

# BSH AIM RAG - Scheduler

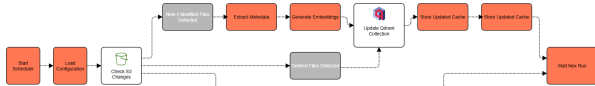
The **BSH AIM RAG Scheduler**, developed by the **AIX team** (formerly CDL), is a scalable and efficient solution for managing the synchronization of files stored in Amazon S3 with **Qdrant**, a high-performance vector database. It supports both **initial indexing** and **incremental updates** by detecting changes in stored files. That means only newly **added**, **modified**, or **deleted** documents get processed — saving time and resources.

## Key Features

				
<b>Change Detection</b>	<b>Metadata Extraction</b>	<b>Error Handling</b>	<b>Configurability</b>	<b>User-Friendly Interface</b>
Uses an S3 housed cache to track file hashes, only processing new or modified files	Automatically extracts metadata from file paths	Includes comprehensive logging and error handling	All important parameters can be set via command-line arguments	A cool UI allowing embedding and indexing of data housed in S3 in clicks.

## Modular Architecture

The **BSH AIM RAG Scheduler** follows a **modular architecture**, where each module is responsible for a specific function:

Module	Responsibility	How It Works	Data Flow	Benefits
Config	Loads configurations from environment variables	<ul style="list-style-type: none"><li>Loading environment variables</li><li>Managing AWS, S3, and Qdrant configurations</li><li>Handling collection-specific settings through CollectionConfig dataclass</li></ul>	<div><ul style="list-style-type: none"><li>main.py reads config, initializes sync manager</li><li>sync_manager coordinates operations:<ul style="list-style-type: none"><li>S3_manager reads files</li><li>metadata_extractor processes paths</li><li>embeddings creates vectors</li><li>vector_store stores in Qdrant</li><li>cache_manager tracks changes</li></ul></li><li>Process repeats for incremental updates</li></ul></div> <div>Flow Chart</div>  <td><p>This modular architecture allows for:</p><ul style="list-style-type: none"><li>Easy addition of new embedding providers</li><li>Simplified collection configuration changes</li><li>Efficient incremental updates</li><li>Reliable deletion handling</li><li>Scalable document processing</li></ul></td>	<p>This modular architecture allows for:</p> <ul style="list-style-type: none"><li>Easy addition of new embedding providers</li><li>Simplified collection configuration changes</li><li>Efficient incremental updates</li><li>Reliable deletion handling</li><li>Scalable document processing</li></ul>
S3Manager	Manages all interactions with Amazon S3	<ul style="list-style-type: none"><li>Listing files in collection-specific folders</li><li>Reading JSON documents from S3</li><li>Streaming documents for processing</li><li>Using folder prefixes to scope operations to specific collections</li></ul>		
MetadataExtractor	Extracts key metadata from file paths	<ul style="list-style-type: none"><li>Parsing S3 keys to extract country information</li><li>Mapping folder structures to metadata attributes</li><li>Providing consistent metadata format for document indexing</li></ul>		
Embedding Provider	Interfaces with multiple embedding models	<ul style="list-style-type: none"><li>Providing unified interface for different embedding providers</li><li>Supporting HuggingFace, Azure OpenAI, and DashScope models</li><li>Converting text to vector embeddings</li><li>Handling model-specific configurations</li></ul>		
VectorStore	Manages interactions with Qdrant, including indexing	<ul style="list-style-type: none"><li>Creating and managing collections</li><li>Converting documents to vector points</li><li>Handling upserts and deletions</li><li>Managing document IDs and metadata</li><li>Providing collection information</li></ul>		
CacheManager	Maintains local file cache for efficient change detection	<ul style="list-style-type: none"><li>Storing document hashes in S3</li><li>Detecting new, modified, and deleted files</li><li>Managing cache files per collection</li><li>Providing efficient change detection</li></ul>		

<b>SyncManager</b>	Orchestrates the entire sync process	<ul style="list-style-type: none"><li>• Coordinating between S3 and Qdrant</li><li>• Managing initial and incremental indexing</li><li>• Handling file deletions</li><li>• Processing documents in batches</li><li>• Maintaining sync state</li></ul>	
<b>Main</b>	Entry point and CLI interface	<ul style="list-style-type: none"><li>• Parsing command-line arguments</li><li>• Loading collection configurations</li><li>• Initializing components</li><li>• Managing sync operations</li><li>• Providing flexible execution options (initial/incremental /specific collection)</li></ul>	

## Collection Configuration

Each collection (e.g., "hr-data-china") is configured in `collections_config.json`, it:

