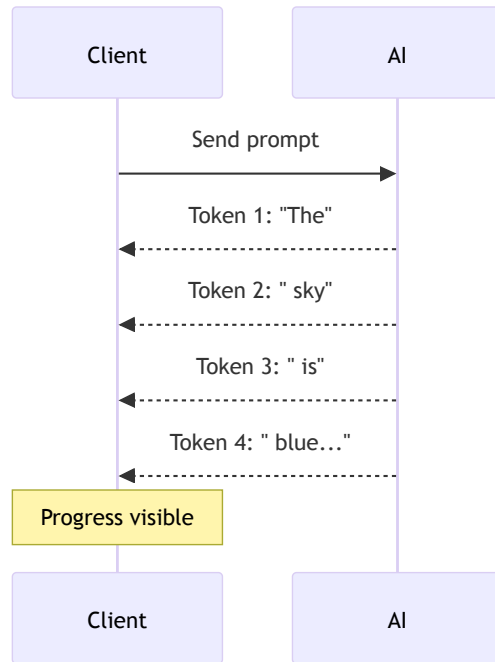


Streaming Architecture Benefits

User Experience

- **Immediate Feedback:** See progress instantly

Technical Benefits

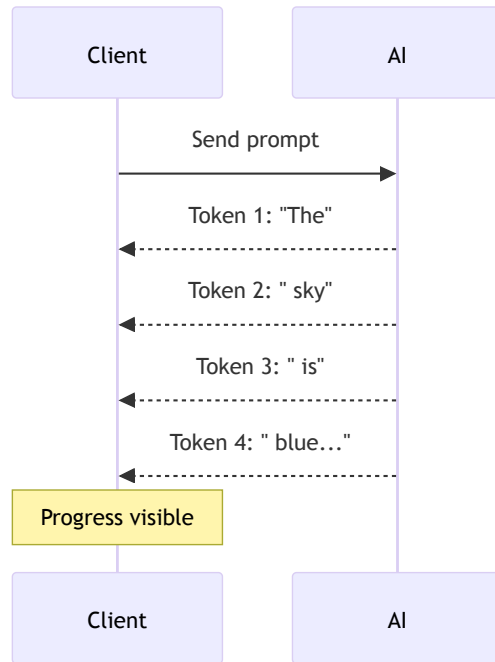


Streaming Architecture Benefits

User Experience

- **Immediate Feedback:** See progress instantly
- **Perceived Speed:** Feels faster

Technical Benefits

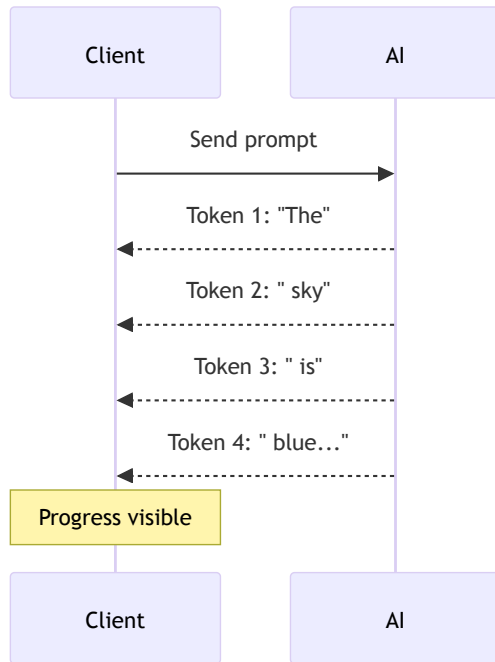


Streaming Architecture Benefits

User Experience

- **Immediate Feedback:** See progress instantly
- **Perceived Speed:** Feels faster
- **Interruptible:** Stop generation early

Technical Benefits

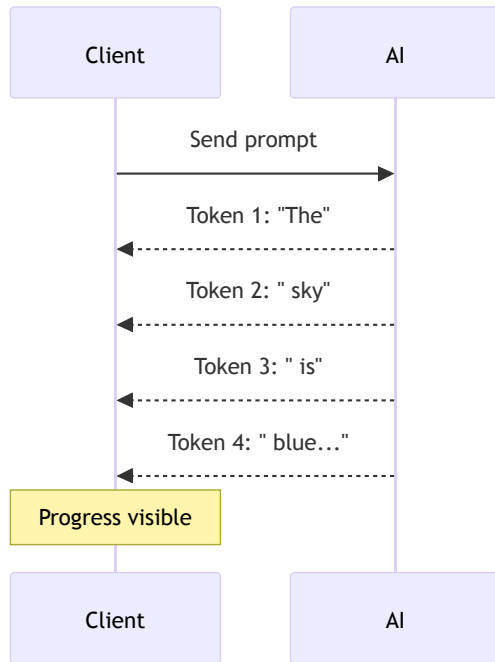


Streaming Architecture Benefits

User Experience

- **Immediate Feedback:** See progress instantly
- **Perceived Speed:** Feels faster
- **Interruptible:** Stop generation early
- **Engaging:** Watch AI "think"

