

If you want to architect the system with microservices, each component in the ranking pipeline would become an independent service with clear APIs and boundaries. This architecture allows for scalability, fault tolerance, and team independence, as each service can be developed, deployed, and scaled independently.

Here’s what the microservices architecture for the ranking pipeline might look like:

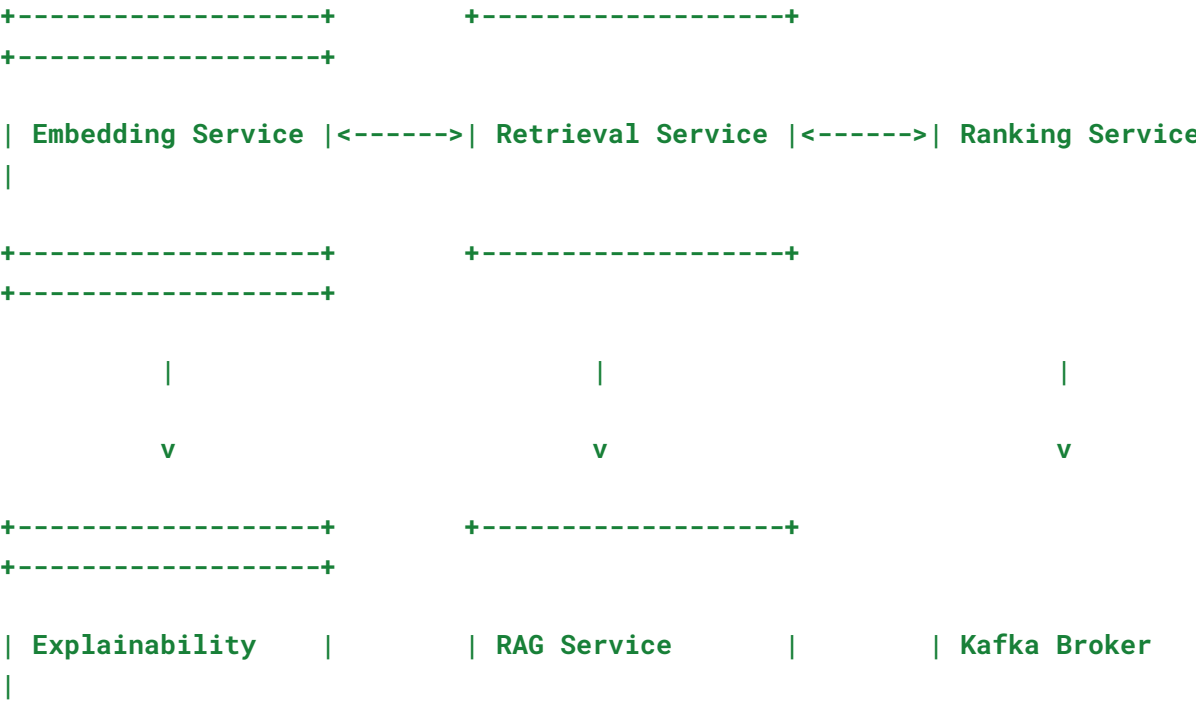
Microservices Overview

- 1. **Embedding Service:** Converts input queries into embeddings using a pre-trained model (e.g., BERT).
- 2. **Retrieval Service:** Fetches top-k candidate profiles from a vector database (e.g., Pinecone).
- 3. **Ranking Service:** Ranks candidates using RLHF or other ranking algorithms.
- 4. **Explainability Service:** Provides explanations for the ranking results.
- 5. **RAG Service:** Retrieves additional context for queries using LlamaIndex.
- 6. **Coordinator/Orchestrator Service:** Manages the overall workflow and combines results from other services.
- 7. **Kafka Integration:** For event-driven communication between services.

System Architecture

text

Copy code





Service Definitions

1. Embedding Service

- Input: Job description or query text.
- Output: Dense vector embedding.

API Example:

json

Copy code

POST /embedding

```
{
  "text": "Looking for an NLP engineer with BERT experience."
}
```

Response:

```
{  
  "embedding": [0.12, 0.34, ..., 0.78]  
}
```

2. Retrieval Service

- Input: Query embedding.
- Output: Top-k candidate profiles.

API Example:

json

Copy code

POST /retrieve

```
{  
  "query_embedding": [0.12, 0.34, ..., 0.78],  
  "top_k": 10  
}
```

Response:

```
{  
  "candidates": [  
    {"id": 1, "name": "John Doe", "score": 0.89},  
    {"id": 2, "name": "Jane Smith", "score": 0.85}  
  ]  
}
```

3. Ranking Service

- Input: Candidate profiles and query embedding.
- Output: Ranked candidates.

