



Building RAG with OpenSearch ML Plugins: From PDFs to Voice-Enabled Search

Kushagra Sharma
Staff Software Engineer @ Genesys



What You'll Learn Today

- Deploy embedding models directly on OpenSearch using ML Commons
- Process PDFs and generate embeddings within OpenSearch
- Build semantic search with knn_vector indices
- Integrate voice search using STT/TTS models
- Create a complete RAG system



Why OpenSearch ML for RAG?

The Problem with Traditional RAG



External Embedding APIs



Vendor Lock-in



Operational Overhead

OpenSearch ML Solution



Native ML Capabilities



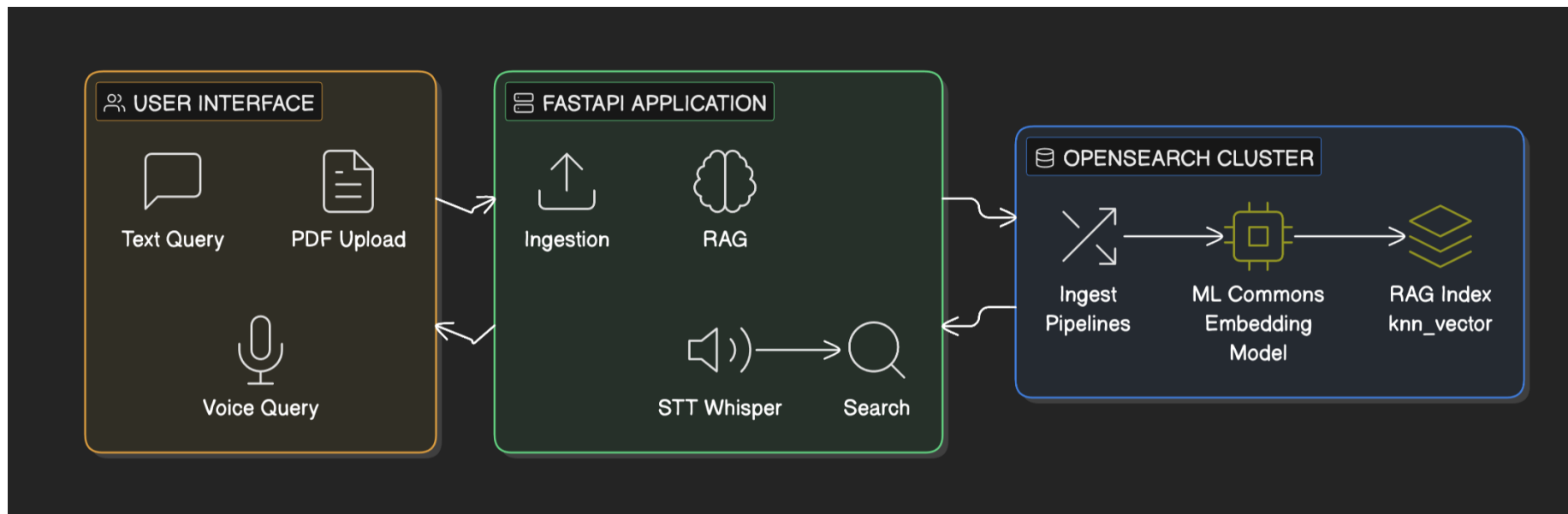
Self-Contained System



Production Ready



Architecture Overview



What is ML Commons?

ML Commons enables deploying and running ML models directly on OpenSearch cluster

- Native model inference
- No external APIs required
- Integrated with OpenSearch ecosystem
- Production-ready and scalable