Technical Benefits



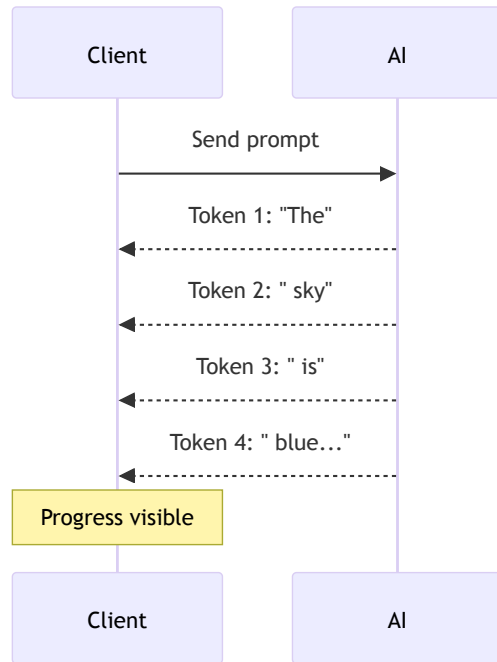
Streaming Architecture Benefits

User Experience

- **Immediate Feedback:** See progress instantly
- **Perceived Speed:** Feels faster
- **Interruptible:** Stop generation early
- **Engaging:** Watch AI "think"

Technical Benefits

- **Lower Latency:** First token quickly



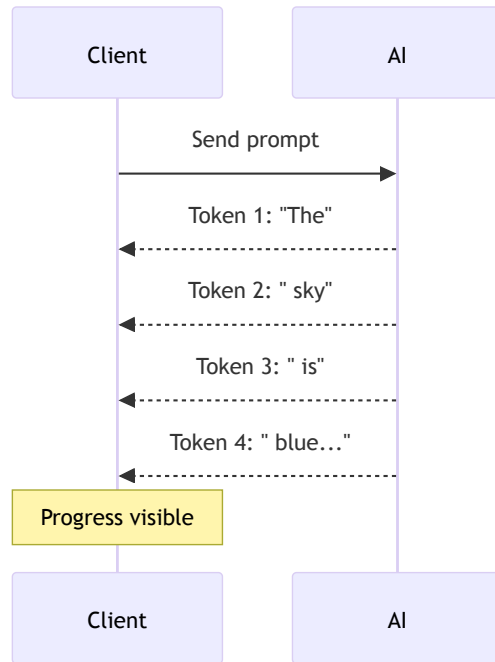
Streaming Architecture Benefits

User Experience

- **Immediate Feedback:** See progress instantly
- **Perceived Speed:** Feels faster
- **Interruptible:** Stop generation early
- **Engaging:** Watch AI "think"

Technical Benefits

- **Lower Latency:** First token quickly
- **Better Resources:** Process as arriving



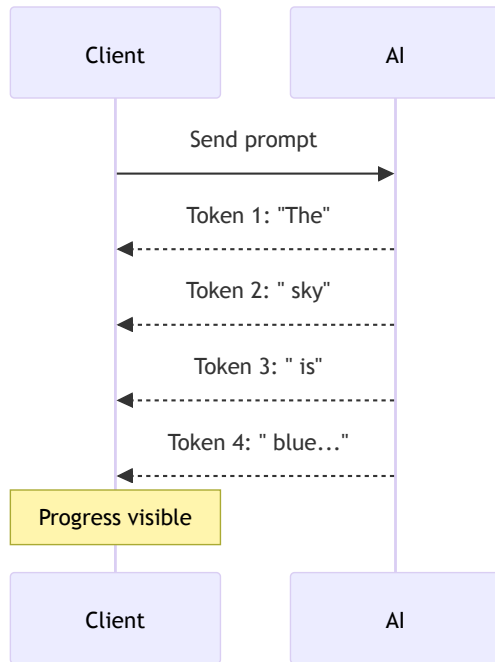
Streaming Architecture Benefits

User Experience

- **Immediate Feedback:** See progress instantly
- **Perceived Speed:** Feels faster
- **Interruptible:** Stop generation early
- **Engaging:** Watch AI "think"

Technical Benefits

- **Lower Latency:** First token quickly
- **Better Resources:** Process as arriving
- **Scalability:** Multiple streams



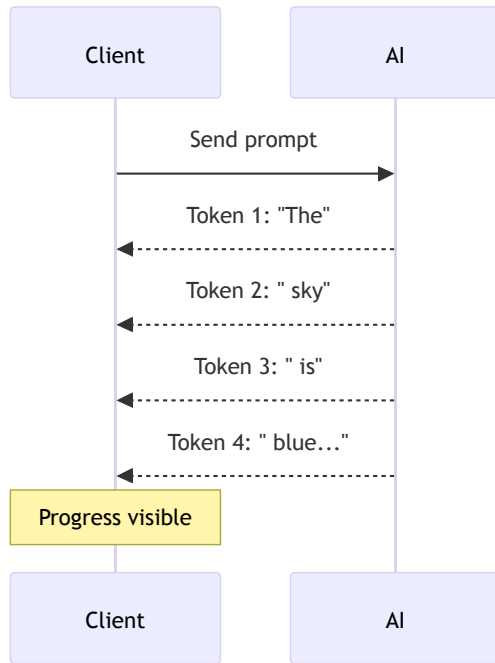
Streaming Architecture Benefits

User Experience

- **Immediate Feedback:** See progress instantly
- **Perceived Speed:** Feels faster
- **Interruptible:** Stop generation early
- **Engaging:** Watch AI "think"