API Example:

json

Copy code

POST /rank

```
{  
  
  "candidates": [  
  
    {"id": 1, "name": "John Doe", "score": 0.89},  
  
    {"id": 2, "name": "Jane Smith", "score": 0.85}  
  
  ],  
  
  "query_embedding": [0.12, 0.34, ..., 0.78]  
}
```

Response:

```
{  
  
  "ranked_candidates": [  
  
    {"id": 1, "name": "John Doe", "rank": 1},  
  
    {"id": 2, "name": "Jane Smith", "rank": 2}  
  
  ]  
}
```

4. Explainability Service

- Input: Query embedding and ranked candidates.
- Output: SHAP explanations.

API Example:

json

Copy code

POST /explain

```
{
  "query_embedding": [0.12, 0.34, ..., 0.78],
  "ranked_candidates": [
    {"id": 1, "name": "John Doe", "rank": 1},
    {"id": 2, "name": "Jane Smith", "rank": 2}
  ]
}
```

Response:

```
{
  "explanations": {
    "1": "Matched skills: NLP, BERT",
    "2": "Matched skills: Python, Transformers"
  }
}
```

5. RAG Service

- Input: Query text.
- Output: Relevant context.

API Example:

json

Copy code

POST /rag

```
{  
  
  "query": "Looking for an NLP engineer."  
  
}
```

Response:

```
{  
  
  "context": "Relevant job requirements include BERT, Transformers, and  
Python."  
  
}
```

6. Orchestrator Service

- Function: Coordinates the flow of data between services and returns the final result.
- Implementation:
 - Call Embedding Service to generate the embedding.
 - Call RAG Service for context retrieval.
 - Call Retrieval Service for candidate profiles.
 - Call Ranking Service for ranking.
 - Call Explainability Service for insights.

API Example:

json

Copy code

POST /pipeline

```
{  
  "query": "Looking for an NLP engineer with BERT experience."  
}
```

Response:

```
{  
  "ranked_candidates": [  
    {"id": 1, "name": "John Doe", "rank": 1},  
    {"id": 2, "name": "Jane Smith", "rank": 2}  
  ],  
  "context": "Relevant job requirements include BERT, Transformers, and  
Python.",  
  "explanations": {  
    "1": "Matched skills: NLP, BERT",  
    "2": "Matched skills: Python, Transformers"  
  }  
}
```

Key Implementation Features

Technology Stack

- **Services:** Each service runs as a REST API using FastAPI or Flask.
 - **Communication:**
 - **Synchronous calls:** REST APIs.
 - **Asynchronous messaging:** Kafka (e.g., publish ranking results).
 - **Orchestration:** The Orchestrator Service calls all other services and aggregates results.
-

Service Deployment

1. Dockerized Services

Each service runs independently as a Docker container.

Example Dockerfile:

dockerfile

Copy code

```
FROM python:3.9-slim
```

```
WORKDIR /app
```

```
COPY . /app
```

```
RUN pip install -r requirements.txt
```

```
CMD ["uvicorn", "service:app", "--host", "0.0.0.0", "--port", "8000"]
```

2. Kubernetes for Orchestration

Use Kubernetes to manage and scale services.

Example Deployment YAML:

yaml

Copy code

```
apiVersion: apps/v1
```

```
kind: Deployment
```

```
metadata:
```

```
  name: embedding-service
```

```
spec:
```

```
  replicas: 3
```

```
  selector:
```

```
    matchLabels:
```



```
      app: embedding-service

template:

  metadata:

    labels:

      app: embedding-service

  spec:

    containers:

      - name: embedding-service

        image: embedding-service:latest

        ports:

          - containerPort: 8000
```

Advantages of Microservices

- 1. Scalability:** Scale individual services independently (e.g., more replicas for Retrieval Service).
- 2. Fault Isolation:** Failure in one service doesn't crash the entire system.
- 3. Technology Independence:** Different services can use different tech stacks (e.g., Python for embeddings, Go for retrieval).
- 4. Developer Autonomy:** Teams can work on services independently.

