# Unlock NextGen Product Search with ML and LLM Innovations
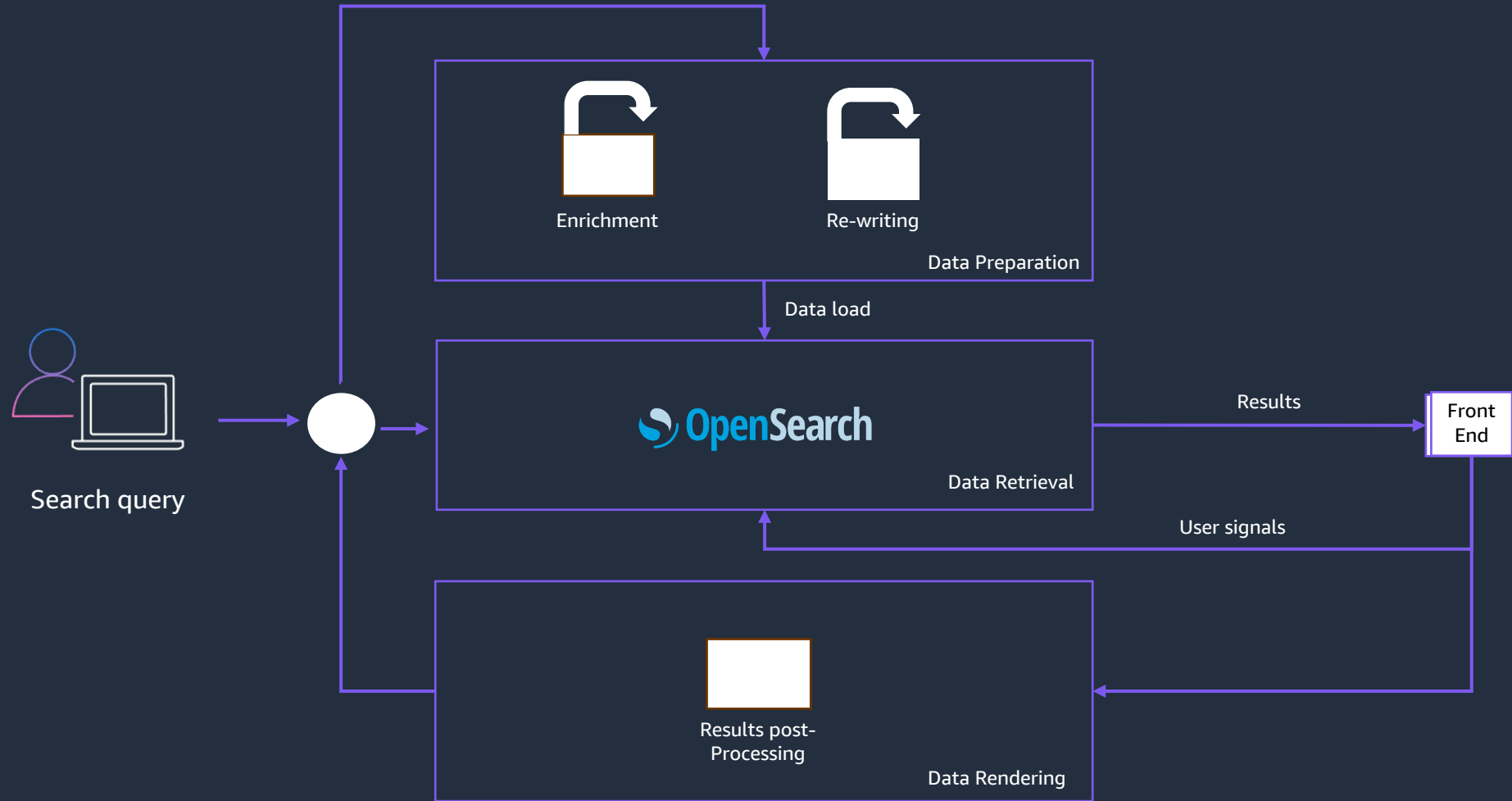
Praveen Mohan Prasad