Technical Benefits

- **Lower Latency:** First token quickly
- **Better Resources:** Process as arriving
- **Scalability:** Multiple streams
- **Error Recovery:** Graceful failures



Demo 8: Retrieval-Augmented Generation

Document-Powered AI

RAG: The Knowledge Problem

RAG: The Knowledge Problem

RAG: The Knowledge Problem

RAG: The Knowledge Problem

RAG: The Knowledge Problem

RAG: The Knowledge Problem

RAG Implementation with LangChain4j

```
1 // Step 1: Load and process documents
2 import static dev.langchain4j.data.document.loader.FileSystemDocumentLoader.loadDocuments;
3
4 List<Document> documents = loadDocuments(
5     toPath("documents/"),
6     glob("*.txt")
7 );
```

RAG Implementation with LangChain4j

```
1 // Step 2: Create vector store and ingest documents
2 InMemoryEmbeddingStore<TextSegment> embeddingStore =
3     new InMemoryEmbeddingStore<>();
4
5 // Magic happens here - documents become searchable vectors
6 EmbeddingStoreIngestor.ingest(documents, embeddingStore);
7
8 ContentRetriever retriever =
9     EmbeddingStoreContentRetriever.from(embeddingStore);
```


RAG Implementation with LangChain4j

```
1 // Step 3: Create AI assistant with document access
2 public interface Assistant {
3     String chat(String userMessage);
4 }
5
6 Assistant assistant = AiServices.builder(Assistant.class)
7     .chatModel(chatModel)
8     .chatMemory(MessageWindowChatMemory.withMaxMessages(10))
9     .contentRetriever(retriever) // The magic ingredient!
10    .build();
```

RAG Implementation with LangChain4j

```
1 // Step 4: Complete EasyRAGDemo implementation
2 public class EasyRAGDemo {
3     public static void main(String[] args) {
4         List<Document> documents = loadDocuments(toPath("documents/"), glob("*.txt"));
5
6         ChatModel chatModel = OpenAiChatModel.builder()
7             .apiKey(System.getenv("OPENAI_API_KEY"))
8             .modelName("gpt-4.1-nano")
9             .build();
10
11         Assistant assistant = AiServices.builder(Assistant.class)
12             .chatModel(chatModel)
13             .chatMemory(MessageWindowChatMemory.withMaxMessages(10))
14             .contentRetriever(createContentRetriever(documents))
15             .build();
16
17         // AI can now answer questions about your documents!
18         startConversationWith(assistant);
19     }
20 }
```

RAG Implementation with LangChain4j

```
1 // Step 4: Complete EasyRAGDemo implementation
2 public class EasyRAGDemo {
3     public static void main(String[] args) {
4         List<Document> documents = loadDocuments(toPath("documents/"), glob("*.txt"));
5
6         ChatModel chatModel = OpenAiChatModel.builder()
7             .apiKey(System.getenv("OPENAI_API_KEY"))
8             .modelName("gpt-4.1-nano")
9             .build();
10
11         Assistant assistant = AiServices.builder(Assistant.class)
12             .chatModel(chatModel)
13             .chatMemory(MessageWindowChatMemory.withMaxMessages(10))
14             .contentRetriever(createContentRetriever(documents))
15             .build();
16
17         // AI can now answer questions about your documents!
18         startConversationWith(assistant);
19     }
20 }
```

Live Demo: Run `EasyRAGDemo.java` - AI answers questions about YOUR documents!

RAG Deep Dive: How It Works

Vector Embeddings

Document Processing

RAG Deep Dive: How It Works

Vector Embeddings

- **Numerical Representation**: Text → vectors
(1536+ dimensions)

Document Processing

RAG Deep Dive: How It Works

Vector Embeddings

- **Numerical Representation**: Text → vectors
(1536+ dimensions)
- **Semantic Similarity**: Similar meaning = closer vectors

Document Processing

RAG Deep Dive: How It Works

Vector Embeddings

- **Numerical Representation**: Text → vectors
(1536+ dimensions)
- **Semantic Similarity**: Similar meaning = closer vectors
- **Searchable**: Find relevant content by vector distance

Document Processing

RAG Deep Dive: How It Works

Vector Embeddings

- **Numerical Representation**: Text → vectors (1536+ dimensions)
- **Semantic Similarity**: Similar meaning = closer vectors
- **Searchable**: Find relevant content by vector distance

Document Processing

- **Chunking**: Split documents into manageable pieces

RAG Deep Dive: How It Works

Vector Embeddings

- **Numerical Representation**: Text → vectors (1536+ dimensions)
- **Semantic Similarity**: Similar meaning = closer vectors
- **Searchable**: Find relevant content by vector distance

Document Processing

- **Chunking**: Split documents into manageable pieces
- **Embedding**: Convert chunks to vectors

RAG Deep Dive: How It Works

Vector Embeddings

- **Numerical Representation**: Text → vectors (1536+ dimensions)
- **Semantic Similarity**: Similar meaning = closer vectors
- **Searchable**: Find relevant content by vector distance

Document Processing

- **Chunking**: Split documents into manageable pieces
- **Embedding**: Convert chunks to vectors
- **Storage**: Index in vector database