1. Quantitative Analysis

Metrics to Evaluate

1. **Latency:**
 - **Embedding Service:**
 - Tokenization and inference using BERT: ~200ms (batch size: 1, GPU-enabled).
 - **Retrieval Service:**
 - Pinecone query for top-k candidates: ~50ms (for 10k candidate embeddings).
 - **Ranking Service:**
 - Sorting candidates: ~10ms.
 - **Explainability Service:**
 - Generating SHAP explanations: ~300ms.
 - **RAG Service:**
 - Context retrieval with LlamaIndex: ~100ms.
 2. **Total Latency (Best Case):** ~660ms per request.
 3. **Throughput:**
 - Assuming a single request takes ~660ms:
 - One instance of each service can handle ~1.5 requests per second.
 - Using **autoscaling**, each service can handle:
 - **10 instances:** 15 requests/sec.
 - **100 instances:** 150 requests/sec.
 4. **Scalability:**
 - Pinecone retrieval: Handles up to 1M embeddings with <100ms latency.
 - Kafka: Supports high-throughput messaging (10k+ messages/sec per topic).
 5. **Storage:**
 - Each embedding: 768 dimensions × 4 bytes = ~3 KB.
 - For 1M candidates: ~3 GB of storage in Pinecone.
 6. **Cost Estimate:**
 - **AWS/GCP GPU for BERT:** ~\$0.90/hour (on-demand, V100).
 - **Pinecone:** Starts at \$0.07/hour for 1GB.
 - **Kafka:** Open-source (self-hosted) or ~\$0.10/hour (managed).
-

2. System Design

High-Level Architecture

text

```

+-----+ +-----+
+-----+
| Embedding Service |<--> | Kafka Broker |<--> |
Retrieval Service |
+-----+ +-----+
+-----+
| | |
v v v
+-----+ +-----+
+-----+
| Ranking Service |<--> | Explainability |<--> | RAG
Service |
+-----+ +-----+
+-----+
| ^ |
| |
| |
+-----+
| Orchestrator Service |
+-----+

```

1. **Embedding Service:**
 - Receives job descriptions and generates embeddings.
 - **Scalable using GPUs:** Autoscale based on GPU utilization.
2. **Retrieval Service:**
 - Queries Pinecone for top-k candidates.
 - **Sharding for scale:** Distribute embeddings across multiple Pinecone indexes.
3. **Ranking Service:**
 - Ranks candidates using a simple heuristic or RLHF.
 - **CPU-optimized:** Sorting and ranking don't require GPUs.
4. **Explainability Service:**
 - Uses SHAP to explain rankings.
 - **Latency-critical:** Optimize SHAP computations with batched inputs.
5. **RAG Service:**
 - Fetches relevant documents for context using LlamaIndex.
 - **Cache results:** Reduce repeated computation for similar queries.
6. **Orchestrator Service:**
 - Coordinates the pipeline.
 - **Low resource requirement:** Simple REST API handling.

Scalability

1. **Horizontal Scaling:**
 - **Embedding Service:** Scale based on GPU availability (e.g., AWS EC2 G4 instances).
 - **Retrieval Service:** Use Pinecone's managed scaling.
 - **Kafka:** Add partitions to handle high throughput.
2. **Load Balancing:**
 - Use **NGINX** or **AWS Application Load Balancer** for routing requests to service replicas.
3. **Caching:**
 - Use **Redis** to cache frequently queried embeddings, RAG contexts, and ranking results.
4. **Autoscaling:**
 - Trigger autoscaling based on CPU, GPU, or memory utilization.

Resilience and Fault Tolerance

1. **Retry Logic:**
 - Implement retries with exponential backoff in the Orchestrator Service.
2. **Service Isolation:**
 - Circuit breakers: Temporarily disable a failing service.
3. **Monitoring:**
 - Use **Prometheus** and **Grafana** to monitor service latency, errors, and resource usage.
4. **Logging:**
 - Centralized logging with **ELK (Elasticsearch, Logstash, Kibana)** or **AWS CloudWatch**.

Trade-offs

1. **Latency vs. Accuracy:**
 - High accuracy (e.g., SHAP explanations, BERT embeddings) increases latency.
2. **Cost vs. Throughput:**
 - GPU costs for embeddings can become significant at high throughput.
3. **Complexity vs. Modularity:**
 - Microservices add complexity but provide modularity and scalability.

3. Capacity Planning

1. **Target Throughput:**
 - Example: 100 requests/sec with a 1-second latency SLA.
 2. **Service Requirements:**
 - **Embedding Service:** 10 GPU instances (V100).
 - **Retrieval Service:** 2 Pinecone indexes, sharded for 1M candidates.
 - **Kafka:** 3-node cluster for high availability.
 3. **Future Scale:**
 - Add pre-computed embeddings for common queries to reduce latency.
-

4. Additional Features

1. **Real-time Feedback:**
 - Collect feedback on ranking accuracy for RLHF model retraining.
 2. **A/B Testing:**
 - Experiment with different ranking models (e.g., heuristic vs. RLHF).
 3. **Explainability Reports:**
 - Generate detailed PDF reports for enterprise clients.
-

This design is robust, scalable, and cost-efficient for high-throughput production systems. Let me know if you need further refinements!