Analytics Specialist TAM
AWS

Hajer Bouafif

Sr. WW Solutions Architect - OpenSearch
AWS

# Search lifecycle



Enrichment     Re-writing

Data Preparation

Data load

Search query

OpenSearch

Data Retrieval

Results

Front End

User signals

Results post-Processing

Data Rendering

# OpenSearch Growth

Apache 2.0 License
Linux Foundation

**>700MM**
OpenSearch project downloads since launch in Q3 2021

**Top 4 search engine**
DB-Engines ranking

**75+**
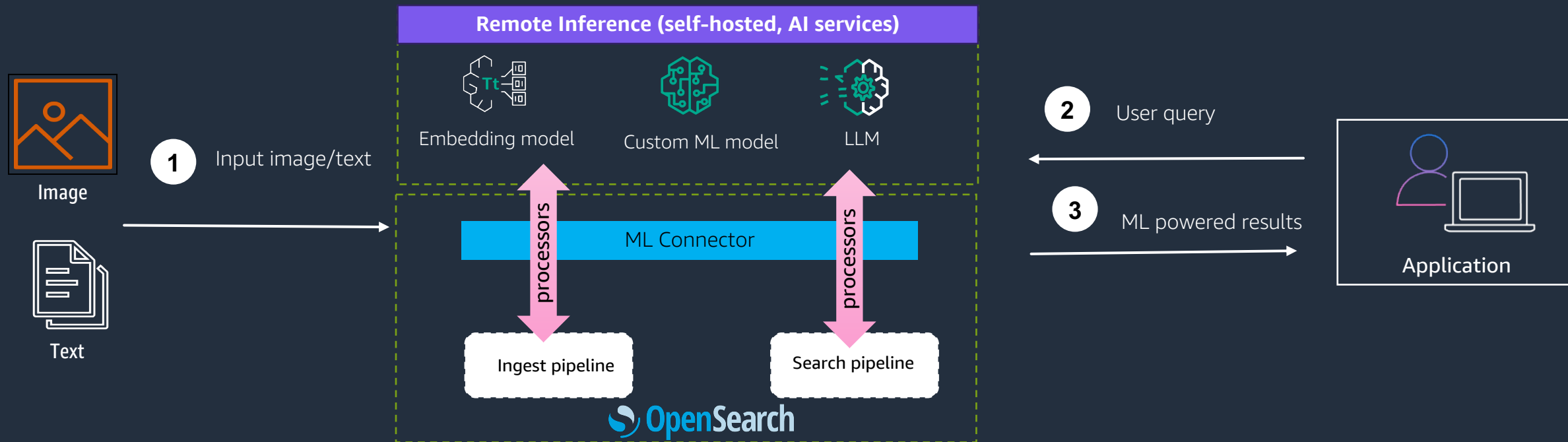Partners and growing

**100s of new features**
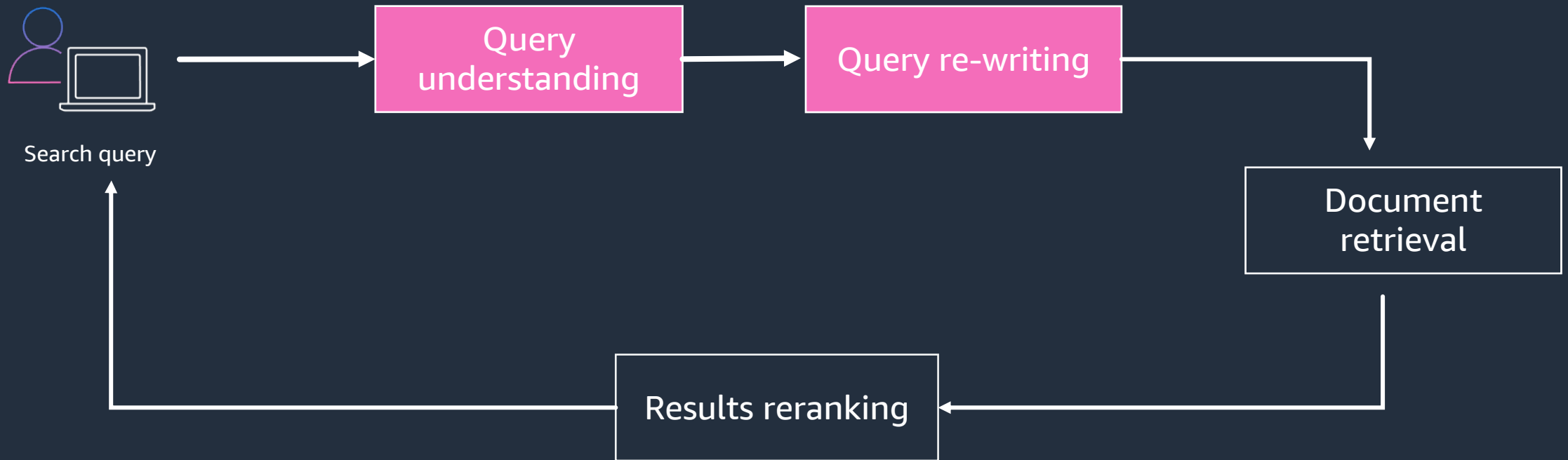32k+ pull requests merged
100+ weekly community contributions