RAG Deep Dive: How It Works

Vector Embeddings

- **Numerical Representation**: Text → vectors (1536+ dimensions)
- **Semantic Similarity**: Similar meaning = closer vectors
- **Searchable**: Find relevant content by vector distance

Document Processing

- **Chunking**: Split documents into manageable pieces
- **Embedding**: Convert chunks to vectors
- **Storage**: Index in vector database
- **Retrieval**: Find most relevant chunks for query

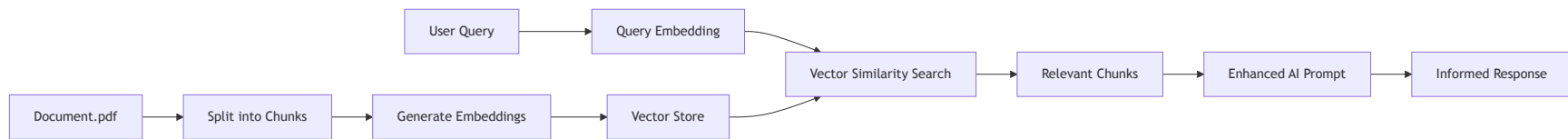
RAG Deep Dive: How It Works

Vector Embeddings

- **Numerical Representation**: Text → vectors (1536+ dimensions)
- **Semantic Similarity**: Similar meaning = closer vectors
- **Searchable**: Find relevant content by vector distance

Document Processing

- **Chunking**: Split documents into manageable pieces
- **Embedding**: Convert chunks to vectors
- **Storage**: Index in vector database
- **Retrieval**: Find most relevant chunks for query



Advanced Topics Overview

Production Considerations

Integration Patterns

Advanced Topics Overview

Production Considerations

- **Vector Databases:** Redis, Pinecone, Weaviate

Integration Patterns

Advanced Topics Overview

Production Considerations

- **Vector Databases:** Redis, Pinecone, Weaviate
- **Chunking Strategies:** Token-based, semantic, overlapping

Integration Patterns

Advanced Topics Overview

Production Considerations

- **Vector Databases:** Redis, Pinecone, Weaviate
- **Chunking Strategies:** Token-based, semantic, overlapping
- **Embedding Models:** OpenAI, sentence-transformers

Integration Patterns

Advanced Topics Overview

Production Considerations

- **Vector Databases:** Redis, Pinecone, Weaviate
- **Chunking Strategies:** Token-based, semantic, overlapping
- **Embedding Models:** OpenAI, sentence-transformers
- **Hybrid Search:** Vector + keyword search

Integration Patterns

Advanced Topics Overview

Production Considerations

- **Vector Databases:** Redis, Pinecone, Weaviate
- **Chunking Strategies:** Token-based, semantic, overlapping
- **Embedding Models:** OpenAI, sentence-transformers
- **Hybrid Search:** Vector + keyword search

Integration Patterns

- **Memory Management:** Conversation history

Advanced Topics Overview

Production Considerations

- **Vector Databases:** Redis, Pinecone, Weaviate
- **Chunking Strategies:** Token-based, semantic, overlapping
- **Embedding Models:** OpenAI, sentence-transformers
- **Hybrid Search:** Vector + keyword search

Integration Patterns

- **Memory Management:** Conversation history
- **Multi-Modal RAG:** Text, images, audio

Advanced Topics Overview

Production Considerations

- **Vector Databases:** Redis, Pinecone, Weaviate
- **Chunking Strategies:** Token-based, semantic, overlapping
- **Embedding Models:** OpenAI, sentence-transformers
- **Hybrid Search:** Vector + keyword search

Integration Patterns

- **Memory Management:** Conversation history
- **Multi-Modal RAG:** Text, images, audio
- **Real-Time Updates:** Dynamic document ingestion

Advanced Topics Overview

Production Considerations

- **Vector Databases:** Redis, Pinecone, Weaviate
- **Chunking Strategies:** Token-based, semantic, overlapping
- **Embedding Models:** OpenAI, sentence-transformers
- **Hybrid Search:** Vector + keyword search

Integration Patterns

- **Memory Management:** Conversation history
- **Multi-Modal RAG:** Text, images, audio
- **Real-Time Updates:** Dynamic document ingestion
- **Security:** Access control, data privacy

Testing Strategy

Smart Cost Control and Quality Assurance

Test Categories for Cost Control

Test Tags

```
1  @Test
2  @Tag("local") // Free
3  void testOllama() {
4      // $0 cost
5  }
6
7  @Test
8  @Tag("cheap") // Low cost
9  void testNano() {
10     // ~$0.001
11 }
12
13 @Test
14 @Tag("demo") // Fast
15 void quickDemo() {
16     // Optimized
17 }
```

Gradle Commands

```
1  # Free tests only
2  ./gradlew testLocal
3
4  # Low-cost tests
5  ./gradlew testCheap
6
7  # Demo tests
8  ./gradlew testDemo
9
10 # Exclude expensive
11 ./gradlew testNotExpensive
```

Test Categories for Cost Control

Test Tags

```
1  @Test
2  @Tag("local") // Free
3  void testOllama() {
4      // $0 cost
5  }
6
7  @Test
8  @Tag("cheap") // Low cost
9  void testNano() {
10     // ~$0.001
11 }
12
13 @Test
14 @Tag("demo") // Fast
15 void quickDemo() {
16     // Optimized
17 }
```