Supported Model Formats

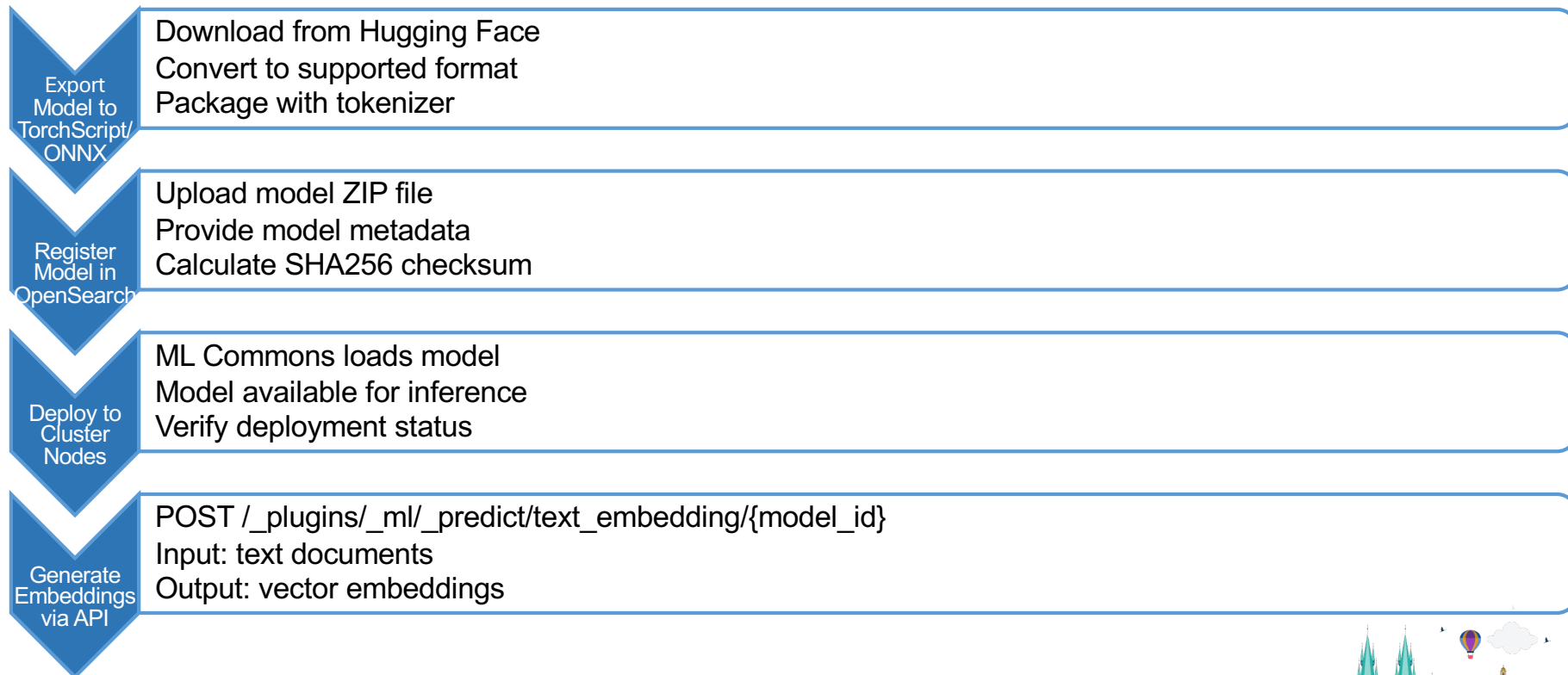
- TorchScript (PyTorch)
- ONNX (Open Neural Network Exchange)
- Both formats supported for flexibility

Model Types Supported

- Text Embedding Models
 - Sentence transformers
 - BERT-based models
 - Multilingual models
- Sparse Encoding Models
 - For keyword-based search
- Cross-Encoder Models
 - For re-ranking results



Deploying an Embedding Model: Step-by-Step



Model Deployment: Demo

```
mark > src > framework > opensearch > register_model.py > ...
3
4
5 from opensearchpy import OpenSearch
6
7 client = OpenSearch([{'host': 'localhost', 'port': 19200}])
8
9 # Register the model
10 response = client.transport.perform_request(
11     'POST',
12     '/_plugins/_ml/models/_register',
13     body={
14         "name": "paraphrase-multilingual-MiniLM-L12-v2",
15         "version": "1.0.0",
16         "model_format": "TORCH_SCRIPT",
17         "model_config": {
18             "model_type": "bert",
19             "embedding_dimension": 384,
20             "framework_type": "sentence_transformers"
21         },
22         "url": "file:///tmp/model.zip",
23         "model_content_hash_value": "4302b02e81aa7b137a381b9ad108e3258077fc3554751b903ae92fe8ada04bd3"
24     }
25 )
```

COMMENTS CODE REFERENCE LOG PROBLEMS 1 OUTPUT DEBUG CONSOLE TERMINAL PORTS

(venv) ksharma@g4767VTY27-D mark_2 %

zsh
Python

Jira: Kushagra Sharma: > No active issue Bitbucket: Kushagra Sharma: 0 0 1 0 0 Amazon Q | GenesysCloud Ln 19, Col 40 Spaces: 4 UTF-8 LF Python 3.13.0 (venv)