**Multiple service providers**
AWS, Oracle, Aiven – Azure and GCP, Bonsai-Azure and GCP

# ML integrations made easy with OpenSearch

# ML-Powered Search Lifecycle



Search query → Query understanding → Query re-writing → Document retrieval → Results reranking → (back to Search query)

# Query understanding starts at ingestion time
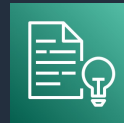


Amazon Rekognition
or
Object detection models

```
{

"description": "Rugged Brown Leather Boots",
"price": "50
"color":"brown",
"category":"Apparel and Accessories",
"objects":"Footwear,Boot,Shoe,Clothing"

}
```

```
{
"Text": "Bob ordered two
sandwiches and three ice cream
cones today from a store in
Seattle."
}
```

Amazon Comprehend
or
NER models

```
{

{"Text": "Bob", "Type": "PERSON" },
{"Text": "two sandwiches", "Type": "QUANTITY" },
{"Text": "three ice cream cones", "Type": "QUANTITY" },
{"Text": "today", "Type": "DATE" },
{"Text": "Seattle", "Type": "LOCATION" }

}
```

**Ingest pipeline: ML inference Processor**

# And applies at search time

## User Query

"Brown leather shoes for men under 50$"

## Search Query

```
{
    "query": {
        "bool": {
            "must": [
                { "match":{ "description":"shoes"}}
                    ],

            "filter": [
                { "term":   { "category": "footwear" }},
                { "range": { "price":{ "lt": 50 }},
                { "term":   { "gender":"male" }},
                { "term":   { "color": "brown"}}
            ]
        }
    }
}
```

"nearby" => "term": {"location": "Berlin"}

"Now" => "term": {"time": "now"}

# Query re-writing using LLMs

**Search Query**

*"Brown leather shoes for men under 50$"*

→

**Structured Query**

```
{
    "query": "shoes",

    "filter": "and
        (
            eq("category", "footwear"),
            lt("price", 50),
            eq("gender", "male")
            eq("color", "brown")
        )"
}
```

Rewrite →

**OpenSearch query DSL**

```
{
    "query": {
        "bool": {
            "must": [
                { "match":{ "description":"shoes"}}
                ],

            "filter": [
                { "term":  { "category": "footwear" }},
                { "range": { "price":{ "lt": 50 }},
                { "term":  { "gender":"male" }},
                { "term":  { "color": "brown"}}
            ]
        }
    }
}
```

GitHub code sample

# Cache to reduce the LLM cost and latency



Lexical/kNN Search

User Query

kNN index

PUT index/_doc

Search latency
3 sec → 600 ms

No

Match

Yes

LLM

Return Structured Query

- Warm up cache on regular intervals
- ISM for auto-deleting the cache (index) based on TTL
- Langchain built-in semantic cache with OpenSearch