Gradle Commands

```
1  # Free tests only
2  ./gradlew testLocal
3
4  # Low-cost tests
5  ./gradlew testCheap
6
7  # Demo tests
8  ./gradlew testDemo
9
10 # Exclude expensive
11 ./gradlew testNotExpensive
```

Strategy: Develop with `testLocal` , validate with `testCheap`

Gradle Test Tasks

Cost-Controlled Testing

```
1  # Free tests only
2  ./gradlew testLocal
3
4  # Low-cost tests
5  ./gradlew testCheap
6
7  # Exclude expensive tests
8  ./gradlew testNotExpensive
9
10 # Quick live demos
11 ./gradlew testDemo
```

Provider-Specific Testing

```
1  # OpenAI API tests only
2  ./gradlew testOpenAI
3
4  # All tests (including expensive)
5  ./gradlew test
6
7  # Custom test selection
8  ./gradlew test --tests "*OpenAi*"
```

Gradle Test Tasks

Cost-Controlled Testing

```
1  # Free tests only
2  ./gradlew testLocal
3
4  # Low-cost tests
5  ./gradlew testCheap
6
7  # Exclude expensive tests
8  ./gradlew testNotExpensive
9
10 # Quick live demos
11 ./gradlew testDemo
```

Provider-Specific Testing

```
1  # OpenAI API tests only
2  ./gradlew testOpenAI
3
4  # All tests (including expensive)
5  ./gradlew test
6
7  # Custom test selection
8  ./gradlew test --tests "*OpenAi*"
```

Smart Development: Use free local models for TDD, validate with cloud models before deployment

Testing Patterns for AI Applications

Unit Testing Approaches

```
1  @Test
2  void shouldParseModelList() {
3      String json = ""
4          {"data": [
5              {"id": "gpt-4", "created": 123}
6          ]}
7      "";
8      ModelList result = gson.fromJson(json, ModelList.class);
9      assertThat(result.data()).hasSize(1);
10 }
11
12 @Test
13 void shouldHandleApiError() {
14     // Test error scenarios
15     assertThrows(RuntimeException.class,
16         () -> service.generateWithInvalidKey());
17 }
```

Integration Testing

```
1  @Test
2  @Tag("local")
3  void shouldStreamTokens() {
4      List<String> tokens = new ArrayList<>();
5
6      ollamaService.generateStreaming("gemma3",
7          "Count to 5",
8          token -> tokens.add(token)
9      );
10
11      assertThat(tokens).isNotEmpty();
12      assertThat(String.join("", tokens))
13          .containsPattern("1.*2.*3.*4.*5");
14 }
```

Quality Assurance for Non-Deterministic Systems

```
1  @Test
2  @Tag("cheap")
3  void shouldAnswerSkyColorQuestion() {
4      String response = chatModel.chat("Why is the sky blue?");
5
6      // Flexible content assertion
7      assertThat(response.toLowerCase())
8          .containsAnyOf("scattering", "wavelength", "atmosphere", "light");
9
10     // Structural assertion
11     assertThat(response).hasSizeGreaterThan(10);
12     assertThat(response).doesNotContain("I don't know");
13 }
```

Quality Assurance for Non-Deterministic Systems

- **Content Assertions:** Check for key concepts, not exact strings

```
1  @Test
2  @Tag("cheap")
3  void shouldAnswerSkyColorQuestion() {
4      String response = chatModel.chat("Why is the sky blue?");
5
6      // Flexible content assertion
7      assertThat(response.toLowerCase())
8          .containsAnyOf("scattering", "wavelength", "atmosphere", "light");
9
10     // Structural assertion
11     assertThat(response).hasSizeGreaterThan(10);
12     assertThat(response).doesNotContain("I don't know");
13 }
```

Quality Assurance for Non-Deterministic Systems

- **Content Assertions:** Check for key concepts, not exact strings
- **Structural Validation:** Verify JSON format, required fields

```
1  @Test
2  @Tag("cheap")
3  void shouldAnswerSkyColorQuestion() {
4      String response = chatModel.chat("Why is the sky blue?");
5
6      // Flexible content assertion
7      assertThat(response.toLowerCase())
8          .containsAnyOf("scattering", "wavelength", "atmosphere", "light");
9
10     // Structural assertion
11     assertThat(response).hasSizeGreaterThan(10);
12     assertThat(response).doesNotContain("I don't know");
13 }
```

Quality Assurance for Non-Deterministic Systems

- **Content Assertions:** Check for key concepts, not exact strings
- **Structural Validation:** Verify JSON format, required fields
- **Error Handling:** Test timeout, rate limits, invalid input

```
1  @Test
2  @Tag("cheap")
3  void shouldAnswerSkyColorQuestion() {
4      String response = chatModel.chat("Why is the sky blue?");
5
6      // Flexible content assertion
7      assertThat(response.toLowerCase())
8          .containsAnyOf("scattering", "wavelength", "atmosphere", "light");
9
10     // Structural assertion
11     assertThat(response).hasSizeGreaterThan(10);
12     assertThat(response).doesNotContain("I don't know");
13 }
```