Document Ingestion Pipeline

PDF Processing

- Extract text from PDFs
- Use ingest-attachment plugin
- Supports multiple formats
- Handles complex documents

Text Chunking






- Split documents into chunks
- Configurable chunk size (500 words)
- Overlap between chunks (50 words)
- Preserve context across chunks

Embedding Generation

- Generate embeddings using ML model
- Use deployed model via ML API
- 384-dimensional vectors
- Index with document metadata

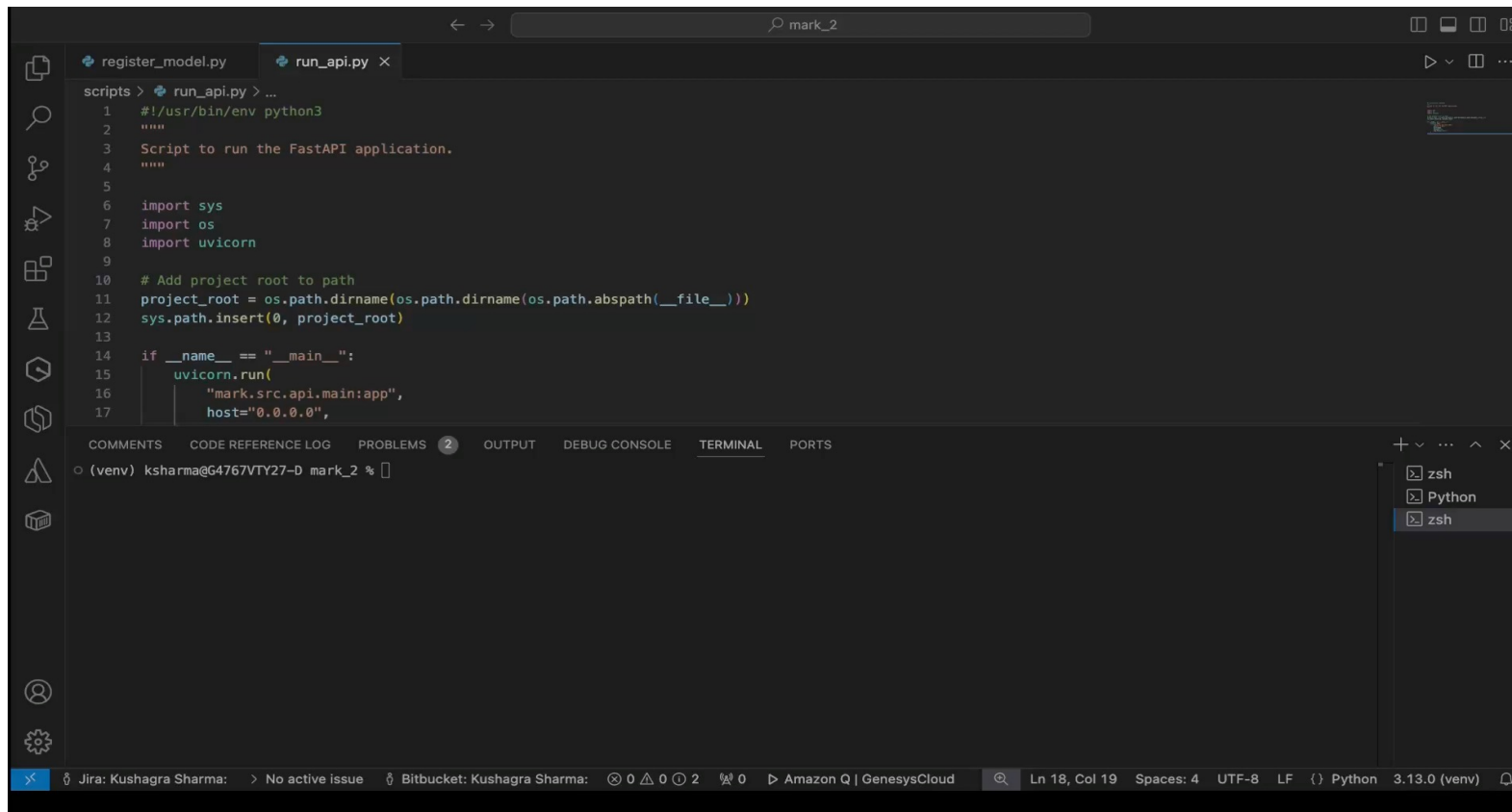


How Document Ingestion Works?