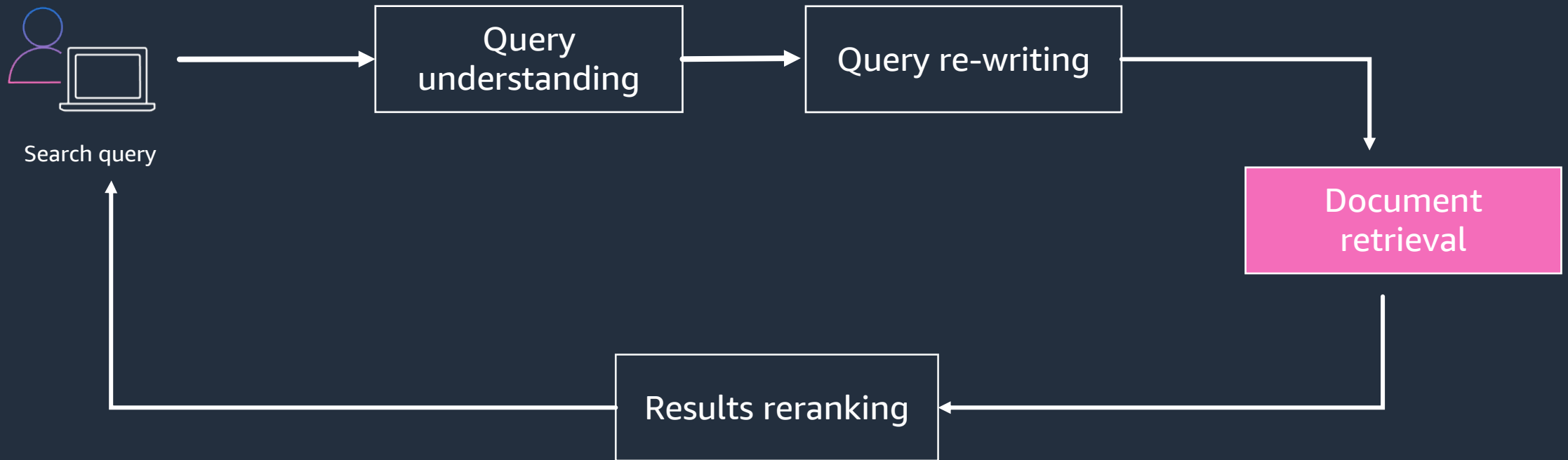# ML-Powered Search Lifecycle

# Retrieval Types

## Sparse retrieval

**BM25**

**Lexical search**
Uses keyword-based matching using TF/IDF.

**Neural sparse search**
Uses `amazon/neural-sparse` embedding models

## Dense retrieval

**KNN/ANN**

**Vector search**
Uses purpose build models for text, image or video embedding models.

**Multimodal search**
Uses embedding models that share embedding space between image and text.

## Hybrid Search

**BM25+KNN**

**Ensemble search**

Uses both traditional keyword-search

combined with vector search

## Conversational search

**(BM25/KNN)+LLM**

**RAG**
Uses LLM to augment the results retrieved from Vector Search.

# Neural sparse search

**Document :**
**"Apple Products are**
**expensive"**

Ingest →

OpenSearch
inverted Index

← Neural sparse query

**Query:**
**"apple**
**headphones"**

```
{
"apple":3.32
"expensive":2.49
"cheap":1.87
"products":1.85
"cost":1.57
"product":1.49
"technology":1.32
"mac":0.59
.
.
"fruit":0.02
"foods":0.01
}
```

Sparse encoding model

```
{
"apple":3.3
"head":2.28
"##phones":2.08
"sound":0.78
"music":0.73
"device":0.6
"nike":0.5
"wireless":0.49
"phone":0.4
"loud":0.39
.
.
.
"hardware":0.05
"version":0.04
}
```

Retrieve with
tokens >= 0.5

Re-rank with
all tokens

Ingest pipeline: Sparse Encoding processor

Search pipeline: Sparse two phase processor

# Improved sparse search performance maintaining relevance

**OpenSearch Sparse models v2**

1. opensearch-neural-sparse-encoding-v2-distill

2. opensearch-neural-sparse-encoding-doc-v2-distill

3. opensearch-neural-sparse-encoding-v2-mini

https://huggingface.co/opensearch-project

| Performance @ CPU | VS Sparse models v1 |
|---|---|
| Ingestion ThroughPut | ↑ 1.74x – 4.18x |
| Search Latency | ↓ 30% |
| BIER: ndcg@10 | ~= |