Quality Assurance for Non-Deterministic Systems

- **Content Assertions:** Check for key concepts, not exact strings
- **Structural Validation:** Verify JSON format, required fields
- **Error Handling:** Test timeout, rate limits, invalid input
- **Performance Bounds:** Response time, token limits

```
1  @Test
2  @Tag("cheap")
3  void shouldAnswerSkyColorQuestion() {
4      String response = chatModel.chat("Why is the sky blue?");
5
6      // Flexible content assertion
7      assertThat(response.toLowerCase())
8          .containsAnyOf("scattering", "wavelength", "atmosphere", "light");
9
10     // Structural assertion
11     assertThat(response).hasSizeGreaterThan(10);
12     assertThat(response).doesNotContain("I don't know");
13 }
```


Quality Assurance for Non-Deterministic Systems

- **Content Assertions:** Check for key concepts, not exact strings
- **Structural Validation:** Verify JSON format, required fields
- **Error Handling:** Test timeout, rate limits, invalid input
- **Performance Bounds:** Response time, token limits
- **Cost Monitoring:** Track API usage in tests

```
1  @Test
2  @Tag("cheap")
3  void shouldAnswerSkyColorQuestion() {
4      String response = chatModel.chat("Why is the sky blue?");
5
6      // Flexible content assertion
7      assertThat(response.toLowerCase())
8          .containsAnyOf("scattering", "wavelength", "atmosphere", "light");
9
10     // Structural assertion
11     assertThat(response).hasSizeGreaterThan(10);
12     assertThat(response).doesNotContain("I don't know");
13 }
```

Course Summary

From HTTP to Intelligent Applications

Your AI Java Journey

Foundation Built

Advanced Capabilities

Your AI Java Journey

Foundation Built

- **HTTP Fundamentals:** Raw API integration patterns

Advanced Capabilities

Your AI Java Journey

Foundation Built

- **HTTP Fundamentals**: Raw API integration patterns
- **Modern Java**: Records, sealed interfaces, pattern matching

Advanced Capabilities

Your AI Java Journey

Foundation Built

- **HTTP Fundamentals**: Raw API integration patterns
- **Modern Java**: Records, sealed interfaces, pattern matching
- **Framework Mastery**: LangChain4j abstractions

Advanced Capabilities

Your AI Java Journey

Foundation Built

- **HTTP Fundamentals**: Raw API integration patterns
- **Modern Java**: Records, sealed interfaces, pattern matching
- **Framework Mastery**: LangChain4j abstractions
- **Cost Management**: Smart testing strategies

Advanced Capabilities

Your AI Java Journey

Foundation Built

- **HTTP Fundamentals**: Raw API integration patterns
- **Modern Java**: Records, sealed interfaces, pattern matching
- **Framework Mastery**: LangChain4j abstractions
- **Cost Management**: Smart testing strategies

Advanced Capabilities

- **Multi-Modal AI**: Text, vision, audio integration

Your AI Java Journey

Foundation Built

- **HTTP Fundamentals**: Raw API integration patterns
- **Modern Java**: Records, sealed interfaces, pattern matching
- **Framework Mastery**: LangChain4j abstractions
- **Cost Management**: Smart testing strategies

Advanced Capabilities

- **Multi-Modal AI**: Text, vision, audio integration
- **Streaming Responses**: Real-time user experiences

Your AI Java Journey

Foundation Built

- **HTTP Fundamentals**: Raw API integration patterns
- **Modern Java**: Records, sealed interfaces, pattern matching
- **Framework Mastery**: LangChain4j abstractions
- **Cost Management**: Smart testing strategies

Advanced Capabilities

- **Multi-Modal AI**: Text, vision, audio integration
- **Streaming Responses**: Real-time user experiences
- **RAG Systems**: Document-powered AI

Your AI Java Journey

Foundation Built

- **HTTP Fundamentals**: Raw API integration patterns
- **Modern Java**: Records, sealed interfaces, pattern matching
- **Framework Mastery**: LangChain4j abstractions
- **Cost Management**: Smart testing strategies

Advanced Capabilities

- **Multi-Modal AI**: Text, vision, audio integration
- **Streaming Responses**: Real-time user experiences
- **RAG Systems**: Document-powered AI
- **Production Patterns**: Error handling, monitoring

Your AI Java Journey

Foundation Built

- **HTTP Fundamentals**: Raw API integration patterns
- **Modern Java**: Records, sealed interfaces, pattern matching
- **Framework Mastery**: LangChain4j abstractions
- **Cost Management**: Smart testing strategies

Result: You can build intelligent Java applications with confidence! 🚀

Advanced Capabilities

- **Multi-Modal AI**: Text, vision, audio integration
- **Streaming Responses**: Real-time user experiences
- **RAG Systems**: Document-powered AI
- **Production Patterns**: Error handling, monitoring

Demo Classes Mastered

Core Demos

- [QuickChatDemo](#) - Fast OpenAI integration
- [TextToSpeechDemo](#) - Audio generation
- [LocalOllamaDemo](#) - Privacy-first AI
- [MultiModelDemo](#) - Provider comparison

Advanced Demos

- [StreamingDemo](#) - Real-time responses
- [ResponsesApiDemo](#) - Raw HTTP patterns
- [ApiComparisonDemo](#) - Framework vs raw
- [EasyRAGDemo](#) - Document Q&A