	User uploads PDF	<ul style="list-style-type: none">→ FastAPI receives file→ Validates file format
	Extract Text Content	<ul style="list-style-type: none">→ OpenSearch Attachment Pipeline→ Extracts text from PDF→ Returns plain text
	Split into Chunks	<ul style="list-style-type: none">→ Text → Chunks (500 words)→ Add overlap (50 words)→ Add metadata (source, page)
	Generate Embeddings	<ul style="list-style-type: none">→ Each chunk → ML API→ POST <code>/_plugins/_ml/_predict/text_embedding/{model_id}</code>→ Returns 384-dimensional vector
	Index in OpenSearch	<ul style="list-style-type: none">→ Chunk + Embedding→ Store in <code>rag_documents</code> index→ <code>knn_vector</code> field stores embedding



Document Ingestion: Demo



```
scripts > run_api.py > ...
1  #!/usr/bin/env python3
2  """
3  Script to run the FastAPI application.
4  """
5
6  import sys
7  import os
8  import uvicorn
9
10 # Add project root to path
11 project_root = os.path.dirname(os.path.dirname(os.path.abspath(__file__)))
12 sys.path.insert(0, project_root)
13
14 if __name__ == "__main__":
15     uvicorn.run(
16         "mark.src.api.main:app",
17         host="0.0.0.0",
```