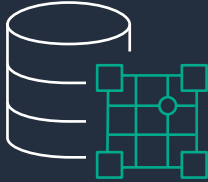https://opensearch.org/blog/neural-sparse-v2-models

# Dense search



Crafted with premium materials, these versatile black sneakers feature a sleek, minimalist look to complement any outfit while providing lasting comfort for urban exploration and adventures.

Ingest

**OpenSearch knn Index**

Neural query

Query: Styling Kicks
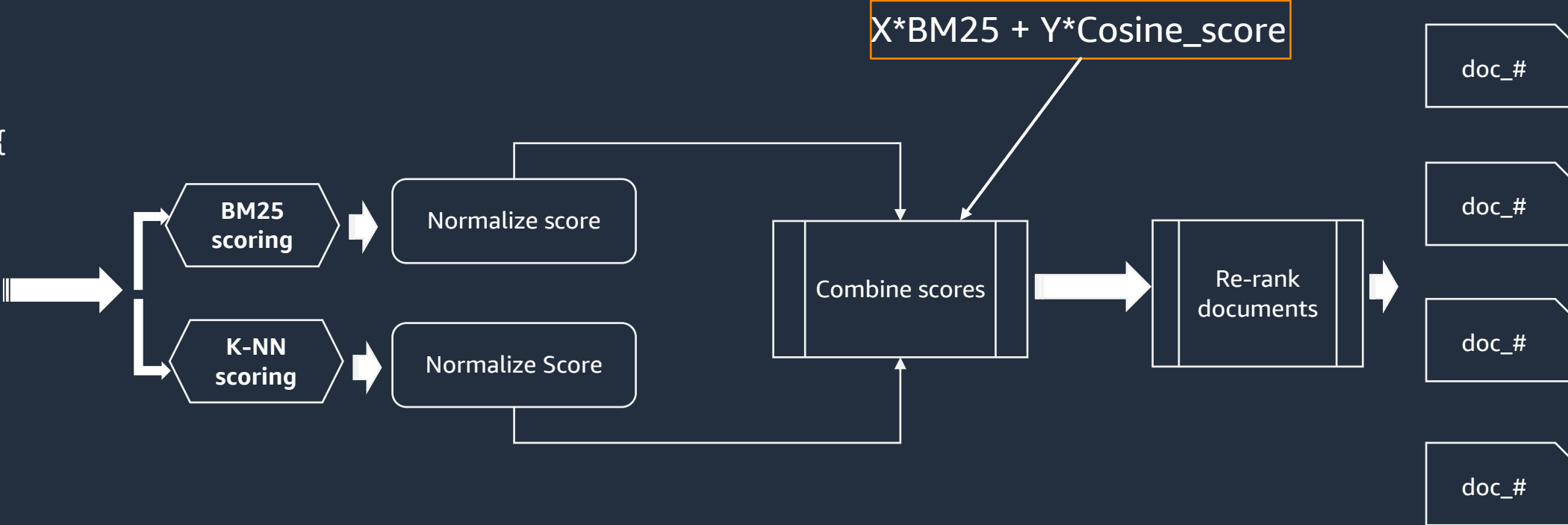
**Dense Embedding model**

Model fine-tuning with synthetic data

**Ingest pipeline: Text/Text-image embedding processor**

# Hybrid search

```
"query": { "hybrid": {
"queries": [ {
    "match":{
      <text query>
  }
},
  {
   "neural":{
    <vector query>
   }
}]
```