Demo Classes Mastered

Core Demos

- [QuickChatDemo](#) - Fast OpenAI integration
- [TextToSpeechDemo](#) - Audio generation
- [LocalOllamaDemo](#) - Privacy-first AI
- [MultiModelDemo](#) - Provider comparison

Advanced Demos

- [StreamingDemo](#) - Real-time responses
- [ResponsesApiDemo](#) - Raw HTTP patterns
- [ApiComparisonDemo](#) - Framework vs raw
- [EasyRAGDemo](#) - Document Q&A

All Ready for Live Demonstration: Each demo is optimized for training delivery

Key Takeaways

Key Takeaways

1. **Understand Both Levels:** Raw HTTP gives you debugging superpowers, frameworks give you productivity

Key Takeaways

1. **Understand Both Levels:** Raw HTTP gives you debugging superpowers, frameworks give you productivity
2. **Cost Management Matters:** Use test categories and local models for development

Key Takeaways

1. **Understand Both Levels**: Raw HTTP gives you debugging superpowers, frameworks give you productivity
2. **Cost Management Matters**: Use test categories and local models for development
3. **Modern Java Shines**: Records, sealed interfaces, and pattern matching are perfect for AI APIs

Key Takeaways

1. **Understand Both Levels**: Raw HTTP gives you debugging superpowers, frameworks give you productivity
2. **Cost Management Matters**: Use test categories and local models for development
3. **Modern Java Shines**: Records, sealed interfaces, and pattern matching are perfect for AI APIs
4. **LangChain4j is Powerful**: Unified API across all providers enables true portability

Key Takeaways

1. **Understand Both Levels**: Raw HTTP gives you debugging superpowers, frameworks give you productivity
2. **Cost Management Matters**: Use test categories and local models for development
3. **Modern Java Shines**: Records, sealed interfaces, and pattern matching are perfect for AI APIs
4. **LangChain4j is Powerful**: Unified API across all providers enables true portability
5. **RAG Changes Everything**: Give AI access to your documents and data

Key Takeaways

1. **Understand Both Levels**: Raw HTTP gives you debugging superpowers, frameworks give you productivity
2. **Cost Management Matters**: Use test categories and local models for development
3. **Modern Java Shines**: Records, sealed interfaces, and pattern matching are perfect for AI APIs
4. **LangChain4j is Powerful**: Unified API across all providers enables true portability
5. **RAG Changes Everything**: Give AI access to your documents and data
6. **Streaming Improves UX**: Real-time responses feel much faster to users

Key Takeaways

1. **Understand Both Levels**: Raw HTTP gives you debugging superpowers, frameworks give you productivity
2. **Cost Management Matters**: Use test categories and local models for development
3. **Modern Java Shines**: Records, sealed interfaces, and pattern matching are perfect for AI APIs
4. **LangChain4j is Powerful**: Unified API across all providers enables true portability
5. **RAG Changes Everything**: Give AI access to your documents and data
6. **Streaming Improves UX**: Real-time responses feel much faster to users
7. **Test Smart, Not Hard**: Use local models for TDD, cloud models for validation

Next Steps: Continue Your AI Journey

Hands-On Practice

Advanced Topics

Next Steps: Continue Your AI Journey

Hands-On Practice

- **Explore the repository:** Try all 8 demo classes

Advanced Topics

Next Steps: Continue Your AI Journey

Hands-On Practice

- **Explore the repository:** Try all 8 demo classes
- **Complete the labs:** 15 progressive exercises

Advanced Topics

Next Steps: Continue Your AI Journey

Hands-On Practice

- **Explore the repository:** Try all 8 demo classes
- **Complete the labs:** 15 progressive exercises
- **Build custom integrations:** Add your own AI providers

Advanced Topics

Next Steps: Continue Your AI Journey

Hands-On Practice

- **Explore the repository:** Try all 8 demo classes
- **Complete the labs:** 15 progressive exercises
- **Build custom integrations:** Add your own AI providers
- **Production deployment:** Scale with real applications

Advanced Topics

Next Steps: Continue Your AI Journey

Hands-On Practice

- **Explore the repository:** Try all 8 demo classes
- **Complete the labs:** 15 progressive exercises
- **Build custom integrations:** Add your own AI providers
- **Production deployment:** Scale with real applications

Advanced Topics

- **Vector databases:** Redis, Pinecone, Weaviate

Next Steps: Continue Your AI Journey

Hands-On Practice

- **Explore the repository:** Try all 8 demo classes
- **Complete the labs:** 15 progressive exercises
- **Build custom integrations:** Add your own AI providers
- **Production deployment:** Scale with real applications

Advanced Topics

- **Vector databases:** Redis, Pinecone, Weaviate
- **Multi-agent systems:** Coordinated AI workflows

Next Steps: Continue Your AI Journey

Hands-On Practice

- **Explore the repository:** Try all 8 demo classes
- **Complete the labs:** 15 progressive exercises
- **Build custom integrations:** Add your own AI providers
- **Production deployment:** Scale with real applications

Advanced Topics

- **Vector databases:** Redis, Pinecone, Weaviate
- **Multi-agent systems:** Coordinated AI workflows
- **Function calling:** Give AI access to your APIs