- Creates or connects to the Qdrant collection with the same name
- Maps to the matching S3 folder
- Performs initial indexing if needed
- Identifies new, modified, and deleted files
- Updates the Qdrant collection accordingly

```
{
  "collections": [
    {
      "name": "hr-data-china",
      "embedding_model": "text-embedding-v3",
      "embedding_type": "alibaba"
    },
    {
      "name": "hr-data-generic",
      "embedding_model": "text-embedding-ada-002",
      "embedding_type": "openai"
    }
  ]
}
```

## Supported Embedding Models

The scheduler supports multiple embedding models from different providers:

Provider		Embedding Models	Dimension	Supported Languages	Supported	Pricing (1M Tokens)	For More
block ed URL	Hugging Face	all-MiniLM-L6-v2	384	<ul style="list-style-type: none"><li>• en</li></ul>	✔	free	<ul style="list-style-type: none"><li>• <a href="#">SentenceTransformers Documentation</a></li><li>• <a href="#">Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks</a></li></ul>
		BAAI/bge-m3	1024	<ul style="list-style-type: none"><li>• en</li><li>• zh</li></ul>	✔	free	<ul style="list-style-type: none"><li>• <a href="#">BGE M3-Embedding: Multi-Lingual, Multi-Functionality, Multi-Granularity Text Embeddings Through Self-Knowledge Distillation</a></li></ul>
block ed URL	Azure Open AI	text-embedding-ada-002	1536	<ul style="list-style-type: none"><li>• en</li><li>• zh</li><li>• ..</li></ul>	✔	\$0.10	<ul style="list-style-type: none"><li>• <a href="#">Vector embedding   OpenAI</a></li></ul>
		text-embedding-3-large			✘	\$0.13	
		text-embedding-3-small			✘	\$0.02	
block ed URL	Alibaba	text-embedding-v3	1024	<ul style="list-style-type: none"><li>• en</li><li>• zh</li></ul>	✔	Free for a limited quota. If the throttling limit is exceeded, your API request fails due to throttling. You must wait for a period of time until the throttling conditions are met before you can call the API again.	

	text-embedding-v2			✗	
	text-embedding-v1			✗	



Each provider implements an **EmbeddingProvider** interface to ensure consistency across different models. Additional models can be added by implementing the following class:

```
import numpy as np
from abc import ABC, abstractmethod

class EmbeddingProvider(ABC):
    @abstractmethod
    def encode(self, text: str) -> np.ndarray:
        pass

    @abstractmethod
    def get_dimension(self) -> int:
        pass
```

## Getting Started

### Sync Commands

This scheduler or indexer tool synchronizes data between Amazon S3 and Qdrant vector database. Below is a comprehensive list of all available commands and their functions.

Command	Format	Description	Example
config	--config PATH	Specifies the path to the collection configuration JSON file that defines S3 buckets and Qdrant collections to synchronize.	python main.py --config configs/collections.json
initial-only	--initial-only	Performs only the initial indexing of data from S3 to Qdrant and then exits without monitoring for changes. Useful for first-time setup.	python main.py --config configs/collections.json --initial-only
incremental-only	--incremental-only	Skips the initial full indexing and only processes files that have changed since the last synchronization. Useful for update operations.	python main.py --config configs/collections.json --incremental-only
collection	--collection NAME	Removes records from Qdrant when corresponding files are deleted from S3. Can be combined with other commands.	python main.py --config configs/collections.json --collection documents

<pre>-- remove-deleted</pre>	<pre>-- remove-deleted</pre>	Removes records from Qdrant when corresponding files are deleted from S3. Can be combined with other commands.	<pre>python main.py --config configs /collections.json --incremental-only -- remove-deleted</pre>
------------------------------	------------------------------	--	---

## Command Combination

You can combine multiple commands for more specific synchronization operations:

Combination	Description	Example
<pre>--config -- collection --initial-only</pre>	Perform initial indexing for a specific collection only	<pre>python main.py --config configs /collections.json --collection aix_cdl_data -- initial-only</pre>
<pre>--config -- incremental-only --remove- deleted</pre>	Process only changes and remove deleted files across all collections	<pre>python main.py --config configs /collections.json --incremental-only --remove- deleted</pre>
<pre>--config -- collection --incremental- only --remove-deleted</pre>	Process only changes and remove deleted files for a specific collection	<pre>python main.py --config configs /collections.json --collection customer_data -- incremental-only --remove-deleted</pre>
<pre>--config configs /collections.json</pre>	<ul style="list-style-type: none"><li>• Perform full synchronization for all collections</li><li>• This initiates the sync process, handles incremental updates, and processes deletions</li></ul>	<pre>python main.py --config configs /collections.json</pre>