X*BM25 + Y*Cosine_score

BM25 scoring → Normalize score

K-NN scoring → Normalize Score

Combine scores → Re-rank documents

doc_#
doc_#
doc_#
doc_#

**Search pipeline: Normalization processor**

# Conversational search



**Search pipeline: RAG processor**

# ML-Powered Search Lifecycle



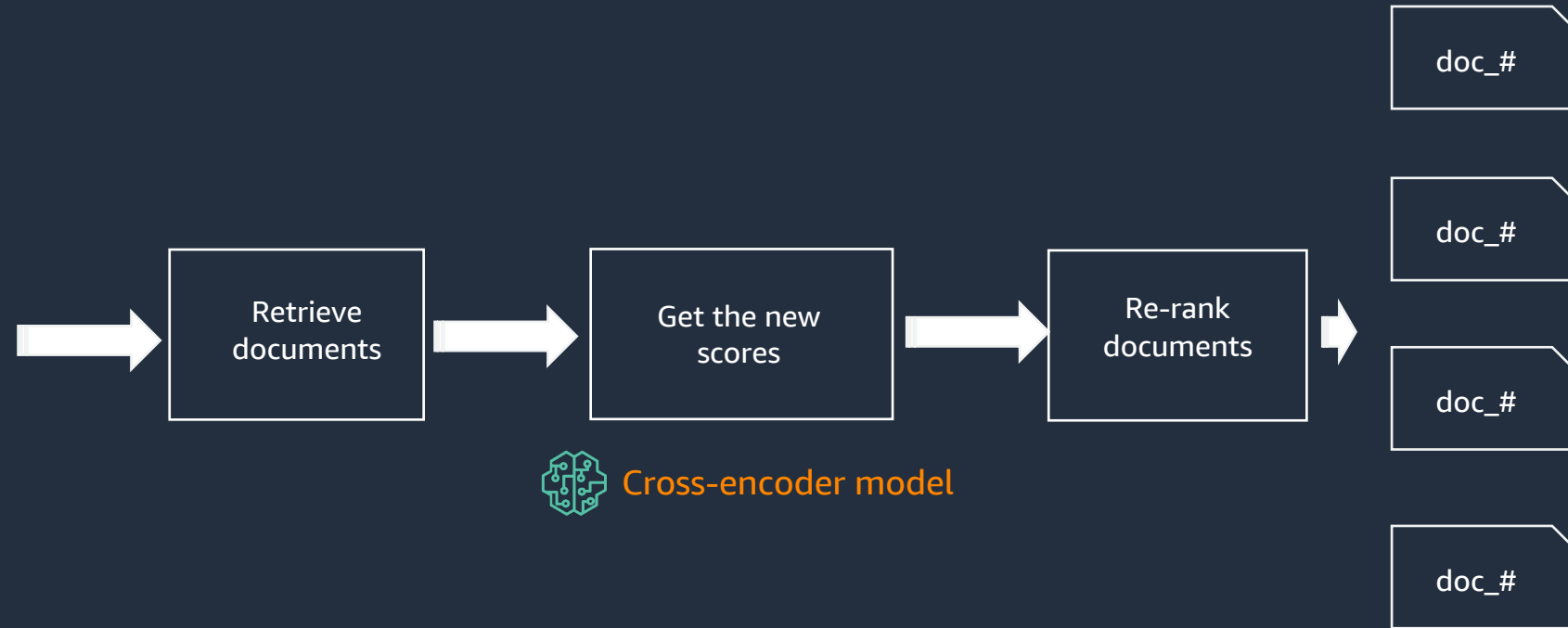Search query → Query understanding → Query re-writing → Document retrieval → Results reranking → (back to Search query)
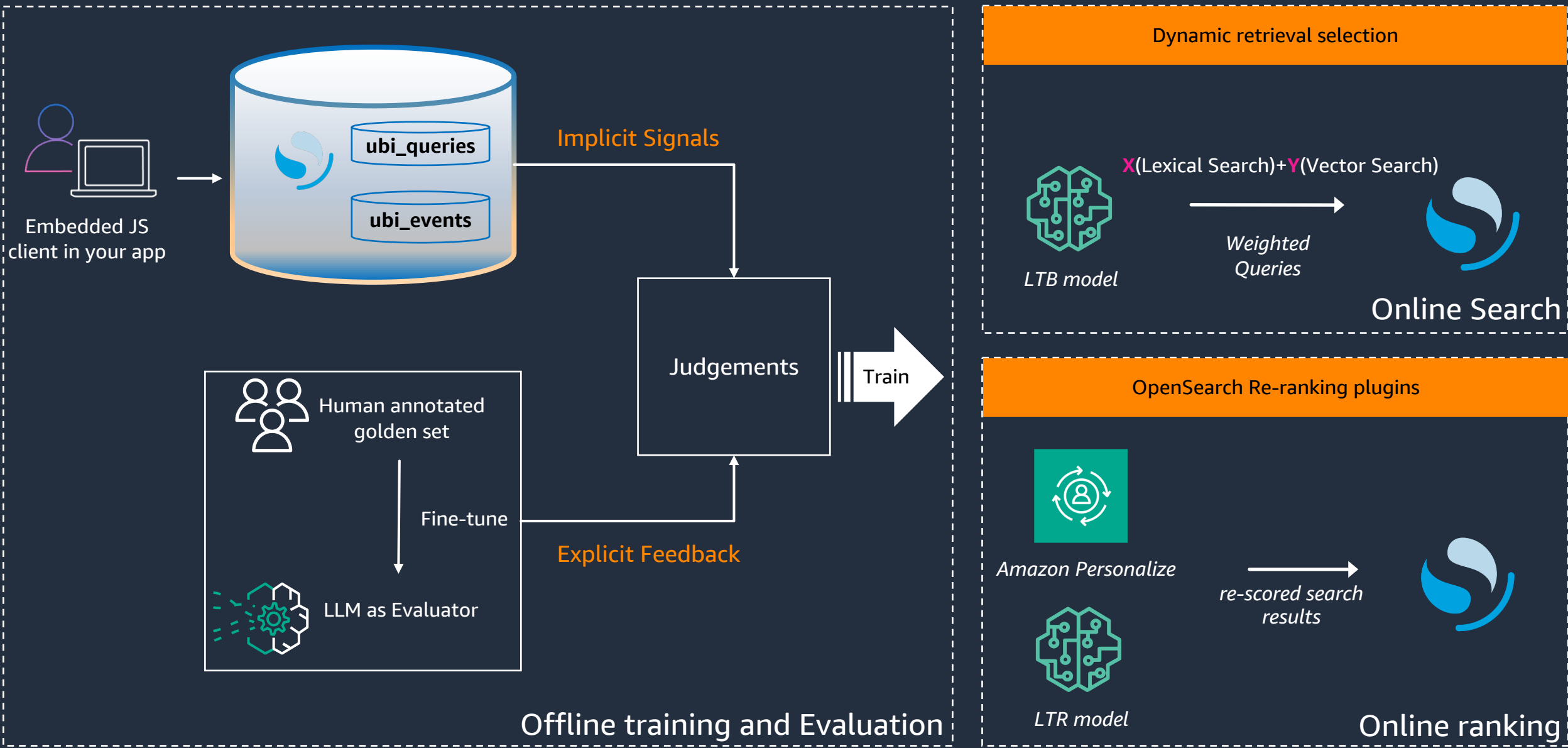
# Reranking with cross-encoder models

```
{
  "query": {
    "match": {
      "text_representation": ""…
    }
  },
  "ext": {
    "rerank": {
      "query_context": {
        "query_text": "Where is Albuquerque?"
      }
    }
  }
}
```

Retrieve documents → Get the new scores → Re-rank documents

Cross-encoder model

doc_#

doc_#

doc_#

doc_#

Search pipeline: Rerank processor

# User Behavior insights to improve your search and ranking

Embedded JS client in your app

ubi_queries

ubi_events

Implicit Signals

Human annotated golden set

Fine-tune

LLM as Evaluator

Explicit Feedback

Judgements

Train

Offline training and Evaluation

## Dynamic retrieval selection

X(Lexical Search)+Y(Vector Search)

LTB model

Weighted Queries

Online Search

## OpenSearch Re-ranking plugins

Amazon Personalize

LTR model

re-scored search results

Online ranking

# Demo

# Thank you!

Hajer Bouafif

linkedin.com/in/hajerbouafif

Praveen Mohan

linkedin.com/in/praveen-mohanprasad