COMMENTS CODE REFERENCE LOG PROBLEMS 2 OUTPUT DEBUG CONSOLE TERMINAL PORTS

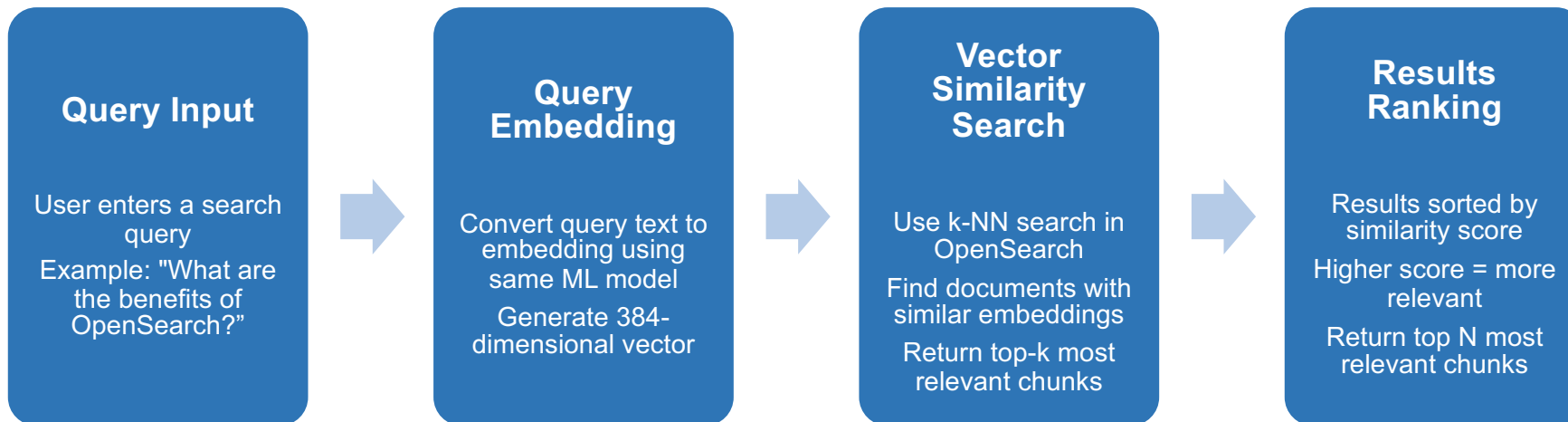
(venv) ksharma@G4767VTY27-D mark_2 %

- zsh
- Python
- zsh

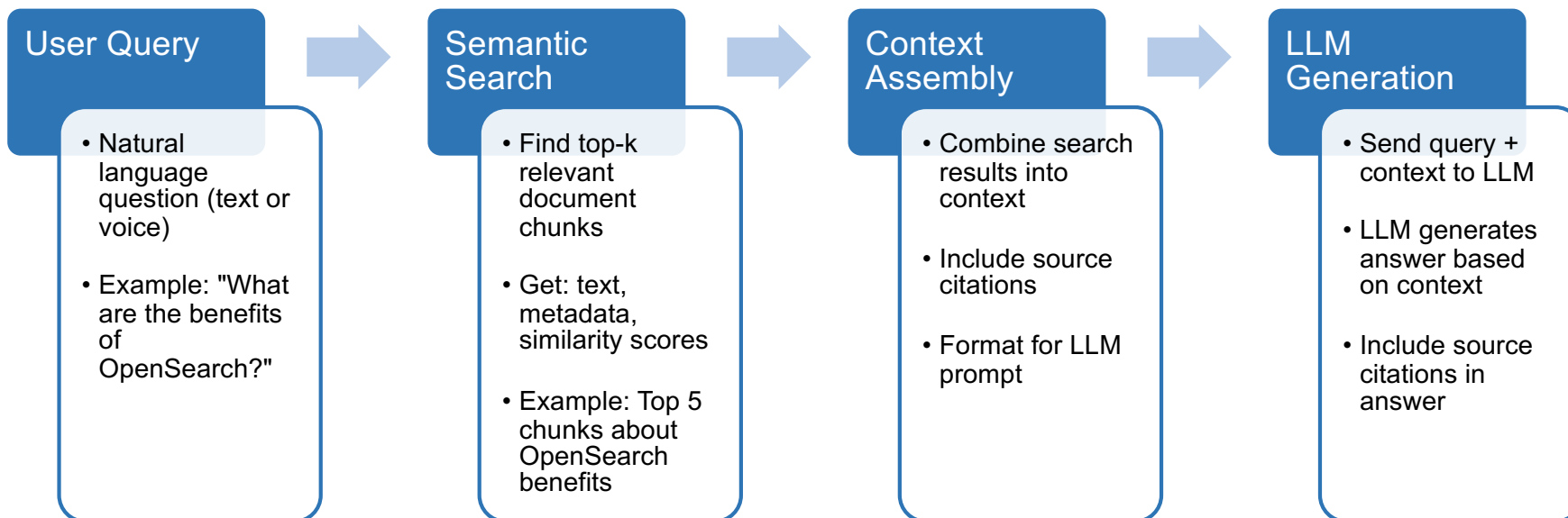
Jira: Kushagra Sharma: > No active issue Bitbucket: Kushagra Sharma: 0 0 0 2 0 Amazon Q | GenesysCloud Ln 18, Col 19 Spaces: 4 UTF-8 LF Python 3.13.0 (venv)



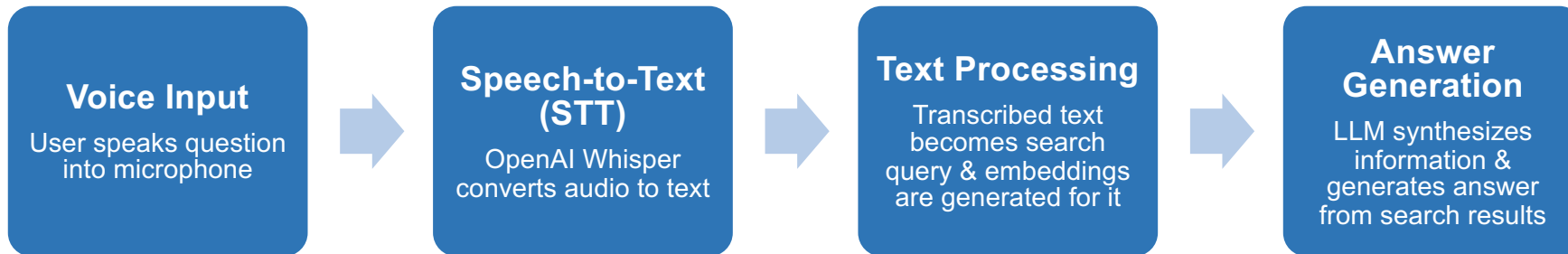
Semantic Search: Finding Relevant Documents