Next Steps: Continue Your AI Journey

Hands-On Practice

- **Explore the repository:** Try all 8 demo classes
- **Complete the labs:** 15 progressive exercises
- **Build custom integrations:** Add your own AI providers
- **Production deployment:** Scale with real applications

Advanced Topics

- **Vector databases:** Redis, Pinecone, Weaviate
- **Multi-agent systems:** Coordinated AI workflows
- **Function calling:** Give AI access to your APIs
- **Custom embeddings:** Domain-specific models

Next Steps: Continue Your AI Journey

Hands-On Practice

- **Explore the repository:** Try all 8 demo classes
- **Complete the labs:** 15 progressive exercises
- **Build custom integrations:** Add your own AI providers
- **Production deployment:** Scale with real applications

Resources:

- [LangChain4j Documentation](#)
- [Course Repository](#)
- [Ollama Models](#)

Advanced Topics

- **Vector databases:** Redis, Pinecone, Weaviate
- **Multi-agent systems:** Coordinated AI workflows
- **Function calling:** Give AI access to your APIs
- **Custom embeddings:** Domain-specific models

Production Considerations

Operational Excellence

Performance Optimization

Production Considerations

Operational Excellence

- **API Cost Monitoring:** Track usage across models

Performance Optimization

Production Considerations

Operational Excellence

- **API Cost Monitoring**: Track usage across models
- **Rate Limiting**: Handle quotas gracefully

Performance Optimization

Production Considerations

Operational Excellence

- **API Cost Monitoring**: Track usage across models
- **Rate Limiting**: Handle quotas gracefully
- **Error Recovery**: Circuit breakers, retries

Performance Optimization

Production Considerations

Operational Excellence

- **API Cost Monitoring:** Track usage across models
- **Rate Limiting:** Handle quotas gracefully
- **Error Recovery:** Circuit breakers, retries
- **Security:** API key management, input validation

Performance Optimization

Production Considerations

Operational Excellence

- **API Cost Monitoring**: Track usage across models
- **Rate Limiting**: Handle quotas gracefully
- **Error Recovery**: Circuit breakers, retries
- **Security**: API key management, input validation

Performance Optimization

- **Model Selection**: Right model for the task

Production Considerations

Operational Excellence

- **API Cost Monitoring**: Track usage across models
- **Rate Limiting**: Handle quotas gracefully
- **Error Recovery**: Circuit breakers, retries
- **Security**: API key management, input validation

Performance Optimization

- **Model Selection**: Right model for the task
- **Caching Strategies**: Cache embeddings and responses

Production Considerations

Operational Excellence

- **API Cost Monitoring**: Track usage across models
- **Rate Limiting**: Handle quotas gracefully
- **Error Recovery**: Circuit breakers, retries
- **Security**: API key management, input validation

Performance Optimization

- **Model Selection**: Right model for the task
- **Caching Strategies**: Cache embeddings and responses
- **Batch Processing**: Optimize API usage

Production Considerations

Operational Excellence

- **API Cost Monitoring**: Track usage across models
- **Rate Limiting**: Handle quotas gracefully
- **Error Recovery**: Circuit breakers, retries
- **Security**: API key management, input validation

Performance Optimization

- **Model Selection**: Right model for the task
- **Caching Strategies**: Cache embeddings and responses
- **Batch Processing**: Optimize API usage
- **Local Models**: Privacy and cost benefits

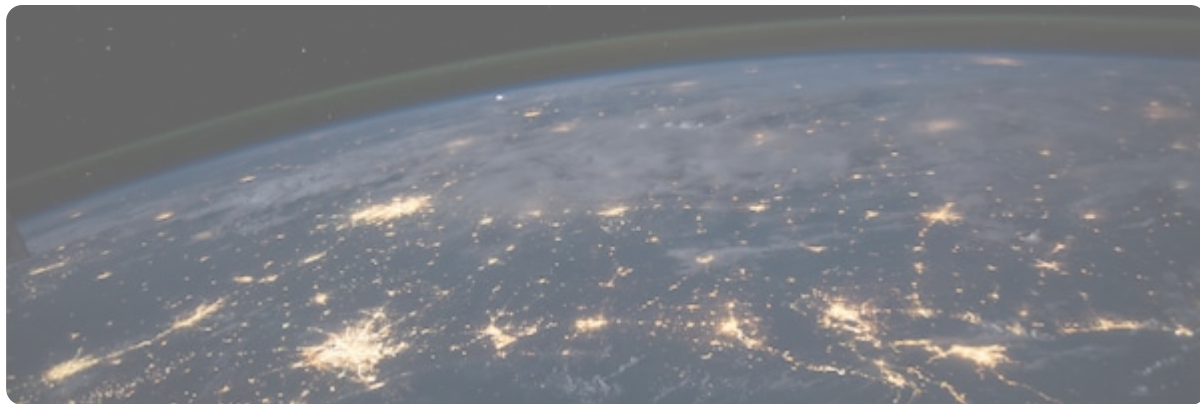
Thank You!

Questions?

Kenneth Kousen

Java Champion, Author, Speaker

kousenit.com | [@kenkousen](https://twitter.com/kenkousen)



Ready to build intelligent Java applications!