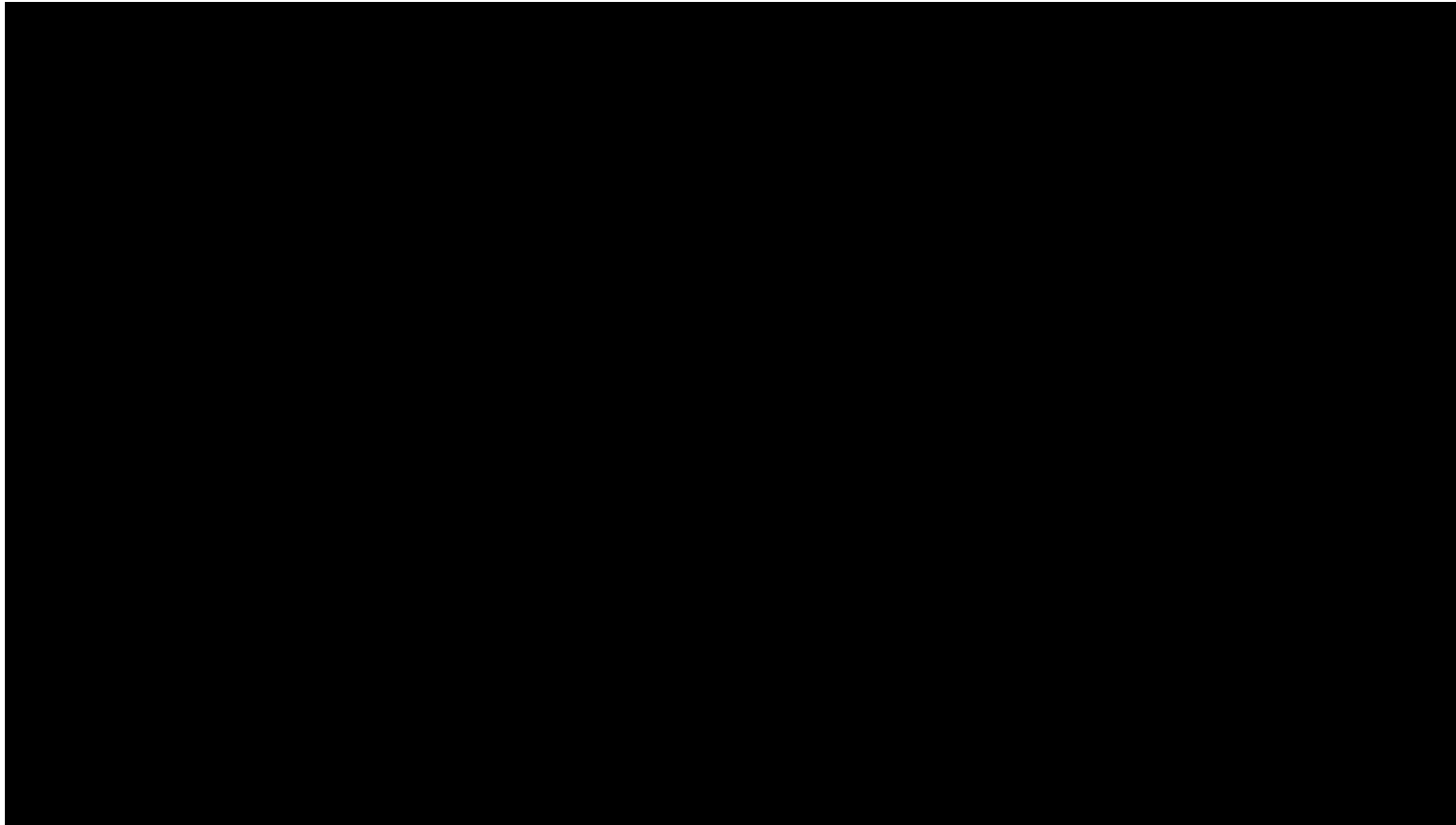
Answer Generation: RAG in Action



Voice-Enabled Search: Making RAG Accessible



Voice-Enabled Semantic Search with RAG: Demo



Key Takeaways

- OpenSearch ML Commons - Deploy embedding models directly on cluster - no external APIs
- Complete RAG Pipeline - PDF ingestion → Semantic search → Answer generation - all in OpenSearch
- Voice-Enabled Search - STT with Whisper, multilingual and accessible
- Self-Contained System - No external dependencies - private, cost-effective, and fully controlled



Benefits : Why this approach matters

- Data Privacy
 - No data sent to external services
 - Compliance-friendly (GDPR, HIPAA, etc.)
- Cost Efficiency
 - No per-request API fees
 - Open-source stack
- Performance
 - Low latency - no network calls
 - Fast embedding generation (<50ms)
- Flexibility
 - Customize models for your domain
 - Control over data and processing



Resources

OpenSearch ML Commons:

<https://docs.opensearch.org/latest/ml-commons-plugin/>

Mistral:

<https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>

Ollama:

<https://docs.ollama.com/>

OpenAI Whisper:

<https://huggingface.co/openai/whisper-base>



Thank you

Contact:

 kushshrma66@gmail.com

 <https://in.linkedin.com/in/kushagra-sharma-07>



OpenSearchCon

EUROPE

